

```
In [ ]: 1  ### 推荐服务
2  '''
3  - 离线推荐
4  - 先召回对召回结果排序
5  - 为每一个用户都进行召回并排序的过程并且把拍好顺序的结果放到数据库中
6  - 如果需要推荐结果的时候 直接到数据库中按照user_id查询，返回推荐结果
7  - 优点 结构比较简单 推荐服务只需要不断计算，把结果保存到数据库中即可
8  - 缺点 实时性查 如果数据1天不更新 1天之内推荐结果一样的，不能反映用户的实时兴趣
9  - 实时推荐
10 - 排序的模型加载好
11 - 召回阶段的结果缓存
12 - 所有用户的特征缓存
13 - 所有物品的特征缓存
14 - 把推荐的服务暴露出去（django flask）需要推荐结果的服务把 用户id 传递过来
15 - 根据id 找到召回结果
16 - 根据id 找到缓存的用户特征
17 - 根据召回结果的物品id 找到物品的特征
18 - 用户特征+物品特征-》逻辑回归模型 就可以预测点击率
19 - 所有召回的物品的点记率都预测并排序 推荐topN
20 - 实时通过LR模型进行排序的好处
21 - 随时修改召回集
22 - 随时调整用户的特征
23 - 当用户需要推荐服务的时候，获取到最新的召回集和用户特征 得到最新的排序结果 更能
24  '''
```

实时产生推荐结果

CTR预测模型+特征==>预测结果==>TOP-N列表

- CTR预测模型在离线阶段已经训练完成，此处仅需加载
- 特征是：用户实时特征 + （该用户对应的）召回集（离线召回+在线召回）中物品的特征
 - 在线召回：用户刚买了某个种类的物品，就随机取出 该类中 若干个物品，放入召回集中
- 预测结果：计算出点击率，排序，得到，例如top20

```
In [2]: 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 = 'OnlineRecommendation'
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 [ ]: 1 # 以下数据来自第八小节
2 ''' 热编码中:
3 "pvalue_level"特征对应关系:
4 +-----+
5 |pvalue_level|pl_onehot_feature|
6 +-----+
7 |          -1|                0.0|
8 |           3|                3.0|
9 |           1|                2.0|
10 |           2|                1.0|
11 +-----+
12
13 "new_user_class_level" 的特征对应关系:
14 +-----+
15 |new_user_class_level|nucl_onehot_feature|
16 +-----+
17 |                -1|                0.0|
18 |                 3|                2.0|
19 |                 1|                4.0|
20 |                 4|                3.0|
21 |                 2|                1.0|
22 +-----+
23 '''
24 pvalue_level_rela = {-1: 0, 3:3, 1:2, 2:1}
25 new_user_class_level_rela = {-1:0, 3:2, 1:4, 4:3, 2:1}
26 '''
27 "cms_group_id"特征对应关系:
28 +-----+
29 |cms_group_id|min(cms_group_id_feature)|
30 +-----+
31 |           7|                9.0|
32 |          11|                6.0|
33 |           3|                0.0|
34 |           8|                8.0|
35 |           0|               12.0|
36 |           5|                3.0|
37 |           6|               10.0|
38 |           9|                5.0|
39 |           1|                7.0|
40 |          10|                4.0|
41 |           4|                1.0|
42 |          12|               11.0|
43 |           2|                2.0|
44 +-----+
45 '''
46 cms_group_id_rela = {
47     7: 9,
48     11: 6,
49     3: 0,
50     8: 8,
51     0: 12,
52     5: 3,
53     6: 10,
54     9: 5,
55     1: 7,
56     10: 4,

```

```

57     4: 1,
58     12: 11,
59     2: 2
60 }
61 '''
62 "final_gender_code"特征对应关系:
63 +-----+
64 |final_gender_code|min(final_gender_code_feature)|
65 +-----+
66 |           1|           1.0|
67 |           2|           0.0|
68 +-----+
69 '''
70 final_gender_code_rela = {1:1, 2:0}
71 '''
72 "age_level"特征对应关系:
73 +-----+
74 |age_level|min(age_level_feature)|
75 +-----+
76 |           3|           0.0|
77 |           0|           6.0|
78 |           5|           2.0|
79 |           6|           5.0|
80 |           1|           4.0|
81 |           4|           1.0|
82 |           2|           3.0|
83 +-----+
84 '''
85 age_level_rela = {3:0, 0:6, 5:2, 6:5, 1:4, 4:1, 2:3}
86
87 '''
88 "shopping_level"特征对应关系:
89 |shopping_level|min(shopping_level_feature)|
90 +-----+
91 |           3|           0.0|
92 |           1|           2.0|
93 |           2|           1.0|
94 +-----+
95 '''
96 shopping_level_rela = {3:0, 1:2, 2:1}
97 '''
98 "occupation"特征对应关系:
99 +-----+
100 |occupation|min(occupation_feature)|
101 +-----+
102 |           0|           0.0|
103 |           1|           1.0|
104 +-----+
105 '''
106 occupation_rela = {0:0, 1:1}
107
108 pid_rela = {
109     "430548_1007": 0,
110     "430549_1007": 1
111 }
112
113 # key是原始数据, value是StringIndexe编码后的值1, 值1经过onehot编码成为独热编码

```



```
In [38]: 1 ## 特征获取
2 import redis
3 import json
4 import pandas as pd
5 from pyspark.ml.linalg import DenseVector
6
7
8 def create_datasets(userId, pid):
9     client_of_recall = redis.StrictRedis(host="192.168.58.100", port=6379, db=9)
10    client_of_features = redis.StrictRedis(host="192.168.58.100", port=6379, db=10)
11    # 获取用户特征
12    user_feature = json.loads(client_of_features.hget("user_features", userId).decode())
13
14    # 获取用户召回集
15    recall_sets = client_of_recall.smembers(userId)
16
17    result = []
18
19
20    # 遍历召回集
21    for adgroupId in recall_sets:
22        adgroupId = int(adgroupId)
23        # 获取该广告的特征值 price
24        ad_feature = json.loads(client_of_features.hget("ad_features", adgroupId).decode())
25
26        features = {}
27        features.update(user_feature)
28        features.update(ad_feature)
29
30        for k,v in features.items():
31            if v is None:
32                features[k] = -1
33
34        features_col = [
35            # 特征值
36            "price",
37            "cms_group_id",
38            "final_gender_code",
39            "age_level",
40            "shopping_level",
41            "occupation",
42            "pid",
43            "pvalue_level",
44            "new_user_class_level"
45        ]
46
47        "cms_group_id", 类别型特征, 约13个分类 ==> 13维
48        "final_gender_code", 类别型特征, 2个分类 ==> 2维
49        "age_level", 类别型特征, 7个分类 ==> 7维
50        "shopping_level", 类别型特征, 3个分类 ==> 3维
51        "occupation", 类别型特征, 2个分类 ==> 2维
52
53
54        price = float(features["price"])
55
56        pid_value = [0 for i in range(2)]#[0,0]
```

```

57     cms_group_id_value = [0 for i in range(13)]
58     final_gender_code_value = [0 for i in range(2)]
59     age_level_value = [0 for i in range(7)]
60     shopping_level_value = [0 for i in range(3)]
61     occupation_value = [0 for i in range(2)]
62     pvalue_level_value = [0 for i in range(4)]
63     new_user_class_level_value = [0 for i in range(5)]
64
65     pid_value[pid_rela[pid]] = 1
66     cms_group_id_value[cms_group_id_rela[int(features["cms_group_id"])]] = 1
67     final_gender_code_value[final_gender_code_rela[int(features["final_gender_code"])]] = 1
68     age_level_value[age_level_rela[int(features["age_level"])]] = 1
69     shopping_level_value[shopping_level_rela[int(features["shopping_level"])]] = 1
70     occupation_value[occupation_rela[int(features["occupation"])]] = 1
71     pvalue_level_value[pvalue_level_rela[int(features["pvalue_level"])]] = 1
72     new_user_class_level_value[new_user_class_level_rela[int(features["new_user_class_level"])]] = 1
73     # print(pid_value)
74     # print(cms_group_id_value)
75     # print(final_gender_code_value)
76     # print(age_level_value)
77     # print(shopping_level_value)
78     # print(occupation_value)
79     # print(pvalue_level_value)
80     # print(new_user_class_level_value)
81
82     vector = DenseVector([price] + pid_value + cms_group_id_value + final_gender_code_value +
83                          + age_level_value + shopping_level_value + occupation_value + pvalue_level_value + new_user_class_level_value)
84
85     result.append((userId, adgroupId, vector))
86
87     return result
88 # 举例看看用户88 广告资源位"430548_1007" 对应的、召回的500条广告 的特征向量
89 create_datasets(88, "430548_1007")

```

...

```

In [34]: 1 # 加载训练好的逻辑回归模型
          2 from pyspark.ml.classification import LogisticRegressionModel
          3 CTR_model = LogisticRegressionModel.load('/models/CTRModel_AllOneHot.obj')

```

```
In [37]: 1 import pandas as pd
2 pdf = pd.DataFrame(create_datasets(8, '430548_1007'), columns=["userId", "adgroupId", "
3 pdf
```

Out[37]:

	userId	adgroupid	features
0	8	568198	[11.0, 1.0, 0.0, 1.0, 0.0, 0.0, 0.0, 0.0, 0.0, ...
1	8	284442	[1.0, 1.0, 0.0, 1.0, 0.0, 0.0, 0.0, 0.0, 0.0, ...
2	8	40366	[4.0, 1.0, 0.0, 1.0, 0.0, 0.0, 0.0, 0.0, 0.0, ...
3	8	446656	[32.0, 1.0, 0.0, 1.0, 0.0, 0.0, 0.0, 0.0, 0.0, ...
4	8	132372	[5.599999904632568, 1.0, 0.0, 1.0, 0.0, 0.0, 0...
...
495	8	238926	[9.800000190734863, 1.0, 0.0, 1.0, 0.0, 0.0, 0...
496	8	13307	[78.0, 1.0, 0.0, 1.0, 0.0, 0.0, 0.0, 0.0, 0.0, ...
497	8	693981	[23.100000381469727, 1.0, 0.0, 1.0, 0.0, 0.0, ...
498	8	258692	[1.600000023841858, 1.0, 0.0, 1.0, 0.0, 0.0, 0...
499	8	708841	[19.600000381469727, 1.0, 0.0, 1.0, 0.0, 0.0, ...

500 rows × 3 columns


```
In [41]: 1 datasets = spark.createDataFrame(pdf)
        2 datasets.show()
```

userId	adgroupId	features
8	568198	[11.0, 1.0, 0.0, 1.0...
8	284442	[1.0, 1.0, 0.0, 1.0,...
8	40366	[4.0, 1.0, 0.0, 1.0,...
8	446656	[32.0, 1.0, 0.0, 1.0...
8	132372	[5.59999990463256...
8	424510	[12.8000001907348...
8	255632	[0.5, 1.0, 0.0, 1.0,...
8	136509	[8.60000038146972...
8	143566	[29.6000003814697...
8	198254	[2.20000004768371...
8	103023	[28.0, 1.0, 0.0, 1.0...
8	262373	[100.0, 1.0, 0.0, 1...
8	63803	[15.0, 1.0, 0.0, 1.0...
8	12052	[59.0, 1.0, 0.0, 1.0...
8	538461	[20.7999992370605...
8	84047	[9.0, 1.0, 0.0, 1.0,...
8	184721	[10.0, 1.0, 0.0, 1.0...
8	157274	[38.0, 1.0, 0.0, 1.0...
8	77436	[29.0, 1.0, 0.0, 1.0...
8	252225	[12.5, 1.0, 0.0, 1.0...

only showing top 20 rows

```
In [42]: 1 prediction = CTR_model.transform(datasets).sort('probability')
        2 prediction.show()
```

userId	adgroupId	features	rawPrediction	probability	prediction
8	202173	[1888.0, 1.0, 0.0, 1.0, ...]	[2.69017894066573...	[0.93644463234420...	0.0
8	241175	[1800.0, 1.0, 0.0, 1.0, ...]	[2.69017975515559...	[0.93644468081943...	0.0
8	247128	[1350.0, 1.0, 0.0, 1.0, ...]	[2.69018392016059...	[0.93644492870359...	0.0
8	788867	[1220.0, 1.0, 0.0, 1.0, ...]	[2.69018512338425...	[0.93644500031440...	0.0
8	730074	[800.0, 1.0, 0.0, 1.0, ...]	[2.69018901072224...	[0.93644523167188...	0.0
8	397105	[800.0, 1.0, 0.0, 1.0, ...]	[2.69018901072224...	[0.93644523167188...	0.0
8	845130	[520.0, 1.0, 0.0, 1.0, ...]	[2.69019160228090...	[0.93644538590977...	0.0
8	295744	[500.0, 1.0, 0.0, 1.0, ...]	[2.69019178739224...	[0.93644539692675...	0.0
8	296690	[418.0, 1.0, 0.0, 1.0, ...]	[2.69019254634870...	[0.93644544209635...	0.0
8	2267	[350.0, 1.0, 0.0, 1.0, ...]	[2.69019317572724...	[0.93644547955404...	0.0
8	603882	[299.0, 1.0, 0.0, 1.0, ...]	[2.69019364776113...	[0.93644550764729...	0.0
8	627350	[299.0, 1.0, 0.0, 1.0, ...]	[2.69019364776113...	[0.93644550764729...	0.0
8	24364	[278.0, 1.0, 0.0, 1.0, ...]	[2.69019384212803...	[0.93644551921510...	0.0
8	270625	[256.0, 1.0, 0.0, 1.0, ...]	[2.69019404575050...	[0.93644553133375...	0.0
8	176624	[248.0, 1.0, 0.0, 1.0, ...]	[2.69019411979503...	[0.93644553574053...	0.0
8	747336	[245.0, 1.0, 0.0, 1.0, ...]	[2.69019414756173...	[0.93644553739308...	0.0
8	24209	[238.0, 1.0, 0.0, 1.0, ...]	[2.69019421235070...	[0.93644554124901...	0.0
8	289624	[228.0, 1.0, 0.0, 1.0, ...]	[2.69019430490637...	[0.93644554675749...	0.0
8	385931	[195.0, 1.0, 0.0, 1.0, ...]	[2.69019461034007...	[0.93644556493546...	0.0
8	44235	[178.0, 1.0, 0.0, 1.0, ...]	[2.69019476768470...	[0.93644557429987...	0.0

only showing top 20 rows

In [46]:

```
1 # 为 8号用户 推荐的top20物品
2 print(prediction.select('adgroupId').head(20))
3 print([i.adgroupId for i in prediction.select('adgroupId').head(20)])
```

```
[Row(adgroupId=202173), Row(adgroupId=241175), Row(adgroupId=247128), Row(adgroupId=788867), Row(adgroupId=397105), Row(adgroupId=730074), Row(adgroupId=845130), Row(adgroupId=295744), Row(adgroupId=296690), Row(adgroupId=2267), Row(adgroupId=627350), Row(adgroupId=603882), Row(adgroupId=24364), Row(adgroupId=270625), Row(adgroupId=176624), Row(adgroupId=747336), Row(adgroupId=24209), Row(adgroupId=289624), Row(adgroupId=385931), Row(adgroupId=44235)]
[202173, 241175, 247128, 788867, 397105, 730074, 845130, 295744, 296690, 2267, 603882, 627350, 24364, 270625, 176624, 747336, 24209, 289624, 385931, 600455]
```