

基于LR的点击率预测模型训练

本小节主要根据广告点击样本数据集(raw_sample)、广告基本特征数据集(ad_feature)、用户基本信息数据集(user_profile)构建出了一个完整的样本数据集，并按日期划分为了训练集(前七天)和测试集(最后一天)，利用逻辑回归进行训练。

训练模型时，通过对类别特征数据进行处理，一定程度达到提高了模型的效果

分析对比两种模型：

1. 训练CTRModel_Normal：直接将对应的特征的特征值组合成对应的特征向量进行训练

```
In [ ]: 1 '''
2 # 剔除冗余、不需要的字段
3 useful_cols = [
4     #
5     # 时间字段，划分训练集和测试集
6     "timestamp",
7     # label目标值字段
8     "clk",
9     # 特征值字段
10    "pid_value",      # 资源位的特征向量
11    "price",          # 广告价格
12    "cms_segid",      # 用户微群ID
13    "cms_group_id",   # 用户组ID
14    "final_gender_code", # 用户性别特征, [1,2]
15    "age_level",      # 年龄等级, 1-
16    "shopping_level",
17    "occupation",
18    "pl_onehot_value",
19    "nucl_onehot_value"
20 ]
21 # 筛选指定字段数据，构建新的数据集
22 datasets_l = datasets.select(*useful_cols)
23 '''
```

2. 训练CTRModel_AllOneHot

- "pid_value", 类别型特征，已被转换为多维特征==> 2维
- "price", 统计型特征 ==> 1维
- "cms_segid", 类别型特征，约97个分类 ==> 1维
- "cms_group_id", 类别型特征，约13个分类 ==> 1维
- "final_gender_code", 类别型特征，2个分类 ==> 1维
- "age_level", 类别型特征，7个分类 ==> 1维
- "shopping_level", 类别型特征，3个分类 ==> 1维
- "occupation", 类别型特征，2个分类 ==> 1维
- "pl_onehot_value", 类别型特征，已被转换为多维特征 ==> 4维
- "nucl_onehot_value" 类别型特征，已被转换为多维特征 ==> 5维

类别性特征都可以考虑进行热独编码，将单一变量变为多变量，相当于增加了相关特征的数量

- "cms_segid", 类别型特征, 约97个分类 ==> 97维 舍弃
- "cms_group_id", 类别型特征, 约13个分类 ==> 13维
- "final_gender_code", 类别型特征, 2个分类 ==> 2维
- "age_level", 类别型特征, 7个分类 ==> 7维
- "shopping_level", 类别型特征, 3个分类 ==> 3维
- "occupation", 类别型特征, 2个分类 ==> 2维

但由于cms_segid分类过多，这里考虑舍弃，避免数据过于稀疏

结论:

对比前面的result_1的预测结果，能发现这里的预测率稍微准确了一点，这里top20里出现了3个点击的，但前面的只出现了1个

因此可见对特征的细化处理，已经帮助我们提高模型的效果的

In []:

1

In []:

1

In []:

```
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 = 'createCTRModelByLR'
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 #     ('spark.executor.instances', 1) # 设置spark executor数量, yarn时起作用
23 )
24 conf.setAll(config)
25
26 spark = SparkSession.builder.config(conf=conf).getOrCreate()
```

```
In [4]: 1 '''
2 raw_sample
3         pid
4 ad_feature
5         price
6 user_profile
7         - cms_segid: 97
8         - cms_group_id: 13
9         - final_gender_code: 2
10        - age_level: 7
11        - shopping_level: 3
12        - occupation: 2
13        - pvalue_level
14        - new_user_class_level
15        '''
```

1.raw_sample - pid

```
In [6]: 1 _raw_sample_df1 = spark.read.csv('/data/raw_sample.csv', header=True)
        2 _raw_sample_df1.show()
        3 _raw_sample_df1.printSchema()
```

user	time_stamp	adgroup_id	pid	nonclk	clk
581738	1494137644	1	430548_1007	1	0
449818	1494638778	3	430548_1007	1	0
914836	1494650879	4	430548_1007	1	0
914836	1494651029	5	430548_1007	1	0
399907	1494302958	8	430548_1007	1	0
628137	1494524935	9	430548_1007	1	0
298139	1494462593	9	430539_1007	1	0
775475	1494561036	9	430548_1007	1	0
555266	1494307136	11	430539_1007	1	0
117840	1494036743	11	430548_1007	1	0
739815	1494115387	11	430539_1007	1	0
623911	1494625301	11	430548_1007	1	0
623911	1494451608	11	430548_1007	1	0
421590	1494034144	11	430548_1007	1	0
976358	1494156949	13	430548_1007	1	0
286630	1494218579	13	430539_1007	1	0
286630	1494289247	13	430539_1007	1	0
771431	1494153867	13	430548_1007	1	0
707120	1494220810	13	430548_1007	1	0
530454	1494293746	13	430548_1007	1	0

only showing top 20 rows

```
root
|-- user: string (nullable = true)
|-- time_stamp: string (nullable = true)
|-- adgroup_id: string (nullable = true)
|-- pid: string (nullable = true)
|-- nonclk: string (nullable = true)
|-- clk: string (nullable = true)
```

```
In [7]: 1 from pyspark.sql.types import StringType, StructField, IntegerType, FloatType, LongType
2 _raw_sample_df2 = _raw_sample_df1.withColumn('user', _raw_sample_df1.user.cast(IntegerType()))
3     withColumn('time_stamp', _raw_sample_df1.time_stamp.cast(LongType())).withColumnRenamed('time_stamp', 'timestamp')
4     withColumn("adgroup_id", _raw_sample_df1.adgroup_id.cast(IntegerType())).withColumnRenamed('adgroup_id', 'adgroup_id')
5     withColumn("pid", _raw_sample_df1.pid.cast(StringType())).withColumnRenamed('pid', 'pid')
6     withColumn("nonclk", _raw_sample_df1.nonclk.cast(IntegerType())).withColumnRenamed('nonclk', 'nonclk')
7     withColumn("clk", _raw_sample_df1.clk.cast(IntegerType())).withColumnRenamed('clk', 'clk')
8 _raw_sample_df2.printSchema()
9 _raw_sample_df2.show()
```

```
root
```

```
-- userId: integer (nullable = true)
-- timestamp: long (nullable = true)
-- adgroupId: integer (nullable = true)
-- pid: string (nullable = true)
-- nonclk: integer (nullable = true)
-- clk: integer (nullable = true)
```

userId	timestamp	adgroupId	pid	nonclk	clk
581738	1494137644	1	430548_1007	1	0
449818	1494638778	3	430548_1007	1	0
914836	1494650879	4	430548_1007	1	0
914836	1494651029	5	430548_1007	1	0
399907	1494302958	8	430548_1007	1	0
628137	1494524935	9	430548_1007	1	0
298139	1494462593	9	430539_1007	1	0
775475	1494561036	9	430548_1007	1	0
555266	1494307136	11	430539_1007	1	0
117840	1494036743	11	430548_1007	1	0
739815	1494115387	11	430539_1007	1	0
623911	1494625301	11	430548_1007	1	0
623911	1494451608	11	430548_1007	1	0
421590	1494034144	11	430548_1007	1	0
976358	1494156949	13	430548_1007	1	0
286630	1494218579	13	430539_1007	1	0
286630	1494289247	13	430539_1007	1	0
771431	1494153867	13	430548_1007	1	0
707120	1494220810	13	430548_1007	1	0
530454	1494293746	13	430548_1007	1	0

```
only showing top 20 rows
```

```
In [8]: 1 from pyspark.ml.feature import OneHotEncoder
2 from pyspark.ml.feature import StringIndexer
3 from pyspark.ml import Pipeline
4 stringindexer = StringIndexer(inputCol='pid', outputCol='pid_feature')
5 encoder = OneHotEncoder(dropLast=False, inputCol='pid_feature', outputCol='pid_value')
6 pipeline = Pipeline(stages=[stringindexer, encoder])
7 pipeline_fit= pipeline.fit(_raw_sample_df2)
8 raw_sample_df = pipeline_fit.transform(_raw_sample_df2)
9 raw_sample_df.show()
10
11 '''pid和特征的对应关系
12 430548_1007: 0
13 430549_1007: 1
14 '''
```

userId	timestamp	adgroupId	pid	nonclk	clk	pid_feature	pid_value
581738	1494137644	1	430548_1007	1	0	0.0	(2, [0], [1.0])
449818	1494638778	3	430548_1007	1	0	0.0	(2, [0], [1.0])
914836	1494650879	4	430548_1007	1	0	0.0	(2, [0], [1.0])
914836	1494651029	5	430548_1007	1	0	0.0	(2, [0], [1.0])
399907	1494302958	8	430548_1007	1	0	0.0	(2, [0], [1.0])
628137	1494524935	9	430548_1007	1	0	0.0	(2, [0], [1.0])
298139	1494462593	9	430539_1007	1	0	1.0	(2, [1], [1.0])
775475	1494561036	9	430548_1007	1	0	0.0	(2, [0], [1.0])
555266	1494307136	11	430539_1007	1	0	1.0	(2, [1], [1.0])
117840	1494036743	11	430548_1007	1	0	0.0	(2, [0], [1.0])
739815	1494115387	11	430539_1007	1	0	1.0	(2, [1], [1.0])
623911	1494625301	11	430548_1007	1	0	0.0	(2, [0], [1.0])
623911	1494451608	11	430548_1007	1	0	0.0	(2, [0], [1.0])
421590	1494034144	11	430548_1007	1	0	0.0	(2, [0], [1.0])
976358	1494156949	13	430548_1007	1	0	0.0	(2, [0], [1.0])
286630	1494218579	13	430539_1007	1	0	1.0	(2, [1], [1.0])
286630	1494289247	13	430539_1007	1	0	1.0	(2, [1], [1.0])
771431	1494153867	13	430548_1007	1	0	0.0	(2, [0], [1.0])
707120	1494220810	13	430548_1007	1	0	0.0	(2, [0], [1.0])
530454	1494293746	13	430548_1007	1	0	0.0	(2, [0], [1.0])

only showing top 20 rows

Out[8]: 'pid和特征的对应关系\n430548_1007: 0\n430549_1007: 1\n'

2.广告基本信息ad_feature - price

```
In [10]: 1 _ad_feature_df = spark.read.csv('/data/ad_feature.csv', header=True)
2         _ad_feature_df.printSchema()
3         _ad_feature_df.show()
```

```
root
|-- adgroup_id: string (nullable = true)
|-- cate_id: string (nullable = true)
|-- campaign_id: string (nullable = true)
|-- customer: string (nullable = true)
|-- brand: string (nullable = true)
|-- price: string (nullable = true)
```

adgroup_id	cate_id	campaign_id	customer	brand	price
63133	6406	83237	1	95471	170.0
313401	6406	83237	1	87331	199.0
248909	392	83237	1	32233	38.0
208458	392	83237	1	174374	139.0
110847	7211	135256	2	145952	32.99
607788	6261	387991	6	207800	199.0
375706	4520	387991	6	NULL	99.0
11115	7213	139747	9	186847	33.0
24484	7207	139744	9	186847	19.0
28589	5953	395195	13	NULL	428.0
23236	5953	395195	13	NULL	368.0
300556	5953	395195	13	NULL	639.0
92560	5953	395195	13	NULL	368.0
590965	4284	28145	14	454237	249.0
529913	4284	70206	14	NULL	249.0
546930	4284	28145	14	NULL	249.0
639794	6261	70206	14	37004	89.9
335413	4284	28145	14	NULL	249.0
794890	4284	70206	14	454237	249.0
684020	6261	70206	14	37004	99.0

only showing top 20 rows

```
In [11]: 1 from pyspark.sql.types import IntegerType, FloatType
2 ad_feature_df = _ad_feature_df.\
3     withColumn("adgroup_id", _ad_feature_df.adgroup_id.cast(IntegerType())).withColumn
4     withColumn("cate_id", _ad_feature_df.cate_id.cast(IntegerType())).withColumnRename
5     withColumn("campaign_id", _ad_feature_df.campaign_id.cast(IntegerType())).withColu
6     withColumn("customer", _ad_feature_df.customer.cast(IntegerType())).withColumnRenam
7     withColumn("brand", _ad_feature_df.brand.cast(IntegerType())).withColumnRenamed("b
8     withColumn("price", _ad_feature_df.price.cast(FloatType()))
9 ad_feature_df.printSchema()
10 ad_feature_df.show()
```

```
root
|-- adgroupId: integer (nullable = true)
|-- cateId: integer (nullable = true)
|-- campaignId: integer (nullable = true)
|-- customerId: integer (nullable = true)
|-- brandId: integer (nullable = true)
|-- price: float (nullable = true)
```

adgroupId	cateId	campaignId	customerId	brandId	price
63133	6406	83237	1	95471	170.0
313401	6406	83237	1	87331	199.0
248909	392	83237	1	32233	38.0
208458	392	83237	1	174374	139.0
110847	7211	135256	2	145952	32.99
607788	6261	387991	6	207800	199.0
375706	4520	387991	6	null	99.0
11115	7213	139747	9	186847	33.0
24484	7207	139744	9	186847	19.0
28589	5953	395195	13	null	428.0
23236	5953	395195	13	null	368.0
300556	5953	395195	13	null	639.0
92560	5953	395195	13	null	368.0
590965	4284	28145	14	454237	249.0
529913	4284	70206	14	null	249.0
546930	4284	28145	14	null	249.0
639794	6261	70206	14	37004	89.9
335413	4284	28145	14	null	249.0
794890	4284	70206	14	454237	249.0
684020	6261	70206	14	37004	99.0

only showing top 20 rows

3. user_profile

- cms_segid: 97
- cms_group_id: 13
- final_gender_code: 2
- age_level: 7
- shopping_level: 3
- occupation: 2

- pvalue_level
- new_user_class_level

```
In [13]: 1 _user_profile_df = spark.read.csv('/data/user_profile.csv', header=True)
2         _user_profile_df.printSchema()
3         _user_profile_df.show()
```

```
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)
```

```
+-----+-----+-----+-----+-----+-----+-----+
|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|
0|
| 523| 5| 2| 2| 2| 1| 3|
1|
| 612| 0| 8| 1| 2| 2| 3|
0| null|
| 1670| 0| 4| 2| 4| null| 1|
0| null|
| 2545| 0| 10| 1| 4| null| 3|
0| null|
| 3644| 49| 6| 2| 6| 2| 3|
0| 2|
| 5777| 44| 5| 2| 5| 2| 3|
0| 2|
| 6211| 0| 9| 1| 3| null| 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| null| 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| null|
| 10812| 0| 4| 2| 4| null| 2|
0| null|
| 10912| 0| 4| 2| 4| 2| 3|
0| null|
| 10996| 0| 5| 2| 5| null| 3|
0| 4|
```

0	11256	8	2	2	2	1	3
0	11310	31	4	2	4	1	3
0		4					
+-----+-----+							
-----+							

only showing top 20 rows

```
In [14]: 1 # 查看每列数据中有没有 'NULL', 如果有, 就不能使用schema, 否则就会使那一整行就变成 null;如
        2 # 注意:"null" 与 "NULL"
        3 [str(c) + ':' + str(_user_profile_df.groupby(c).count().show()) for c in _user_profile_df.columns]
        4 # 根据结果可见仅仅有null, 而没有NULL, 因此可以使用schema
```

	userid	count
	505039	1
	577511	1
	627835	1
	692974	1
	742322	1
	746750	1
	777511	1
	800757	1
	878358	1
	976473	1
	1141237	1
	34635	1
	265095	1
	308633	1
	344922	1
	472235	1
	618816	1

```
In [15]: 1 from pyspark.sql.types import StructType, StructField, IntegerType
2 schema = StructType([
3     StructField("userId", IntegerType()),
4     StructField("cms_segid", IntegerType()),
5     StructField("cms_group_id", IntegerType()),
6     StructField("final_gender_code", IntegerType()),
7     StructField("age_level", IntegerType()),
8     StructField("pvalue_level", IntegerType()),
9     StructField("shopping_level", IntegerType()),
10    StructField("occupation", IntegerType()),
11    StructField("new_user_class_level", IntegerType())
12 ])
13 _user_profile_df1 = spark.read.csv('/data/user_profile.csv', header=True, schema=schema)
14 _user_profile_df1.printSchema()
15 _user_profile_df1.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|
|occupation|new_user_class_level|
+-----+-----+-----+-----+-----+-----+-----+
| 234| 0| 5| 2| 5| null| 3|
0| 3|
| 523| 5| 2| 2| 2| 1| 3|
1| 2|
| 612| 0| 8| 1| 2| 2| 3|
0| null|
| 1670| 0| 4| 2| 4| null| 1|
0| null|
| 2545| 0| 10| 1| 4| null| 3|
0| null|
| 3644| 49| 6| 2| 6| 2| 3|
0| 2|
| 5777| 44| 5| 2| 5| 2| 3|
0| 2|
| 6211| 0| 9| 1| 3| null| 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| null| 3|
0| 4|
```

0	9510	55	8	1	2	2	2
0	10122	33	4	2	4	2	3
0	10549	0	4	2	4	2	3
0		null					
0	10812	0	4	2	4	null	2
0		null					
0	10912	0	4	2	4	2	3
0		null					
0	10996	0	5	2	5	null	3
0		4					
0	11256	8	2	2	2	1	3
0		3					
0	11310	31	4	2	4	1	3
0		4					
+-----+-----+-----+-----+-----+-----+-----+							
-----+-----+-----+-----+							

only showing top 20 rows

[3], [1.0])							
6823	43	5	2	5	2	3	
0	1	1.0	(4, [1], [1.0])		4.0	(5,	
[4], [1.0])							
6972	5	2	2	2	2	3	
1	2	1.0	(4, [1], [1.0])		1.0	(5,	
[1], [1.0])							
9293	0	5	2	5	-1	3	
0	4	0.0	(4, [0], [1.0])		3.0	(5,	
[3], [1.0])							
9510	55	8	1	2	2	2	
0	2	1.0	(4, [1], [1.0])		1.0	(5,	
[1], [1.0])							
10122	33	4	2	4	2	3	
0	2	1.0	(4, [1], [1.0])		1.0	(5,	
[1], [1.0])							
10549	0	4	2	4	2	3	
0	-1	1.0	(4, [1], [1.0])		0.0	(5,	
[0], [1.0])							
10812	0	4	2	4	-1	2	
0	-1	0.0	(4, [0], [1.0])		0.0	(5,	
[0], [1.0])							
10912	0	4	2	4	2	3	
0	-1	1.0	(4, [1], [1.0])		0.0	(5,	
[0], [1.0])							
10996	0	5	2	5	-1	3	
0	4	0.0	(4, [0], [1.0])		3.0	(5,	
[3], [1.0])							
11256	8	2	2	2	1	3	
0	3	2.0	(4, [2], [1.0])		2.0	(5,	
[2], [1.0])							
11310	31	4	2	4	1	3	
0	4	2.0	(4, [2], [1.0])		3.0	(5,	
[3], [1.0])							
+-----+-----+-----+-----+-----+-----+-----+							
-----+-----+-----+-----+-----+-----+-----+							
-----+							

only showing top 20 rows

```
In [17]: 1 # 找出两者的映射关系 max min 都是一个值!
2 user_profile_df.groupby('pvalue_level').max('pl_onehot_feature').show()
3 user_profile_df.groupBy("new_user_class_level").max("nucl_onehot_feature").show()
```

pvalue_level	max(pl_onehot_feature)
-1	0.0
3	3.0
1	2.0
2	1.0

new_user_class_level	max(nucl_onehot_feature)
-1	0.0
3	2.0
1	4.0
4	3.0
2	1.0

4.raw_sample表(包含userId和广告Id) 合并user_profile和ad_feature表

Dataframe数据合并: [pyspark.sql.DataFrame.join](https://spark.apache.org/docs/latest/api/python/pyspark.sql.html?highlight=join#pyspark.sql.DataFrame.join)
<https://spark.apache.org/docs/latest/api/python/pyspark.sql.html?highlight=join#pyspark.sql.DataFrame.join>

不同合并方式介绍 (<https://stackoverflow.com/questions/38549/what-is-the-difference-between-inner-join-and-outer-join>)

```
In [28]: 1 # 由此可见 三张表的useid和adid个数是不一样的
2 print(raw_sample_df.count())
3 print(ad_feature_df.count())
4 print(user_profile_df.count())
5 print('*'*10)
6 print(raw_sample_df.groupBy('adgroupId').count().count())
7 print(ad_feature_df.groupBy('adgroupId').count().count())
8 print(raw_sample_df.groupBy('userId').count().count())
9 print(user_profile_df.groupBy('userId').count().count())
```

```
26557961
846811
1061768
*****
846811
846811
1141729
1061768
```



```
In [20]: 1 condition = [raw_sample_df.adgroupId == ad_feature_df.adgroupId]
2         _ = raw_sample_df.join(ad_feature_df, on=condition, how='outer')
3
4         condition2 = [_.userId == user_profile_df.userId]
5         datasets = _.join(user_profile_df, on=condition2, how='outer')
6
7         datasets.printSchema()
8         print(datasets.count())
```

```
root
|-- userId: integer (nullable = true)
|-- timestamp: long (nullable = true)
|-- adgroupId: integer (nullable = true)
|-- pid: string (nullable = true)
|-- nonclk: integer (nullable = true)
|-- clk: integer (nullable = true)
|-- pid_feature: double (nullable = true)
|-- pid_value: vector (nullable = true)
|-- adgroupId: integer (nullable = true)
|-- cateId: integer (nullable = true)
|-- campaignId: integer (nullable = true)
|-- customerId: integer (nullable = true)
|-- brandId: integer (nullable = true)
|-- price: float (nullable = true)
|-- 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)
|-- nucl_onehot_feature: double (nullable = true)
|-- nucl_onehot_value: vector (nullable = true)
```

26557961

1. 训练CTRModel Normal: 直接将对应的特征的特征值组合成对应的特征向量进行训练

```
In [21]: 1 # 延申学习: 作为条件的那些列不能被select
2         # datasets.select('nucl_onehot_feature')
```

```
In [23]: 1 # 剔除冗余、不需要的字段
2 useful_cols = [
3     # 时间字段, 划分训练集和测试集
4     "timestamp",
5     # label目标值字段
6     "clk",
7     # 特征值字段
8     "pid_value",      # 资源位的特征向量
9     "price",          # 广告价格
10    "cms_segid",       # 用户微群ID
11    "cms_group_id",    # 用户组ID
12    "final_gender_code", # 用户性别特征, [1,2]
13    "age_level",       # 年龄等级, 1-
14    "shopping_level",
15    "occupation",
16    "pl_onehot_value",
17    "nucl_onehot_value"
18 ]
19 datasets_1 = datasets.select(*[useful_cols])
20 datasets_1.printSchema()
```

```
root
|-- timestamp: long (nullable = true)
|-- clk: integer (nullable = true)
|-- pid_value: vector (nullable = true)
|-- price: float (nullable = true)
|-- cms_segid: integer (nullable = true)
|-- cms_group_id: integer (nullable = true)
|-- final_gender_code: integer (nullable = true)
|-- age_level: integer (nullable = true)
|-- shopping_level: integer (nullable = true)
|-- occupation: integer (nullable = true)
|-- pl_onehot_value: vector (nullable = true)
|-- nucl_onehot_value: vector (nullable = true)
```

```
In [25]: 1 # 三张表行数不同, 合并后肯定有空值, 要去掉空值
2 # str类型的空值(null或者NULL)不能被dropna()掉
3 # str类型的NULL转化为 非str类型后, show()会显示null
4 datasets_1=datasets_1.dropna()
5 print("剔除空值数据后, 还剩: ", datasets_1.count())
```

剔除空值数据后, 还剩: 25029435

根据特征字段计算出特征向量, 并划分出训练数据集和测试数据集

```
In [ ]: 1 ##### 根据特征字段计算出特征向量, 并划分出训练数据集和测试数据集# 延申学习: 找到最大时间
2 # datasets_1.orderBy('timestamp', ascending=False).show()
3 # #1494691186
```

```
In [32]: 1 # 根据特征字段 计算出特征向量, 并划分出 训练数据集和测试数据集
2 from pyspark.ml.feature import VectorAssembler
3 datasets_1 = VectorAssembler().setInputCols(usable_cols[2:]).setOutputCol('features')
4 #训练数据集
5 train_datasets_1 = datasets_1.filter(datasets_1.timestamp <= (1494691186-24*60*60))
6 #测试数据集
7 test_datasets_1 = datasets_1.where(datasets_1.timestamp > (1494691186-24*60*60))
8 # 所有特征的特征向量已经汇总在features字段中
9 train_datasets_1.show(5)
10 test_datasets_1.show(5)
```

...

创建逻辑回归训练器, 并训练模型: [LogisticRegression](https://spark.apache.org/docs/latest/api/python/pyspark.ml.html?highlight=logisticregression#pyspark.ml.classification.LogisticRegression)

(<https://spark.apache.org/docs/latest/api/python/pyspark.ml.html?highlight=logisticregression#pyspark.ml.classification.LogisticRegression>)、
[LogisticRegressionModel](https://spark.apache.org/docs/latest/api/python/pyspark.ml.html?highlight=logisticregression#pyspark.ml.classification.LogisticRegressionModel) (<https://spark.apache.org/docs/latest/api/python/pyspark.ml.html?highlight=logisticregression#pyspark.ml.classification.LogisticRegressionModel>)

```
In [33]: 1 from pyspark.ml.classification import LogisticRegression
2 lr = LogisticRegression()
3 model = lr.setLabelCol('clk').setFeaturesCol('features').fit(train_datasets_1)
4 # model.save(hadoop上)
5
6 # 练习使用下面的类, 使用已经训练好的模型
7 from pyspark.ml.classification import LogisticRegressionModel
8 model = LogisticRegressionModel.load('/models/CTRModel_Normal.obj')
9 result_1 = model.transform(test_datasets_1)
10 result_1.show()
```

...

```
In [34]: 1 result_1.select('clk','price','probability','prediction').sort('probability').show(100)
2 # 预测的前20个, 命中了3个
```

...

```
In [35]: 1 # 只查看样本中点击的被实际点击的条目的预测情况
2 result_1.select('clk','price','probability','prediction').where('clk==1').sort('probability').show(100)
3 # 默认按照概率的50%进行分类, 大于0.5预测为1, 小于0.5预测为0
```

...

2. 训练CTRModel_AllOneHot

- "pid_value", 类别型特征, 已被转换为多维特征==> 2维
- "price", 统计型特征 ==> 1维
- "cms_segid", 类别型特征, 约97个分类 ==> 1维
- "cms_group_id", 类别型特征, 约13个分类 ==> 1维
- "final_gender_code", 类别型特征, 2个分类 ==> 1维
- "age_level", 类别型特征, 7个分类 ==> 1维
- "shopping_level", 类别型特征, 3个分类 ==> 1维
- "occupation", 类别型特征, 2个分类 ==> 1维
- "pl_onehot_value", 类别型特征, 已被转换为多维特征 ==> 4维

- "nucl_onehot_value" 类别型特征, 已被转换为多维特征 ==> 5维

类别性特征都可以考虑进行热独编码, 将单一变量变为多变量, 相当于增加了相关特征的数量

- "cms_segid", 类别型特征, 约97个分类 ==> 97维 舍弃
- "cms_group_id", 类别型特征, 约13个分类 ==> 13维
- "final_gender_code", 类别型特征, 2个分类 ==> 2维
- "age_level", 类别型特征, 7个分类 ==> 7维
- "shopping_level", 类别型特征, 3个分类 ==> 3维
- "occupation", 类别型特征, 2个分类 ==> 2维

但由于cms_segid分类过多, 这里考虑舍弃, 避免数据过于稀疏

```
In [42]: 1 # 首先查看每个特征列 有几个不同的值, n个不同的值就可以转换为n维
2 #=====耗时4min=====
3 [str(c) + '特征的种类个数:' + str(datasets_1.groupby(c).count().count()) for c in datasets_1.columns]
```

```
Out[42]: ['pid_value特征的种类个数:2',
'cms_segid特征的种类个数:97',
'cms_group_id特征的种类个数:13',
'final_gender_code特征的种类个数:2',
'age_level特征的种类个数:7',
'shopping_level特征的种类个数:3',
'occupation特征的种类个数:2',
'pl_onehot_value特征的种类个数:4',
'nucl_onehot_value特征的种类个数:5']
```

```
In [43]: 1 datasets_1.first()
```

```
Out[43]: Row(timestamp=1494261938, clk=0, pid_value=SparseVector(2, {1: 1.0}), price=108.0, cms_segid=0, cms_group_id=11, final_gender_code=1, age_level=5, shopping_level=3, occupation=0, pl_onehot_value=SparseVector(4, {0: 1.0}), nucl_onehot_value=SparseVector(5, {1: 1.0}), features=SparseVector(18, {1: 1.0, 2: 108.0, 4: 11.0, 5: 1.0, 6: 5.0, 7: 3.0, 9: 1.0, 14: 1.0}))
```

```

In [53]: 1 # 先将下列五列数据转为字符串类型，以便于进行热独编码
2 # - "cms_group_id", 类别型特征，约13个分类 ==> 13
3 # - "final_gender_code", 类别型特征，2个分类 ==> 2
4 # - "age_level", 类别型特征，7个分类 ==> 7
5 # - "shopping_level", 类别型特征，3个分类 ==> 3
6 # - "occupation", 类别型特征，2个分类 ==> 2
7 ## datasets 三张表合并生成的表
8 datasets_2 = datasets.withColumn("cms_group_id", datasets.cms_group_id.cast(StringType))\
9     .withColumn("final_gender_code", datasets.final_gender_code.cast(StringType()))\
10    .withColumn("age_level", datasets.age_level.cast(StringType()))\
11    .withColumn("shopping_level", datasets.shopping_level.cast(StringType()))\
12    .withColumn("occupation", datasets.occupation.cast(StringType()))
13 useful_cols_2 = [
14     # 时间值，划分训练集和测试集
15     "timestamp",
16     # label目标值
17     "clk",
18     # 特征值
19     "price",
20     "cms_group_id",
21     "final_gender_code",
22     "age_level",
23     "shopping_level",
24     "occupation",
25     "pid_value",
26     "pl_onehot_value",
27     "nucl_onehot_value"
28 ]
29 datasets_2 = datasets_2.select(*useful_cols_2)
30 datasets_2 = datasets_2.dropna()

```

```

In [54]: 1 from pyspark.ml.feature import OneHotEncoder
2 from pyspark.ml.feature import StringIndexer
3 from pyspark.ml import Pipeline
4
5 def oneHotEncoder(col1, col2, col3, data):
6     stringindexer = StringIndexer(inputCol=col1, outputCol=col2)
7     encoder = OneHotEncoder(dropLast=False, inputCol=col2, outputCol=col3)
8     pipeline = Pipeline(stages=[stringindexer, encoder])
9     pipeline_fit = pipeline.fit(data)
10    return pipeline_fit.transform(data)
11 # 对以下5个特征进行独热编码
12 # - "cms_group_id", 类别型特征，约13个分类 ==> 13
13 # - "final_gender_code", 类别型特征，2个分类 ==> 2
14 # - "age_level", 类别型特征，7个分类 ==> 7
15 # - "shopping_level", 类别型特征，3个分类 ==> 3
16 # - "occupation", 类别型特征，2个分类 ==> 2
17 datasets_2 = oneHotEncoder("cms_group_id", "cms_group_id_feature", "cms_group_id_value", datasets_2)
18 datasets_2 = oneHotEncoder("final_gender_code", "final_gender_code_feature", "final_gender_code_value", datasets_2)
19 datasets_2 = oneHotEncoder("age_level", "age_level_feature", "age_level_value", datasets_2)
20 datasets_2 = oneHotEncoder("shopping_level", "shopping_level_feature", "shopping_level_value", datasets_2)
21 datasets_2 = oneHotEncoder("occupation", "occupation_feature", "occupation_value", datasets_2)

```

```
In [55]: 1 # onehot编码完成后, 查看下对应关系
2 # min max都是一个值, 比如 1的min和max是一个值
3 datasets_2.groupBy("cms_group_id").min("cms_group_id_feature").show()
4 datasets_2.groupBy("final_gender_code").min("final_gender_code_feature").show()
5 datasets_2.groupBy("age_level").min("age_level_feature").show()
6 datasets_2.groupBy("shopping_level").min("shopping_level_feature").show()
7 datasets_2.groupBy("occupation").min("occupation_feature").show()
```

...

```
In [ ]: 1 # 独热编码后, 特征字段不再是之前的字段, 重新定义字段
2 feature_cols = [
3     # 特征值
4     "price",
5     "cms_group_id_value",
6     "final_gender_code_value",
7     "age_level_value",
8     "shopping_level_value",
9     "occupation_value",
10    "pid_value",
11    "pl_onehot_value",
12    "nucl_onehot_value"
13 ]
14 # 根据特征字段计算出特征向量, 并划分出训练数据集和测试数据集
15 from pyspark.ml.feature import VectorAssembler
16 datasets_2 = VectorAssembler().setInputCols(feature_cols).setOutputCol('features').tra
```

```
In [64]: 1 # 训练样本集
2 train_datasets_2 = datasets_2.filter(datasets_2.timestamp<=(1494691186-24*60*60))
3 # 测试样本集
4 test_datasets_2 = datasets_2.where(datasets_2.timestamp>(1494691186-24*60*60))
5 train_datasets_2.printSchema()
6 train_datasets_2.first()
7 # features=SparseVector(39, {0: 108.0, 7: 1.0, 15: 1.0, 18: 1.0, 23: 1.0, 26: 1.0, 29:
```

```
root
|-- timestamp: long (nullable = true)
|-- clk: integer (nullable = true)
|-- price: float (nullable = true)
|-- cms_group_id: string (nullable = true)
|-- final_gender_code: string (nullable = true)
|-- age_level: string (nullable = true)
|-- shopping_level: string (nullable = true)
|-- occupation: string (nullable = true)
|-- pid_value: vector (nullable = true)
|-- pl_onehot_value: vector (nullable = true)
|-- nucl_onehot_value: vector (nullable = true)
|-- cms_group_id_feature: double (nullable = true)
|-- cms_group_id_value: vector (nullable = true)
|-- final_gender_code_feature: double (nullable = true)
|-- final_gender_code_value: vector (nullable = true)
|-- age_level_feature: double (nullable = true)
|-- age_level_value: vector (nullable = true)
|-- shopping_level_feature: double (nullable = true)
|-- shopping_level_value: vector (nullable = true)
|-- occupation_feature: double (nullable = true)
|-- occupation_value: vector (nullable = true)
|-- features: vector (nullable = true)
```

```
Out[64]: Row(timestamp=1494261938, clk=0, price=108.0, cms_group_id='11', final_gender_code='1',
age_level='5', shopping_level='3', occupation='0', pid_value=SparseVector(2, {1: 1.0}),
pl_onehot_value=SparseVector(4, {0: 1.0}), nucl_onehot_value=SparseVector(5, {1: 1.0}),
cms_group_id_feature=6.0, cms_group_id_value=SparseVector(13, {6: 1.0}), final_gender_c
ode_feature=1.0, final_gender_code_value=SparseVector(2, {1: 1.0}), age_level_feature=
2.0, age_level_value=SparseVector(7, {2: 1.0}), shopping_level_feature=0.0, shopping_le
vel_value=SparseVector(3, {0: 1.0}), occupation_feature=0.0, occupation_value=SparseVec
tor(2, {0: 1.0}), features=SparseVector(39, {0: 108.0, 7: 1.0, 15: 1.0, 18: 1.0, 23: 1.
0, 26: 1.0, 29: 1.0, 30: 1.0, 35: 1.0}))
```

```
In [67]: 1 #=====时间:=====
2 # 创建逻辑回归训练器,并训练模型
3 from pyspark.ml.classification import LogisticRegression
4 lr2 = LogisticRegression()
5 model2 = lr2.setLabelCol('clk').setFeaturesCol('features').fit(train_datasets_2)
6 # models.save(hadoop上)
7 from pyspark.ml.classification import LogisticRegressionModel
8 model2 = LogisticRegressionModel.load('/models/CTRModel_AllOneHot.obj')
9 result_2 = model2.transform(test_datasets_2)
10 result_2.select('clk','price','probability','prediction').sort('probability').show(100)
11 # 对比前面的result_1的预测结果,能发现这里的预测率稍微准确了一点,这里top20里出现了3个
12 # 因此可见对特征的细化处理,已经帮助我们提高模型的效果的
```

clk	price	probability	prediction
0	1.0E8	[0.855244188928558, 0.1447558110714421]	0.0
0	1.0E8	[0.883531437621234, 0.11646856237876606]	0.0
0	1.0E8	[0.8916980898561577, 0.10830191014384229]	0.0
1	5.5555556E7	[0.9251174396034961, 0.07488256039650386]	0.0
0	179.01	[0.9323995173830968, 0.0676004826169032]	0.0
1	159.0	[0.9323995290566156, 0.06760047094338446]	0.0
0	118.0	[0.9323995529753702, 0.06760044702462979]	0.0
0	688.0	[0.9345150616595344, 0.0654849383404656]	0.0
0	339.0	[0.9345152593362689, 0.0654847406637311]	0.0
0	335.0	[0.9345152616019017, 0.06548473839809842]	0.0
0	220.0	[0.9345153267388108, 0.06548467326118919]	0.0
0	176.0	[0.9345153516607427, 0.06548464833925725]	0.0
0	158.0	[0.9345153618560761, 0.065484638143924]	0.0
0	158.0	[0.9345153618560761, 0.065484638143924]	0.0
1	149.0	[0.934515366953742, 0.06548463304625793]	0.0
0	122.5	[0.9345153819635345, 0.06548461803646553]	0.0

```
In [68]: 1 result_2.where('clk=1').select("clk", "price", "probability", "prediction").orderBy("price")
2 # 从该结果也可以看出, result_2的点击率预测率普遍要比result_1高出一点点
```

clk	price	probability	prediction
1	5.5555556E7	[0.9251174396034961, 0.07488256039650386]	0.0
1	159.0	[0.9323995290566156, 0.06760047094338446]	0.0
1	149.0	[0.934515366953742, 0.06548463304625793]	0.0
1	8888.0	[0.9349439274648473, 0.0650560725351527]	0.0
1	138.0	[0.9349441477080421, 0.065055852291958]	0.0
1	35.0	[0.9349442056925659, 0.06505579430743408]	0.0
1	519.0	[0.934948638706219, 0.06505136129378104]	0.0
1	478.0	[0.9349486617859604, 0.06505133821403952]	0.0
1	349.0	[0.9349487344026585, 0.06505126559734156]	0.0
1	348.0	[0.9349487349655783, 0.06505126503442173]	0.0
1	316.0	[0.9349487529790108, 0.06505124702098909]	0.0
1	298.0	[0.9349487631115648, 0.06505123688843525]	0.0
1	298.0	[0.9349487631115648, 0.06505123688843525]	0.0
1	199.0	[0.9349488188405846, 0.06505118115941552]	0.0
1	199.0	[0.9349488188405846, 0.06505118115941552]	0.0
1	198.0	[0.9349488194035036, 0.06505118059649637]	0.0

In []:

1