

预处理behavior_log数据集

检查数据格式，以及对数据进行简单的统计运算(不同用户不同行为的次数)

1. 建立spark的session

```
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 PYSPARK_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['PYSPARK_PYTHON'] = PYSPARK_PYTHON
7 os.environ['PYSPARK_DRIVER_PYTHON'] = PYSPARK_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 = 'preprocessingBehaviorLog'
14 SPARK_URL = 'spark://192.168.58.100:7077'
15
16 conf = SparkConf() # 创建spark config对象
17 config = (
18     ("spark.app.name", SPARK_APP_NAME), # 设置启动的spark的app名称, 没有提供, 将随
19     ("spark.executor.memory", "2g"), # 设置该app启动时占用的内存用量, 默认1g
20     ("spark.master", SPARK_URL), # spark master的地址
21     ("spark.executor.cores", "2"), # 设置spark executor使用的CPU核心数
22     # 以下三项配置, 可以控制执行器数量
23     # ("spark.dynamicAllocation.enabled", True),
24     # ("spark.dynamicAllocation.initialExecutors", 1), # 1个执行器
25     # ("spark.shuffle.service.enabled", True)
26     # ('spark.sql.pivotMaxValues', '99999'), # 当需要pivot DF, 且值很多时, 需要修改,
27 )
28 # 查看更详细配置及说明: https://spark.apache.org/docs/latest/configuration.html
29
30 conf.setAll(config)
31
32 # 利用config对象, 创建spark session
33 spark = SparkSession.builder.config(conf=conf).getOrCreate()
```

2. 从hdfs上加载要处理的数据

该数据集有20G, 这里仅截取其中一部分

2.1 在hdfs上查看待处理数据

```
In [3]: 1 !hadoop fs -ls /data
        2 # test1.csv是behavior_log数据集截取的一部分
```

Found 5 items

```
-rw-r--r--  1 root supergroup  31286431 2020-11-05 18:06 /data/ad_feature.csv
drwxr-xr-x  - root supergroup           0 2020-11-06 15:38 /data/models
-rw-r--r--  1 root supergroup 1088060964 2020-11-05 18:06 /data/raw_sample.csv
-rw-r--r--  1 root supergroup  33812570 2020-12-12 18:39 /data/test1.csv
-rw-r--r--  1 root supergroup  24056588 2020-11-05 18:06 /data/user_profile.csv
```

2.2 将待处理数据加载到spark, 查看各列数据的类型

```
In [5]: 1 df = spark.read.csv('/data/test1.csv', header=True)
        2 df.show()
```

```
+-----+-----+-----+-----+
|  user|time_stamp|btag| cate| brand|
+-----+-----+-----+-----+
|558157|1493741625|  pv| 6250| 91286|
|558157|1493741626|  pv| 6250| 91286|
|558157|1493741627|  pv| 6250| 91286|
|728690|1493776998|  pv|11800| 62353|
|332634|1493809895|  pv| 1101|365477|
|857237|1493816945|  pv| 1043|110616|
|619381|1493774638|  pv|   385|428950|
|467042|1493772641|  pv| 8237|301299|
|467042|1493772644|  pv| 8237|301299|
|991528|1493780710|  pv| 7270|274795|
|991528|1493780712|  pv| 7270|274795|
|991528|1493780712|  pv| 7270|274795|
|991528|1493780712|  pv| 7270|274795|
|991528|1493780714|  pv| 7270|274795|
|991528|1493780765|  pv| 7270|274795|
|991528|1493780714|  pv| 7270|274795|
|991528|1493780765|  pv| 7270|274795|
|991528|1493780764|  pv| 7270|274795|
|991528|1493780633|  pv| 7270|274795|
|991528|1493780764|  pv| 7270|274795|
+-----+-----+-----+-----+
```

only showing top 20 rows

```
In [6]: 1 df.printSchema()
```

```
root
|-- user: string (nullable = true)
|-- time_stamp: string (nullable = true)
|-- btag: string (nullable = true)
|-- cate: string (nullable = true)
|-- brand: string (nullable = true)
```

2.3 转换成可以处理的数据类型

由上步知各列都是字符串型，将用于数字处理的类型转化成int和long类型

```
In [7]: 1 from pyspark.sql.types import StructType, StructField, StringType, IntegerType, LongType
2 schema = StructType([
3     StructField('userId', IntegerType()),
4     StructField('timestamp', LongType()),
5     StructField('btag', StringType()),
6     StructField('cateId', IntegerType()),
7     StructField('brandId', IntegerType())
8 ])#给每列数据 重命名+变换数据类型，相当于Row()
9 behavior_log_df=spark.read.csv('/data/test1.csv', header=True, schema=schema)
10 behavior_log_df
```

Out[7]: DataFrame[userId: int, timestamp: bigint, btag: string, cateId: int, brandId: int]

```
In [8]: 1 behavior_log_df.show()
```

```
+-----+-----+-----+-----+
|userId|timestamp|btag|cateId|brandId|
+-----+-----+-----+-----+
|558157|1493741625|pv|6250|91286|
|558157|1493741626|pv|6250|91286|
|558157|1493741627|pv|6250|91286|
|728690|1493776998|pv|11800|62353|
|332634|1493809895|pv|1101|365477|
|857237|1493816945|pv|1043|110616|
|619381|1493774638|pv|385|428950|
|467042|1493772641|pv|8237|301299|
|467042|1493772644|pv|8237|301299|
|991528|1493780710|pv|7270|274795|
|991528|1493780712|pv|7270|274795|
|991528|1493780712|pv|7270|274795|
|991528|1493780712|pv|7270|274795|
|991528|1493780714|pv|7270|274795|
|991528|1493780765|pv|7270|274795|
|991528|1493780714|pv|7270|274795|
|991528|1493780765|pv|7270|274795|
|991528|1493780764|pv|7270|274795|
|991528|1493780633|pv|7270|274795|
|991528|1493780764|pv|7270|274795|
+-----+-----+-----+-----+
only showing top 20 rows
```

```
In [9]: 1 # 再次查看数据类型是否转成我们可以处理的类型
2 behavior_log_df.printSchema()
```

```
root
|-- userId: integer (nullable = true)
|-- timestamp: long (nullable = true)
|-- btag: string (nullable = true)
|-- cateId: integer (nullable = true)
|-- brandId: integer (nullable = true)
```

3. spark处理数据

3.1 分析数据集字段的类型和格式

- 查看是否有空值
- 查看每列数据的类型
- 查看每列数据的类别情况

```
In [12]: 1 behavior_log_df.groupby('userId').count().show()
```

```
+-----+-----+
|  userId|count|
+-----+-----+
|  668074|    7|
|  550961|    4|
|  313148|    2|
|  285987|    5|
|  369968|    1|
|  328078|    1|
|  309510|    1|
|  201031|    4|
| 1097471|   15|
| 1137432|    1|
|  511225|    2|
|  396208|   13|
|  913040|    9|
|  604834|    2|
|  535839|    1|
| 147711|   12|
|  304448|    4|
|  815397|    4|
|  387467|    1|
| 1141195|    5|
+-----+-----+
```

only showing top 20 rows

```
In [10]: 1 # 查看总共用多少用户
2 # 第一个count(), 返回的是一个dataframe, 这里的count计算的是每一个分组的个数, 即每个用
3 # 完整的数据集有113w个用户
4 print('查看userId的数据情况: ', behavior_log_df.groupby('userId').count().count())
```

查看userId的数据情况: 242109

```
In [8]: 1 behavior_log_df.groupby('btag').count().show()
2 print('查看btag的数据情况:', behavior_log_df.groupby('btag').count().collect())#collect

+----+-----+
|btag| count|
+----+-----+
| buy| 13668|
| fav| 12807|
| cart| 22142|
| pv| 951383|
+----+-----+
```

查看btag的数据情况: [Row(btag='buy', count=13668), Row(btag='fav', count=12807), Row(btag='cart', count=22142), Row(btag='pv', count=951383)]

```
In [9]: 1 print('查看cateId的数据情况', behavior_log_df.groupby('cateId').count().count())
2 # 商品类别
```

查看cateId的数据情况 6044

```
In [10]: 1 print('查看brandId的数据情况', behavior_log_df.groupby('brandId').count().count())
2 # 商品品牌类别
```

查看brandId的数据情况 49179

```
In [13]: 1 # pandas中选择某列使用df['某列的名字'], 但是sparksql不能这样用, 要使用sql语句, df.select
2 # 可以使用df['某列的名字'].cast(某种数据类型如Longtype())
3 print('判断数据是否有空值:')
4 print('原始数据有%d行'%behavior_log_df.count())
5 #dropna()-某行数据有一个值是空值, 就将该行删除, 注意该方法不去掉str类型的null和NULL
6 print('去掉空值后数据有%d行'%behavior_log_df.dropna().count())
7 # 根据结构可知: 没有空值
```

判断数据是否有空值:
原始数据有1000000行
去掉空值后数据有1000000行

3.2 统计每个用户对各类商品、各类品牌的pv、fav、cart、buy数量, 建立透视表

pyspark.sql.GrouppedData.pivot

pivot透视操作, 把某列里的字段值转换成行并进行聚合运算

如果透视的字段中的不同属性值超过10000个, 则需要设置spark.sql.pivotMaxValues, 否则计算过程中会出现错误。文档介绍 (<https://spark.apache.org/docs/latest/api/python/pyspark.sql.html?highlight=pivot#pyspark.sql.GrouppedData.pivot>)。

```
In [14]: 1 # 统计每个用户对各类商品的pv、fav、cart、buy数量
2 # 建立透视表: | userId|cateId| pv| fav|cart| buy|
3 cate_count_df=behavior_log_df.groupby(behavior_log_df.userId,behavior_log_df.cateId).p
4 # cate_count_df.show()
5 # tmp.where('userId=38456').show()
6 # 条件选择where
```

```
In [15]: 1 cate_count_df.printSchema()
```

```
root
|-- userId: integer (nullable = true)
|-- cateId: integer (nullable = true)
|-- cv: long (nullable = true)
|-- fav: long (nullable = true)
|-- cart: long (nullable = true)
|-- buy: long (nullable = true)
```

```
In [13]: 1 # 统计每个用户对各个品牌的pv、fav、cart、buy数量
2 # 建立透视表 | userId|brandId| cv| fav|cart| buy|
3 brand_count_df=behavior_log_df.groupby(behavior_log_df.userId,behavior_log_df.brandId).p
4 brand_count_df.show()
```

```
+-----+-----+-----+-----+-----+
| userId|brandId|  cv| fav|cart| buy|
+-----+-----+-----+-----+-----+
| 454702|  98931|null|null|null|null|
| 755246| 237649|null|null|null|null|
|  59109| 293023|null|null|null|null|
| 671713| 435218|null|null|null|null|
| 915329| 383166|null|  1|null|null|
| 161038| 247861|null|null|null|null|
|  20612| 184921|null|null|null|null|
|1078294| 113336|null|null|null|null|
| 353660|  39211|null|null|null|null|
| 726597| 102457|null|null|null|null|
| 286584| 186296|null|null|null|null|
|  60577|  93403|null|null|null|null|
|1055041| 393034|null|null|null|null|
|1056336| 352813|null|null|null|null|
| 970730| 231690|null|null|null|null|
| 743396| 310648|null|null|null|null|
| 633672| 139148|null|null|null|null|
| 921779| 168156|null|null|null|null|
| 383183| 185225|null|null|null|null|
| 187100| 238536|null|null|null|null|
+-----+-----+-----+-----+-----+
```

only showing top 20 rows

3.3 将透视表先存储起来

由于运算时间比较长，所以这里先将结果存储在hdfs，供后续其他操作使用 写入数据时才开始计算

```
In [ ]: 1 cate_count_df.write.csv("cate_count.csv", header=True)
        2 brand_count_df.write.csv("brand_count.csv", header=True)
```