

12. 实时日志行为处理

根据实时的用户日志行为做出相应的反馈：**实时更新特征、实时更新召回集**

例，用户的浏览、收藏、加购物车、购买等行为记录到日志后，可能包含如下信息：

- 时间
- 地点
- 用户Id
- 商品Id
- 类别Id
- 品牌Id
- 商品价格

以上信息，就目前我们持有的数据而言，能产生实时影响大概只有分为两种类：

- 对用户的基本信息产生影响的数据：**地点(根据当前低点定位来判断用户的消费环境/等级)、购买行为的商品价格**
- 对用户召回结果产生影响的数据：**商品的类别、品牌**

因此现在假设日志格式为："时间,地点,用户ID,商品ID,类别ID,品牌ID,商品价格"

```
In [1]: 1 # spark配置信息
2 from pyspark import SparkConf
3 from pyspark.sql import SparkSession
4
5 SPARK_APP_NAME = "processingOnlineData"
6 SPARK_URL = "yarn"
7
8 conf = SparkConf() # 创建spark config对象
9 config = (
10     ("spark.app.name", SPARK_APP_NAME), # 设置启动的spark的app名称，没有提供，将随
11     ("spark.executor.memory", "2g"), # 设置该app启动时占用的内存用量，默认1g
12     ("spark.executor.cores", "2"), # 设置spark executor使用的CPU核心数
13     ("spark.executor.instances", 1) # 设置spark executor数量，yarn时起作用
14 )
15 # 查看更详细配置及说明: https://spark.apache.org/docs/latest/configuration.html
16 #
17 conf.setAll(config)
18 # 利用config对象，创建spark session
19 spark = SparkSession.builder.config(conf=conf).getOrCreate()
```

```
In [7]: 1 # 注意：初次安装并运行时，由于使用了kafka，所以会自动下载一系列的依赖jar包，会耗费一定
2
3 from pyspark.streaming.kafka import KafkaUtils
4 from pyspark.streaming import StreamingContext
5
6 ssc = StreamingContext(spark.sparkContext, 2)
7
8 kafkaParams = {"metadata.broker.list": "192.168.58.100:9092"}
9 dstream = KafkaUtils.createDirectStream(ssc, ["mytopic"], kafkaParams)
```

In [6]: 1 !hadoop fs -ls /models

Found 3 items

```
drwxr-xr-x - root supergroup          0 2020-12-12 21:11 /models/CTRModel_AllOneHot.obj
drwxr-xr-x - root supergroup          0 2020-12-12 21:12 /models/CTRModel_Normal.obj
drwxr-xr-x - root supergroup          0 2020-12-12 21:12 /models/userCateRatingALSModel.obj
```

```
In [9]: 1 ##### 获取广告和类别的对应关系
2 # 从HDFS中加载广告基本信息数据，返回spark dataframe对象
3 df = spark.read.csv("/data/ad_feature.csv", header=True)
4
5 # 注意：由于本数据集中存在NULL字样的数据，无法直接设置schema，只能先将NULL类型的数据处理掉
6
7 from pyspark.sql.types import StructType, StructField, IntegerType, FloatType
8
9 # 替换掉NULL字符串，替换掉
10 df = df.replace("NULL", "-1")
11
12 # 更改df表结构：更改列类型和列名称
13 ad_feature_df = df.\
14     withColumn("adgroup_id", df.adgroup_id.cast(IntegerType())).withColumnRenamed("adgroup_id", "adgroup_id")\
15     withColumn("cate_id", df.cate_id.cast(IntegerType())).withColumnRenamed("cate_id", "cate_id")\
16     withColumn("campaign_id", df.campaign_id.cast(IntegerType())).withColumnRenamed("campaign_id", "campaign_id")\
17     withColumn("customer", df.customer.cast(IntegerType())).withColumnRenamed("customer", "customer")\
18     withColumn("brand", df.brand.cast(IntegerType())).withColumnRenamed("brand", "brand")\
19     withColumn("price", df.price.cast(FloatType()))
20
21 # 这里我们只需要adgroupId、和cateId
22 _ = ad_feature_df.select("adgroupId", "cateId")
23 # 由于这里数据集其实很少，所以我们再直接转成Pandas dataframe来处理，把数据载入内存
24 pdf = _.toPandas()
25
26
27 # 手动释放一些内存
28 del df
29 del ad_feature_df
30 del _
31 import gc
32 gc.collect()
```

Out[9]: 34

```
In [10]: 1 def m(e):
2         # 当前设定日志数据格式: "时间, 地点, 用户ID, 商品ID, 类别ID, 品牌ID, 商品价格"
3         # 用逗号分割
4         return e[1].split(",")
5
6         import redis
7         import json
8         import numpy as np
9
10        client1 = redis.StrictRedis(host="192.168.58.100", port=6379, db=10)
11        client2 = redis.StrictRedis(host="192.168.58.100", port=6379, db=9)#离线召回集也9库中
12
13        def f(rdd):
14            print("foreach", rdd.collect())
15            for r in rdd.collect():
16                userId = r[2]
17                location_level = r[1]    # 取值范围1-4
18
19                new_user_class_level = location_level if int(location_level) in [1,2,3,4] else
20                data = json.loads(client1.hget("user_features", userId))
21                data["new_user_class_level"] = new_user_class_level    # 注意: 该需求只是假设
22                client1.hset("user_features", userId, json.dumps(data))
23
24                cateId = r[4]
25                # 用户买了哪个类的商品, 就从该类商品中随机抽出50个商品, 在线召回
26                ad_list = pdf.where(pdf.cateId==int(cateId)).dropna().adgroupId.astype(np.int)
27                if ad_list.size > 0:
28                    # 随机抽出当前类别50个广告, 进行在线召回
29                    ret = set(np.random.choice(ad_list, 50))
30                    # 更新到redis 中
31                    client2.sadd(userId, *ret)
```

```
In [11]: 1 dstream.map(m).foreachRDD(f)
```

[illegible]

```
1 ssc.stop()
```

1