

## 分析并预处理user\_profile数据集

### 缺失值处理方案：

1. 填充方案：结合用户的其他特征值，利用随机森林算法进行预测；但产生了大量人为构建的数据，一定程度上增加了数据的噪音
2. 把变量映射到高维空间：如pvalue\_level的1维数据，转换成是否1、是否2、是否3、是否缺失的4维数据；这样保证了所有原始数据不变，同时能提高精确度，但这样会导致数据变得比较稀疏，如果样本量很小，反而会导致样本效果较差，因此也不能滥用

### 最终方案2的效果好

In [ ]:

```
1 '''
2 特征选取
3 我该选取那些特征呢？
4 除了前面处理的pvalue_level和new_user_class_level需要作为特征以外，（能体现出用户的购买
5 分类特征值：个数
6 - cms_segid: 97
7 - cms_group_id: 13
8 - final_gender_code: 2
9 - age_level: 7
10 - shopping_level: 3
11 - occupation: 2
12 - pvalue_level
13 - new_user_class_level
14 - price
15 根据经验，以上几个分类特征都一定程度能体现用户在购物方面的特征，且类别都较少，都可以用
16 '''
```

```
In [1]: 1 import os
2 # 配置pyspark和spark driver运行时 使用的python解释器
3 JAVA_HOME = '/root/bigdata/jdk'
4 PYSARK_PYTHON = '/miniconda2/envs/py365/bin/python'
5 # 当存在多个版本时, 不指定很可能会导致出错
6 os.environ['PYSARK_PYTHON'] = PYSARK_PYTHON
7 os.environ['PYSARK_DRIVER_PYTHON'] = PYSARK_PYTHON
8 os.environ['JAVA_HOME'] = JAVA_HOME
9 # 配置spark信息
10 from pyspark import SparkConf
11 from pyspark.sql import SparkSession
12
13 SPARK_APP_NAME = 'preprocessingUserProfile'
14 SPARK_URL = 'spark://192.168.58.100:7077'
15
16 conf = SparkConf()
17 config = (
18     ('spark.app.name', SPARK_APP_NAME),
19     ('spark.executor.memory', '2g'),
20     ('spark.master', SPARK_URL),
21     ('spark.executor.cores', '2')
22 )
23 conf.setAll(config)
24
25 spark = SparkSession.builder.config(conf=conf).getOrCreate()
```

```
In [2]: 1 df = spark.read.csv('/data/user_profile.csv',header=True)
2 # 发现pvalue_level和new_user_class_level存在空值:
3 # (注意此处的null表示空值,而如果是NULL,则往往表示是一个字符串)
4 # 因此直接利用schema就可以加载进该数据,无需替换null值
5 df.show()
```

userid	cms_segid	cms_group_id	final_gender_code	age_level	pvalue_level	shopping_level	occupation	new_user_class_level
234	0	5	2	5	null	3		
523	5	2	2	2	1	3		
612	0	8	1	2	2	3		
1670	0	4	2	4	null	1		
2545	0	10	1	4	null	3		
3644	49	6	2	6	2	3		
5777	44	5	2	5	2	3		
6211	0	9	1	3	null	3		
6355	2	1	2	1	1	3		
6823	43	5	2	5	2	3		
6972	5	2	2	2	2	3		
9293	0	5	2	5	null	3		
9510	55	8	1	2	2	2		
10122	33	4	2	4	2	3		
10549	0	4	2	4	2	3		
10812	0	4	2	4	null	2		
10912	0	4	2	4	2	3		
10996	0	5	2	5	null	3		
11256	8	2	2	2	1	3		
11310	31	4	2	4	1	3		

only showing top 20 rows

In [3]:

```
1 df.printSchema()

root
|-- userid: string (nullable = true)
|-- cms_segid: string (nullable = true)
|-- cms_group_id: string (nullable = true)
|-- final_gender_code: string (nullable = true)
|-- age_level: string (nullable = true)
|-- pvalue_level: string (nullable = true)
|-- shopping_level: string (nullable = true)
|-- occupation: string (nullable = true)
|-- new_user_class_level : string (nullable = true)
```

In [7]:

```
1 ## 延申学习dropna()
2 ## str的空值不能通过dropna() 去除掉
3 # print(df.select('pvalue_level').count())
4 # df.dropna()
5 # print(df.select('pvalue_level').count())
```

1061768

1061768

In [4]:

```
1 # 注意：这里的null会直接被pyspark识别为None数据，也就是na数据，所以这里可以直接利用sch
2 # 注意：如果数据集中存在NULL字样的数据，无法直接设置schema，只能先将NULL类型的数据处理
3 # 如果直接schema，下图会变成下图：
4 # +-----+-----+-----+-----+-----+-----+
5 # |adgroup_id|cate_id|campaign_id|customer| brand|price|
6 # +-----+-----+-----+-----+-----+-----+
7 # |   375706|   4520|   387991|      6|  NULL| 99.0|
8 #
9 # +-----+-----+-----+-----+-----+-----+
10 # |adgroupId|cateId|campaignId|customerId|brandId|price|
11 # +-----+-----+-----+-----+-----+-----+
12 # |      null|  null|      null|      null|  null| null|
13 #处理方式如下：
14 # from pyspark.sql.types import IntegerType, FloatType
15 # ad_feature_df = df.\
16 #     withColumn("adgroup_id", df.adgroup_id.cast(IntegerType())).withColumnRenamed("a
17 #     withColumn("cate_id", df.cate_id.cast(IntegerType())).withColumnRenamed("cate_id
18 #     withColumn("campaign_id", df.campaign_id.cast(IntegerType())).withColumnRenamed
19 #     withColumn("customer", df.customer.cast(IntegerType())).withColumnRenamed("custo
20 #     withColumn("brand", df.brand.cast(IntegerType())).withColumnRenamed("brand", "br
21 #     withColumn("price", df.price.cast(FloatType()))
```

```
In [3]: 1 # 当然了，我们只能看到部分数值，所以有可能还有NULL，只是没有被我们看到，所以最好使用上
2 from pyspark.sql.types import StructType, StructField, IntegerType
3 schema = StructType([
4     StructField("userId", IntegerType()),
5     StructField("cms_segid", IntegerType()),
6     StructField("cms_group_id", IntegerType()),
7     StructField("final_gender_code", IntegerType()),
8     StructField("age_level", IntegerType()),
9     StructField("pvalue_level", IntegerType()),
10    StructField("shopping_level", IntegerType()),
11    StructField("occupation", IntegerType()),
12    StructField("new_user_class_level", IntegerType())
13 ])
14 user_profile_df = spark.read.csv('/data/user_profile.csv', header=True, schema=schema)
15 user_profile_df.printSchema()
16 user_profile_df.show()
```

```
root
```

```
-- userId: integer (nullable = true)
-- cms_segid: integer (nullable = true)
-- cms_group_id: integer (nullable = true)
-- final_gender_code: integer (nullable = true)
-- age_level: integer (nullable = true)
-- pvalue_level: integer (nullable = true)
-- shopping_level: integer (nullable = true)
-- occupation: integer (nullable = true)
-- new_user_class_level: integer (nullable = true)
```

```
+-----+-----+-----+-----+-----+-----+-----+
+-----+
|userId|cms_segid|cms_group_id|final_gender_code|age_level|pvalue_level|shopping_level|
+-----+
|      234|      0|      5|      2|      5|      null|      3|
+-----+
|      500|      0|      5|      2|      5|      null|      3|
+-----+
```

```
In [6]: 1 print('每个特征包含的种类数: ')
2 print(['特征' + str(c) + '的个数' + str(user_profile_df.groupby(c).count().count()) for c in user_profile_df.columns])
3 print('表中总共有%d行'%user_profile_df.count())
4 print('每列各自的行数: ')
5 [str(c) + '列的行数: ' + str(user_profile_df.select(c).dropna().count()) for c in user_profile_df.columns]
```

每个特征包含的种类数:

['特征cms\_segid的个数97', '特征cms\_group\_id的个数13', '特征final\_gender\_code的个数2', '特征age\_level的个数7', '特征pvalue\_level的个数4', '特征shopping\_level的个数3', '特征occupation的个数2', '特征new\_user\_class\_level的个数5']

表中总共有1061768行

每列各自的行数:

```
Out[6]: ['userId列的行数: 1061768',
'cms_segid列的行数: 1061768',
'cms_group_id列的行数: 1061768',
'final_gender_code列的行数: 1061768',
'age_level列的行数: 1061768',
'pvalue_level列的行数: 485851',
'shopping_level列的行数: 1061768',
'occupation列的行数: 1061768',
'new_user_class_level列的行数: 716848']
```

```
In [7]: 1 user_profile_df.groupby('final_gender_code').count().show()
```

```
+-----+-----+
|final_gender_code| count|
+-----+-----+
|                | 1|377517|
|                | 2|684251|
+-----+-----+
```

## 1. 处理缺失值，两列数据有空值

```
In [8]: 1 # 有缺失值的特征
        2 user_profile_df.groupby('pvalue_level').count().show()
        3 user_profile_df.groupby('new_user_class_level').count().show()
```

```
+-----+-----+
|pvalue_level| count|
+-----+-----+
|          null|575917|
|              1|154436|
|              3| 37759|
|              2|293656|
+-----+-----+
```

```
+-----+-----+
|new_user_class_level| count|
+-----+-----+
|                  null|344920|
|                      1| 80548|
|                      3|173047|
|                      4|138833|
|                      2|324420|
+-----+-----+
```

## 1.1 先不删除确实严重的特征

### 缺失值数据处理

注意，一般情况下：

- 缺失率低于10%：可直接进行相应的填充，如默认值、均值、算法拟合等等；
- 高于10%：往往会考虑舍弃该特征
- 特征处理，如1维转多维

但根据我们的经验，我们的广告推荐其实和用户的消费水平、用户所在城市等级都有比较大的关联，因此在这里pvalue\_level、new\_user\_class\_level都是比较重要的特征，我们不考虑舍弃

### 缺失值处理方案：

1. 填充方案：结合用户的其他特征值，利用随机森林算法进行预测；但产生了大量人为构建的数据，一定程度上增加了数据的噪音
2. 把变量映射到高维空间：如pvalue\_level的1维数据，转换成是否1、是否2、是否3、是否缺失的4维数据；这样保证了所有原始数据不变，同时能提高精确度，但这样会导致数据变得比较稀疏，如果样本量很小，反而会导致样本效果较差，因此也不能滥用

#### 1.1.1 填充方案之 利用随机森林对缺失值进行预测

##### a. 利用随机森林对pvalue\_level的缺失值进行预测

```
In [10]: 1 # 随机森林填充缺失值
2 # 随机森林模型的输入值类型必须是 LabeledPoint类型
3 from pyspark.mllib.regression import LabeledPoint
4 # dropna() 整行删除, 某行有一个空值 就删掉整行
5 # mlib是基于rdd的, 因此后续处理过程需要转化成rdd
6 train_data = user_profile_df.dropna(subset=['pvalue_level']).rdd.\
7     map(lambda r:LabeledPoint(r.pvalue_level-1,\
8         # final_gender_code原数据中是1、2, 随机森林模型使用时要变成0, 1
9         [r.cms_segid, r.cms_group_id, r.final_gender_code-1, r.age_level, r.shopping_level]))
10 train_data.collect()
```

```
Out[10]: [LabeledPoint(0.0, [5.0, 2.0, 1.0, 2.0, 2.0, 1.0]),
LabeledPoint(1.0, [0.0, 8.0, 0.0, 2.0, 2.0, 0.0]),
LabeledPoint(1.0, [49.0, 6.0, 1.0, 6.0, 2.0, 0.0]),
LabeledPoint(1.0, [44.0, 5.0, 1.0, 5.0, 2.0, 0.0]),
LabeledPoint(0.0, [2.0, 1.0, 1.0, 1.0, 2.0, 0.0]),
LabeledPoint(1.0, [43.0, 5.0, 1.0, 5.0, 2.0, 0.0]),
LabeledPoint(1.0, [5.0, 2.0, 1.0, 2.0, 2.0, 1.0]),
LabeledPoint(1.0, [55.0, 8.0, 0.0, 2.0, 1.0, 0.0]),
LabeledPoint(1.0, [33.0, 4.0, 1.0, 4.0, 2.0, 0.0]),
LabeledPoint(1.0, [0.0, 4.0, 1.0, 4.0, 2.0, 0.0]),
LabeledPoint(1.0, [0.0, 4.0, 1.0, 4.0, 2.0, 0.0]),
LabeledPoint(0.0, [8.0, 2.0, 1.0, 2.0, 2.0, 0.0]),
LabeledPoint(0.0, [31.0, 4.0, 1.0, 4.0, 2.0, 0.0]),
LabeledPoint(1.0, [20.0, 3.0, 1.0, 3.0, 2.0, 0.0]),
LabeledPoint(1.0, [33.0, 4.0, 1.0, 4.0, 2.0, 0.0]),
LabeledPoint(1.0, [36.0, 5.0, 1.0, 5.0, 0.0, 0.0]),
LabeledPoint(1.0, [20.0, 3.0, 1.0, 3.0, 2.0, 0.0]),
LabeledPoint(1.0, [8.0, 2.0, 1.0, 2.0, 2.0, 0.0]),
LabeledPoint(1.0, [0.0, 4.0, 1.0, 4.0, 2.0, 0.0]),
LabeledPoint(0.0, [5.0, 2.0, 1.0, 2.0, 2.0, 1.0])]
```

```
In [11]: 1 # # 延申学习: LabeledPoint 与 稀疏矩阵SparseVector
2 # from pyspark.mllib.linalg import SparseVector
3 # from pyspark.mllib.regression import LabeledPoint
4
5 # # Create a labeled point with a positive label and a dense feature vector.
6 # pos = LabeledPoint(1.0, [1.0, 0.0, 3.0])
7
8 # # Create a labeled point with a negative label and a sparse feature vector.
9 # neg = LabeledPoint(0.0, SparseVector(3, [0, 2], [1.0, 3.0]))
```

```
In [12]: 1 from pyspark.mllib.tree import RandomForest
2 # 训练分类模型
3 # 参数1 训练的数据
4 #参数2 目标值的分类个数 0,1,2
5 #参数3 特征中是否包含分类的特征 {2:2,3:7} {2:2} 表示 在特征中 第三个特征是分类的: 有两
6 #参数4 随机森林中 树的棵数
7 model = RandomForest.trainClassifier(train_data, 3, {2:2, 3:7, 4:3, 5:2}, 5)
```



```
In [13]: 1 # 看看上上个cell中哪个特征是分类特征
2 # user_profile_df.groupby('cms_segid').count().show()# 种类太多
3 # user_profile_df.groupby('cms_group_id').count().show()# 种类太多
4 # user_profile_df.groupby('final_gender_code').count().show()# 1、2 ->0、1
5 # user_profile_df.groupby('age_level').count().show() # 0-6
6 # user_profile_df.groupby('shopping_level').count().show() # 1、2、3 -> 0、1、2
7 # user_profile_df.groupby('occupation').count().show() # 0、1
```

```
In [14]: 1 # 取出pvalue_level缺失的那些行
2 pl_na_df = user_profile_df.fillna(-1).where('pvalue_level!=-1')
```

```
In [15]: 1 # 随机森林模型 使用的是rdd, 准备输入模型的rdd
2 rdd = pl_na_df.rdd.map(lambda r: (r.cms_segid, r.cms_group_id, r.final_gender_code-1,
3 predicts = model.predict(rdd)
4 print(predicts.take(20))
5 print('预测值总数',predicts.count())
6 # 这里注意predict参数, 如果是预测多个, 那么参数必须是直接有列表构成的rdd参数, 而不能是
7 # 因此这里经过map函数处理, 将每一行数据转换为普通的列表数据
```

[1.0, 1.0, 1.0, 1.0, 1.0, 1.0, 1.0, 1.0, 0.0, 1.0, 1.0, 1.0, 1.0, 1.0, 0.0, 1.0, 0.0, 1.0, 1.0, 1.0]

预测值总数 575917

```
In [16]: 1 # 这里数据量比较小, 直接转换为pandas dataframe来处理, 因为方便, 但注意如果数据量较大不
2 temp = predicts.map(lambda x:int(x)).collect()
3 pdf = pl_na_df.toPandas()
4 import numpy as np
5 # 在pandas df的基础上直接替换掉列数据
6 pdf["pvalue_level"] = np.array(temp) + 1 # 注意+1 还原预测值
7 pdf
8
```

```
Out[16]:
```

	userId	cms_segid	cms_group_id	final_gender_code	age_level	pvalue_level	shopping_l
0	234	0	5	2	5	2	
1	1670	0	4	2	4	2	
2	2545	0	10	1	4	2	
3	6211	0	9	1	3	2	
4	9293	0	5	2	5	2	
...	...	...	...	...	...	...	...
575912	1137329	0	4	2	4	2	
575913	1137582	0	9	1	3	2	
575914	1137955	0	3	2	3	1	
575915	1138545	0	4	2	4	2	
575916	1139632	0	7	1	1	2	

575917 rows × 9 columns

```
In [17]: 1 # 与非缺失数据进行拼接, 完成pvalue_level的缺失值预测
2 new_user_profile_df = user_profile_df.dropna(subset=["pvalue_level"]).unionAll(spark.
3 new_user_profile_df.show()
4 # 注意: unionAll的使用, 两个df的表结构必须完全一样
```

...

### b. 利用随机森林对new\_user\_class\_level的缺失值进行预测

```
In [18]: 1 from pyspark.mllib.regression import LabeledPoint
2
3 # 选出new_user_class_level全部的
4 train_data2 = user_profile_df.dropna(subset=["new_user_class_level"]).rdd.map(
5     lambda r:LabeledPoint(r.new_user_class_level - 1,\
6                             [r.cms_segid, r.cms_group_id, r.final_gender_code-1, r.age_
7 ])
8 from pyspark.mllib.tree import RandomForest
9 model2 = RandomForest.trainClassifier(train_data2, 4, {2:2, 3:7, 4:3, 5:2}, 5)
```

```
In [19]: 1 nul_na_df = user_profile_df.na.fill(-1).where("new_user_class_level=-1")
2 nul_na_df.show(10)
3
4 def row(r):
5     return r.cms_segid, r.cms_group_id, r.final_gender_code, r.age_level, r.shopping_
6
7 rdd2 = nul_na_df.rdd.map(row)
8 predicts2 = model2.predict(rdd2)
9 predicts2.take(20)
```

```
+-----+-----+-----+-----+-----+-----+-----+
+-----+-----+-----+-----+
|userId|cms_segid|cms_group_id|final_gender_code|age_level|pvalue_level|shopping_level|
|occupation|new_user_class_level|
+-----+-----+-----+-----+-----+-----+-----+
+-----+-----+-----+
| 612| 0| 8| 1| 2| 2|
3| 0| -1|
| 1670| 0| 4| 2| 4| -1|
1| 0| -1|
| 2545| 0| 10| 1| 4| -1|
3| 0| -1|
| 10549| 0| 4| 2| 4| 2|
3| 0| -1|
| 10812| 0| 4| 2| 4| -1|
2| 0| -1|
| 10912| 0| 4| 2| 4| 2|
3| 0| -1|
| 12620| 0| 4| 2| 4| -1|
0| 0| 1|
```

**总结:** 可以发现由于这两个字段的缺失过多, 所以预测出来的值已经大大失真, 但如果缺失率在10%以下, 这种方法是比较有效的一种

因此我们接下来采用将变量映射到高维空间的方法来处理数据，即将缺失项也当做一个单独的特征来对待，保证数据的原始性

由于该思想正好和热独编码实现方法一样，因此这里直接使用热独编码方式处理数据

```
In [23]: 1 from pyspark.ml.feature import OneHotEncoder
2 from pyspark.ml.feature import StringIndexer
3 from pyspark.ml import Pipeline
4 from pyspark.sql.types import StringType
5 user_profile_df = user_profile_df.fillna(-1)
6 user_profile_df.show()
7 # 热独编码时，必须先将待处理字段 转为 字符串类型才可处理
8 user_profile_df = user_profile_df.withColumn('pvalue_level',user_profile_df.pvalue_level
9     .withColumn('new_user_class_level',user_profile_df.new_user_class_level.cast(StringType))
10 user_profile_df.printSchema()
11 # 对pvalue_level进行热独编码，求值
12 stringindexer = StringIndexer(inputCol='pvalue_level',outputCol='pl_onehot_feature')
13 encoder = OneHotEncoder(dropLast=False,inputCol='pl_onehot_feature',outputCol='pl_onehot_value')
14 pipeline = Pipeline(stages=[stringindexer,encoder])
15 pipeline_fit = pipeline.fit(user_profile_df)
16 user_profile_df2 = pipeline_fit.transform(user_profile_df)
17 # pl_onehot_value列的值为稀疏向量，存储热独编码的结果
18 user_profile_df2.printSchema()
19 user_profile_df2.show()
```

```
+-----+-----+-----+-----+-----+-----+-----+
+-----+-----+
|userId|cms_segid|cms_group_id|final_gender_code|age_level|pvalue_level|shopping_level|occupation|new_user_class_level|
+-----+-----+-----+-----+-----+-----+-----+
| 234|0|5|2|5|-1|
3|0|3|
| 523|5|2|2|2|1|
3|1|2|
| 612|0|8|1|2|2|
3|0|-1|
| 1670|0|4|2|4|-1|
1|0|-1|
| 2545|0|10|1|4|-1|
3|0|-1|
| 3644|49|6|2|6|2|
3|0|2|
| 5777|44|5|2|5|2|
3|0|2|
| 6211|0|9|1|3|-1|
3|0|2|
| 6355|2|1|2|1|1|
3|0|4|
| 6823|43|5|2|5|2|
3|0|1|
| 6972|5|2|2|2|2|
3|1|2|
| 9293|0|5|2|5|-1|
3|0|4|
| 9510|55|8|1|2|2|
2|0|2|
| 10122|33|4|2|4|2|
3|0|2|
| 10549|0|4|2|4|2|
3|0|-1|
| 10812|0|4|2|4|-1|
```

```

2|      0|      -1|
| 10912|      0|      4|      2|      4|      2|
3|      0|      -1|
| 10996|      0|      5|      2|      5|      -1|
3|      0|      4|
| 11256|      8|      2|      2|      2|      1|
3|      0|      3|
| 11310|      31|      4|      2|      4|      1|
3|      0|      4|
+-----+-----+-----+-----+-----+-----+
+-----+-----+

```

only showing top 20 rows

```

root
|-- userId: integer (nullable = true)
|-- cms_segid: integer (nullable = true)
|-- cms_group_id: integer (nullable = true)
|-- final_gender_code: integer (nullable = true)
|-- age_level: integer (nullable = true)
|-- pvalue_level: string (nullable = true)
|-- shopping_level: integer (nullable = true)
|-- occupation: integer (nullable = true)
|-- new_user_class_level: string (nullable = true)

```

```

root
|-- userId: integer (nullable = true)
|-- cms_segid: integer (nullable = true)
|-- cms_group_id: integer (nullable = true)
|-- final_gender_code: integer (nullable = true)
|-- age_level: integer (nullable = true)
|-- pvalue_level: string (nullable = true)
|-- shopping_level: integer (nullable = true)
|-- occupation: integer (nullable = true)
|-- new_user_class_level: string (nullable = true)
|-- pl_onehot_feature: double (nullable = true)
|-- pl_onehot_value: vector (nullable = true)

```

```

+-----+-----+-----+-----+-----+-----+-----+
+-----+-----+-----+-----+-----+-----+
|userId|cms_segid|cms_group_id|final_gender_code|age_level|pvalue_level|shopping_level|
|occupation|new_user_class_level|pl_onehot_feature|pl_onehot_value|
+-----+-----+-----+-----+-----+-----+-----+
+-----+-----+-----+-----+-----+-----+
|      234|      0|      5|      2|      5|      -1|
3|      0|      3|      0.0| (4, [0], [1.0])|
|      523|      5|      2|      2|      2|      1|
3|      1|      2|      2.0| (4, [2], [1.0])|
|      612|      0|      8|      1|      2|      2|
3|      0|      -1|      1.0| (4, [1], [1.0])|
|      1670|      0|      4|      2|      4|      -1|
1|      0|      -1|      0.0| (4, [0], [1.0])|
|      2545|      0|      10|      1|      4|      -1|
3|      0|      -1|      0.0| (4, [0], [1.0])|
|      3644|      49|      6|      2|      6|      2|
3|      0|      2|      1.0| (4, [1], [1.0])|
|      5777|      44|      5|      2|      5|      2|
3|      0|      2|      1.0| (4, [1], [1.0])|

```

	6211	0	9	1	3	-1
3	0		2	0.0	(4, [0], [1.0])	
	6355	2	1	2	1	1
3	0		4	2.0	(4, [2], [1.0])	
	6823	43	5	2	5	2
3	0		1	1.0	(4, [1], [1.0])	
	6972	5	2	2	2	2
3	1		2	1.0	(4, [1], [1.0])	
	9293	0	5	2	5	-1
3	0		4	0.0	(4, [0], [1.0])	
	9510	55	8	1	2	2
2	0		2	1.0	(4, [1], [1.0])	
	10122	33	4	2	4	2
3	0		2	1.0	(4, [1], [1.0])	
	10549	0	4	2	4	2
3	0		-1	1.0	(4, [1], [1.0])	
	10812	0	4	2	4	-1
2	0		-1	0.0	(4, [0], [1.0])	
	10912	0	4	2	4	2
3	0		-1	1.0	(4, [1], [1.0])	
	10996	0	5	2	5	-1
3	0		4	0.0	(4, [0], [1.0])	
	11256	8	2	2	2	1
3	0		3	2.0	(4, [2], [1.0])	
	11310	31	4	2	4	1
3	0		4	2.0	(4, [2], [1.0])	
+-----+-----+-----+-----+-----+-----+-----+						
+-----+-----+-----+-----+-----+-----+-----+						

only showing top 20 rows

```
In [29]: 1 # 使用热编码转换new_user_class_level的一维数据为多维
2 stringindexer = StringIndexer(inputCol='new_user_class_level', outputCol='nucl_onehot_
3 encoder = OneHotEncoder(dropLast=False, inputCol='nucl_onehot_feature', outputCol='nuc
4 pipeline = Pipeline(stages=[stringindexer, encoder])
5 pipeline_fit = pipeline.fit(user_profile_df2)
6 user_profile_df3 = pipeline_fit.transform(user_profile_df2)
7 user_profile_df3.printSchema()
8 user_profile_df3.show()
```

```
In [30]: 1 # 如何将用户特征合并?
2 #例如我们要合并上表中的三列为一列: "age_level", "pl_onehot_value", "nucl_onehot_value"
3 from pyspark.ml.feature import VectorAssembler
4 feature_df = VectorAssembler().setInputCols(["age_level", "pl_onehot_value", "nucl_onehot_value"])
5 .transform(user_profile_df3)
6 feature_df.show()
```

```
+-----+-----+-----+-----+-----+-----+-----+
+-----+-----+-----+-----+-----+-----+-----+
+-----+-----+-----+-----+-----+-----+-----+
|userId|cms_segid|cms_group_id|final_gender_code|age_level|pvalue_level|shopping_level|occupation|new_user_class_level|pl_onehot_feature|pl_onehot_value|nucl_onehot_feature|nucl_onehot_value|features|
+-----+-----+-----+-----+-----+-----+-----+
| 234| 0| 5| 2| 5| -1|
3| 0| 3| 0.0| (4, [0], [1.0])|
2.0| (5, [2], [1.0])| (10, [0, 1, 7], [5.0, ...|
| 523| 5| 2| 2| 2| 1|
3| 1| 2| 2.0| (4, [2], [1.0])|
1.0| (5, [1], [1.0])| (10, [0, 3, 6], [2.0, ...|
| 612| 0| 8| 1| 2| 2|
3| 0| -1| 1.0| (4, [1], [1.0])|
0.0| (5, [0], [1.0])| (10, [0, 2, 5], [2.0, ...|
| 1670| 0| 4| 2| 4| -1|
1| 0| -1| 0.0| (4, [0], [1.0])|
0.0| (5, [0], [1.0])| (10, [0, 1, 5], [4.0, ...|
| 2545| 0| 10| 1| 4| -1|
3| 0| -1| 0.0| (4, [0], [1.0])|
0.0| (5, [0], [1.0])| (10, [0, 1, 5], [4.0, ...|
| 3644| 49| 6| 2| 6| 2|
3| 0| 2| 1.0| (4, [1], [1.0])|
1.0| (5, [1], [1.0])| (10, [0, 2, 6], [6.0, ...|
| 5777| 44| 5| 2| 5| 2|
3| 0| 2| 1.0| (4, [1], [1.0])|
1.0| (5, [1], [1.0])| (10, [0, 2, 6], [5.0, ...|
| 6211| 0| 9| 1| 3| -1|
3| 0| 2| 0.0| (4, [0], [1.0])|
1.0| (5, [1], [1.0])| (10, [0, 1, 6], [3.0, ...|
| 6355| 2| 1| 2| 1| 1|
3| 0| 4| 2.0| (4, [2], [1.0])|
3.0| (5, [3], [1.0])| (10, [0, 3, 8], [1.0, ...|
| 6823| 43| 5| 2| 5| 2|
3| 0| 1| 1.0| (4, [1], [1.0])|
4.0| (5, [4], [1.0])| (10, [0, 2, 9], [5.0, ...|
| 6972| 5| 2| 2| 2| 2|
3| 1| 2| 1.0| (4, [1], [1.0])|
1.0| (5, [1], [1.0])| (10, [0, 2, 6], [2.0, ...|
| 9293| 0| 5| 2| 5| -1|
3| 0| 4| 0.0| (4, [0], [1.0])|
3.0| (5, [3], [1.0])| (10, [0, 1, 8], [5.0, ...|
| 9510| 55| 8| 1| 2| 2|
2| 0| 2| 1.0| (4, [1], [1.0])|
1.0| (5, [1], [1.0])| (10, [0, 2, 6], [2.0, ...|
| 10122| 33| 4| 2| 4| 2|
3| 0| 2| 1.0| (4, [1], [1.0])|
```



```

1.0|      (5, [1], [1.0]) | (10, [0, 2, 6], [4.0, ... |
| 10549|      0|      4|
3|      0|      -1|
0.0|      (5, [0], [1.0]) | (10, [0, 2, 5], [4.0, ... |
| 10812|      0|      4|
2|      0|      -1|
0.0|      (5, [0], [1.0]) | (10, [0, 1, 5], [4.0, ... |
| 10912|      0|      4|
3|      0|      -1|
0.0|      (5, [0], [1.0]) | (10, [0, 2, 5], [4.0, ... |
| 10996|      0|      5|
3|      0|      4|
3.0|      (5, [3], [1.0]) | (10, [0, 1, 8], [5.0, ... |
| 11256|      8|      2|
3|      0|      3|
2.0|      (5, [2], [1.0]) | (10, [0, 3, 7], [2.0, ... |
| 11310|      31|      4|
3|      0|      4|
3.0|      (5, [3], [1.0]) | (10, [0, 3, 8], [4.0, ... |

```

```

2|      4|      2|
1.0|      (4, [1], [1.0]) |
2|      4|      -1|
0.0|      (4, [0], [1.0]) |
2|      4|      2|
1.0|      (4, [1], [1.0]) |
2|      5|      -1|
0.0|      (4, [0], [1.0]) |
2|      2|      1|
2.0|      (4, [2], [1.0]) |
2|      4|      1|
2.0|      (4, [2], [1.0]) |

```

```

+-----+-----+-----+-----+-----+-----+
+-----+-----+-----+-----+-----+-----+
+-----+-----+-----+-----+-----+-----+

```

only showing top 20 rows

In [31]: 1 feature\_df.select('features').show()

```

+-----+
|      features|
+-----+
| (10, [0, 1, 7], [5.0, ... |
| (10, [0, 3, 6], [2.0, ... |
| (10, [0, 2, 5], [2.0, ... |
| (10, [0, 1, 5], [4.0, ... |
| (10, [0, 1, 5], [4.0, ... |
| (10, [0, 2, 6], [6.0, ... |
| (10, [0, 2, 6], [5.0, ... |
| (10, [0, 1, 6], [3.0, ... |
| (10, [0, 3, 8], [1.0, ... |
| (10, [0, 2, 9], [5.0, ... |
| (10, [0, 2, 6], [2.0, ... |
| (10, [0, 1, 8], [5.0, ... |
| (10, [0, 2, 6], [2.0, ... |
| (10, [0, 2, 6], [4.0, ... |
| (10, [0, 2, 5], [4.0, ... |
| (10, [0, 1, 5], [4.0, ... |
| (10, [0, 2, 5], [4.0, ... |

```

```
In [ ]: 1 # 我该选取那些特征呢?
2 # 除了前面处理的pvalue_level和new_user_class_level需要作为特征以外，（能体现出用户的购
3 # 分类特征值：个数
4 # - cms_segid: 97
5 # - cms_group_id: 13
6 # - final_gender_code: 2
7 # - age_level: 7
8 # - shopping_level: 3
9 # - occupation: 2
10 # - pvalue_level
11 # - new_user_class_level
12 # - price
13 # 根据经验，以上几个分类特征都一定程度能体现用户在购物方面的特征，且类别都较少，都可以
```