

10.离线数据缓存之离线特征

目的:

缓存用户和广告的离线特征, 用于实时推荐

上节已经为每个用户缓存了500个商品, 这些商品是从最感兴趣的类别中随机选的, 使用这些商品的离线特征利用LR算出点击率, 再排序。

```
In [ ]: 1 # spark配置信息
2 from pyspark import SparkConf
3 from pyspark.sql import SparkSession
4
5 SPARK_APP_NAME = "cacheOfflineFeatures"
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.master", SPARK_URL), # spark master的地址
13     ("spark.executor.cores", "2"), # 设置spark executor使用的CPU核心数
14     ("spark.executor.instances", 1) # 设置spark executor数量, yarn时起作用
15 )
16 # 查看更详细配置及说明: https://spark.apache.org/docs/latest/configuration.html
17 #
18 conf.setAll(config)
19
20 # 利用config对象, 创建spark session
21 spark = SparkSession.builder.config(conf=conf).getOrCreate()
```

```
In [64]: 1 # "pid", 广告资源位, 属于场景特征, 也就是说, 每一种广告通常是可以防止在多种资源位下的
2 # 因此这里对于pid, 应该是由广告系统发起推荐请求时, 向推荐系统明确要推荐的用户是谁, 以
3 # 这样如果有多个资源位, 那么每个资源位都会对应相应的一个推荐列表
4
5 # 需要进行缓存的特征值
6
7 feature_cols_from_ad = [
8     "price" # 来自广告基本信息中
9 ]
10
11 # 用户特征
12 feature_cols_from_user = [
13     "cms_group_id",
14     "final_gender_code",
15     "age_level",
16     "shopping_level",
17     "occupation",
18     "pvalue_level",
19     "new_user_class_level"
20 ]
```

从HDFS中加载广告基本信息数据

```
In [66]: 1 _ad_feature_df = spark.read.csv("/data/ad_feature.csv", header=True)
2
3 # 更改表结构, 转换为对应的数据类型
4 from pyspark.sql.types import StructType, StructField, IntegerType, FloatType
5
6 # 替换掉NULL字符串
7 _ad_feature_df = _ad_feature_df.replace("NULL", "-1")
8
9 # 更改df表结构: 更改列类型和列名称
10 ad_feature_df = _ad_feature_df.\
11     withColumn("adgroup_id", _ad_feature_df.adgroup_id.cast(IntegerType())).withColumnRenamed("adgroup_id", "adgroup_id")\
12     withColumn("cate_id", _ad_feature_df.cate_id.cast(IntegerType())).withColumnRenamed("cate_id", "cate_id")\
13     withColumn("campaign_id", _ad_feature_df.campaign_id.cast(IntegerType())).withColumnRenamed("campaign_id", "campaign_id")\
14     withColumn("customer", _ad_feature_df.customer.cast(IntegerType())).withColumnRenamed("customer", "customer")\
15     withColumn("brand", _ad_feature_df.brand.cast(IntegerType())).withColumnRenamed("brand", "brand")\
16     withColumn("price", _ad_feature_df.price.cast(FloatType()))
17
18 def foreachPartition(partition):
19
20     import redis
21     import json
22     client = redis.StrictRedis(host="192.168.58.100", port=6379, db=10)
23
24     for r in partition:
25         data = {
26             "price": r.price
27         }
28         # 转成json字符串再保存, 能保证数据再次倒出来时, 能有效的转换成python类型
29         client.hset("ad_features", r.adgroupId, json.dumps(data))
30
31 ad_feature_df.foreachPartition(foreachPartition)
```

从HDFS加载用户基本信息数据

```
In [67]: 1 from pyspark.sql.types import StructType, StructField, StringType, IntegerType, LongType
2
3 # 构建表结构schema对象
4 schema = StructType([
5     StructField("userId", IntegerType()),
6     StructField("cms_segid", IntegerType()),
7     StructField("cms_group_id", IntegerType()),
8     StructField("final_gender_code", IntegerType()),
9     StructField("age_level", IntegerType()),
10    StructField("pvalue_level", IntegerType()),
11    StructField("shopping_level", IntegerType()),
12    StructField("occupation", IntegerType()),
13    StructField("new_user_class_level", IntegerType())
14 ])
15 # 利用schema从hdfs加载
16 user_profile_df = spark.read.csv("/data/user_profile.csv", header=True, schema=schema)
17 user_profile_df
```

Out[67]: DataFrame[userId: int, cms_segid: int, cms_group_id: int, final_gender_code: int, age_level: int, pvalue_level: int, shopping_level: int, occupation: int, new_user_class_level: int]

```
In [68]: 1 def foreachPartition2(partition):
2
3     import redis
4     import json
5     client = redis.StrictRedis(host="192.168.58.100", port=6379, db=10)
6
7     for r in partition:
8         data = {
9             "cms_group_id": r.cms_group_id,
10            "final_gender_code": r.final_gender_code,
11            "age_level": r.age_level,
12            "shopping_level": r.shopping_level,
13            "occupation": r.occupation,
14            "pvalue_level": r.pvalue_level,
15            "new_user_class_level": r.new_user_class_level
16        }
17        # 转成json字符串再保存, 能保证数据再次倒出来时, 能有效的转换成python类型
18        client.hset("user_features", r.userId, json.dumps(data))
19
20 user_profile_df.foreachPartition(foreachPartition2)
```

In []:

1