

离线埋点日志处理

日志是有固定格式的，使用正则匹配将日志转换为我们处理过的spark sql dataframe

前面已经把点击流日志存储hdfs中，这里再实现对曝光日志的采集，但注意曝光日志只需要发送到hdfs即可

```
/root/bigdata/flume/conf/exposure_log_hdfs.properties :
```

```
# Name the components on this agent
a1.sources = r1
a1.sinks = k1
a1.channels = c1

# Describe/configure the source
a1.sources.r1.type = exec
a1.sources.r1.command = tail -F /root/meiduoSourceCode/logs/exposure.log
a1.sources.r1.channels = c1

a1.sources.r1.interceptors = t1
a1.sources.r1.interceptors.t1.type = timestamp

# Use a channel which buffers events in memory
a1.channels.c1.type = memory
a1.channels.c1.capacity = 1000
a1.channels.c1.transactionCapacity = 100

a1.sinks.k1.type = hdfs
a1.sinks.k1.channel = c1
a1.sinks.k1.hdfs.path = hdfs://localhost:9000/project2-meiduo-rs/logs/exposure/%y-%m-%d
a1.sinks.k1.hdfs.userLocalTimeStamp = true
a1.sinks.k1.hdfs.filePrefix = exposure-
a1.sinks.k1.hdfs.fileType = DataStream
a1.sinks.k1.hdfs.writeFormat = Text
a1.sinks.k1.hdfs.round = true
a1.sinks.k1.hdfs.roundValue = 10
a1.sinks.k1.hdfs.roundUnit = minute
```

启动flume: `flume-ng agent -f /root/bigdata/flume/conf/exposure_log_hdfs.properties -n a1`

```

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 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 = "processingSKUMetadata"
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, 指一
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, 指一台虚拟
22     ("hive.metastore.uris", "thrift://localhost:9083"), # 配置hive元数据的访问, 否
23
24     # 以下三项配置, 可以控制执行器数量
25     ("spark.dynamicAllocation.enabled", True),
26     ("spark.dynamicAllocation.initialExecutors", 1), # 1个执行器
27     ("spark.shuffle.service.enabled", True)
28     ("spark.sql.pivotMaxValues", '99999'), # 当需要pivot DF, 且值很多时, 需要修改, 默
29 )
30 # 查看更详细配置及说明: https://spark.apache.org/docs/latest/configuration.html
31
32 conf.setAll(config)
33
34 # 利用config对象, 创建spark session
35 spark = SparkSession.builder.config(conf=conf).enableHiveSupport().getOrCreate()

```

```

In [2]: 1 !hadoop fs -ls /meiduo_mall/logs/click-trace

Found 2 items
drwxr-xr-x - root supergroup 0 2020-12-23 23:04 /meiduo_mall/logs/click-trace/20-12-23
drwxr-xr-x - root supergroup 0 2020-12-24 11:36 /meiduo_mall/logs/click-trace/20-12-24

```

```
In [7]: 1 date = '20-12-24'
2 click_trace = spark.read.csv('/meiduo_mall/logs/click-trace/%s'%date)
3 click_trace.show(truncate=False)
```

```
|_c0|
|
```

```
|2018/12/01 02:41:57: exposure_timesteamp<1543603102> exposure_loc<detail> timesteamp
<1543603317> behavior<pv> uid<1> sku_id<1> cate_id<1> stay_time<60>|
|2018/12/01 02:43:18: exposure_timesteamp<1543603102> exposure_loc<detail> timesteamp
<1543603398> behavior<pv> uid<1> sku_id<1> cate_id<1> stay_time<60>|
|2018/12/01 02:43:19: exposure_timesteamp<1543603102> exposure_loc<detail> timesteamp
<1543603399> behavior<pv> uid<1> sku_id<1> cate_id<1> stay_time<60>|
|2018/12/01 02:43:20: exposure_timesteamp<1543603102> exposure_loc<detail> timesteamp
<1543603400> behavior<pv> uid<1> sku_id<1> cate_id<1> stay_time<60>|
|2018/12/01 02:43:21: exposure_timesteamp<1543603102> exposure_loc<detail> timesteamp
<1543603401> behavior<pv> uid<1> sku_id<1> cate_id<1> stay_time<60>|
|2018/12/01 13:13:28: exposure_timesteamp<1543641203> exposure_loc<detail> timesteamp
<1543641208> behavior<pv> uid<1> sku_id<1> cate_id<1> stay_time<60>|
|2018/12/01 13:16:01: exposure_timesteamp<1543641203> exposure_loc<detail> timesteamp
<1543641201> behavior<pv> uid<1> sku_id<1> cate_id<1> stay_time<60>|
```

```
In [8]: 1 date = '20-12-24'
2 exposure = spark.read.csv('/meiduo_mall/logs/exposure/%s'%date)
3 exposure.show(truncate=False)
```

```
+-----+
+-----+
|_c0|
+-----+
+-----+
|2018/11/30 03:19:41: exposure_timesteamp<1543519181> exposure_loc<detail> uid<1> sku_i
d<1> cate_id<1>|
|2018/11/30 13:18:13: exposure_timesteamp<1543555093> exposure_loc<detail> uid<1> sku_i
d<1> cate_id<1>|
|2018/11/30 13:18:13: exposure_timesteamp<1543555093> exposure_loc<detail> uid<1> sku_i
d<1> cate_id<1>|
|2018/12/01 01:47:59: exposure_timesteamp<1543600079> exposure_loc<detail> uid<1> sku_i
d<1> cate_id<1>|
|2018/12/01 01:49:20: exposure_timesteamp<1543600079> exposure_loc<detail> uid<1> sku_i
d<1> cate_id<1>|
|2018/12/01 01:52:53: exposure_timesteamp<1543600373> exposure_loc<detail> uid<1> sku_i
d<1> cate_id<1>|
|2018/12/01 01:57:47: exposure_timesteamp<1543600666> exposure_loc<detail> uid<1> sku_i
d<1> cate_id<1>|
|2020/12/23 21:21:21: exposure_timesteamp<1608729679> exposure_loc<detail> uid<1> sku_i
d<1> cate_id<1>|
|2020/12/24 11:09:20: exposure_timesteamp<1608779359> exposure_loc<detail> uid<1> sku_i
d<1> cate_id<1>|
|2020/12/24 11:09:20: exposure_timesteamp<1608779359> exposure_loc<detail> uid<1> sku_i
d<1> cate_id<1>|
|2018/11/30 03:19:41: exposure_timesteamp<1543519181> exposure_loc<detail> uid<1> sku_i
d<1> cate_id<1>|
|2018/11/30 13:18:13: exposure_timesteamp<1543555093> exposure_loc<detail> uid<1> sku_i
d<1> cate_id<1>|
|2018/11/30 13:18:13: exposure_timesteamp<1543555093> exposure_loc<detail> uid<1> sku_i
d<1> cate_id<1>|
|2018/12/01 01:47:59: exposure_timesteamp<1543600079> exposure_loc<detail> uid<1> sku_i
d<1> cate_id<1>|
|2018/12/01 01:49:20: exposure_timesteamp<1543600079> exposure_loc<detail> uid<1> sku_i
d<1> cate_id<1>|
|2018/12/01 01:52:53: exposure_timesteamp<1543600373> exposure_loc<detail> uid<1> sku_i
d<1> cate_id<1>|
|2018/12/01 01:57:47: exposure_timesteamp<1543600666> exposure_loc<detail> uid<1> sku_i
d<1> cate_id<1>|
|2020/12/23 21:21:21: exposure_timesteamp<1608729679> exposure_loc<detail> uid<1> sku_i
d<1> cate_id<1>|
|2020/12/24 11:09:20: exposure_timesteamp<1608779359> exposure_loc<detail> uid<1> sku_i
d<1> cate_id<1>|
|2020/12/24 11:09:20: exposure_timesteamp<1608779359> exposure_loc<detail> uid<1> sku_i
d<1> cate_id<1>|
+-----+
+-----+
```

```
In [39]: 1 # 延伸学习：1. 正则表达式匹配 2. 日志如何转化spark sql dataframe
2 # import re
3 # s = '2018/12/01 02:35:13: exposure_timesteamp<1543601846> exposure_loc<detail> times
4
5 # match = re.search("\
6 # exposure_timesteamp<(P<exposure_timesteamp>.*)> \
7 # exposure_loc<(P<exposure_loc>.*)> \
8 # timesteamp<(P<timesteamp>.*)> \
9 # behavior<(P<behavior>.*)> \
10 # uid<(P<uid>.*)> \
11 # sku_id<(P<sku_id>.*)> \
12 # cate_id<(P<cate_id>.*)> \
13 # stay_time<(P<stay_time>.*)>", s)
14 # match.groupdict()
15
16 # {'exposure_timesteamp': '1543601846',
17 #  'exposure_loc': 'detail',
18 #  'timesteamp': '1543602913',
19 #  'behavior': 'pv',
20 #  'uid': '1',
21 #  'sku_id': '1',
22 #  'cate_id': '1',
23 #  'stay_time': '60'}
```

```

In [12]: 1 # 点击日志的有用信息转换成df
2 import re
3 from pyspark.sql import Row
4
5 def map(row):
6     match = re.search("\
7     exposure_timesteamp<(P<exposure_timesteamp>.*)> \
8     exposure_loc<(P<exposure_loc>.*)> \
9     timesteamp<(P<timesteamp>.*)> \
10    behavior<(P<behavior>.*)> \
11    uid<(P<uid>.*)> \
12    sku_id<(P<sku_id>.*)> \
13    cate_id<(P<cate_id>.*)> \
14    stay_time<(P<stay_time>.*)>", row._c0)
15
16    row = Row(exposure_timesteamp=match.group("exposure_timesteamp"),
17              exposure_loc=match.group("exposure_loc"),
18              timesteamp=match.group("timesteamp"),
19              behavior=match.group("behavior"),
20              uid=match.group("uid"),
21              sku_id=match.group("sku_id"),
22              cate_id=match.group("cate_id"),
23              stay_time=match.group("stay_time"))
24
25    return row
26 click_trace.rdd.map(map).toDF().show()

```

behavior	cate_id	exposure_loc	exposure_timesteamp	sku_id	stay_time	timesteamp	uid
pv	1	detail	1543603102	1	60	1543603317	1
pv	1	detail	1543603102	1	60	1543603398	1
pv	1	detail	1543603102	1	60	1543603399	1
pv	1	detail	1543603102	1	60	1543603400	1
pv	1	detail	1543603102	1	60	1543603401	1
pv	1	detail	1543641203	1	60	1543641208	1
pv	1	detail	1543641203	1	60	1543641361	1
pv	1	detail	1608780903	1	60	1608780903	1
pv	1	detail	1608780903	1	60	1608780903	1
pv	1	detail	1608780903	1	60	1608780903	1
pv	1	detail	1608780903	1	60	1608780903	1
pv	1	detail	1543641203	1	60	1543641363	1
pv	1	detail	1608729679	1	60	1608730353	1
pv	1	detail	1608779359	1	60	1608779362	1
pv	1	detail	1608779359	1	60	1608779362	1
pv	1	detail	1608780822	1	60	1608780822	1
pv	1	detail	1608780838	1	60	1608780838	1
pv	1	detail	1543603102	1	60	1543603316	1
pv	1	detail	1543603102	1	60	1543603316	1
pv	1	detail	1543603102	1	60	1543603316	1

only showing top 20 rows

```
In [11]: 1 # 曝光日志的有用信息转换成df
2 import re
3 from pyspark.sql import Row
4
5 def map(row):
6     match = re.search("\
7 exposure_timesteamp<(P<exposure_timesteamp>.*)> \
8 exposure_loc<(P<exposure_loc>.*)> \
9 uid<(P<uid>.*)> \
10 sku_id<(P<sku_id>.*)> \
11 cate_id<(P<cate_id>.*)>", row._c0)
12     match.group
13
14     row = Row(exposure_timesteamp=match.group("exposure_timesteamp"),
15               exposure_loc=match.group("exposure_loc"),
16               uid=match.group("uid"),
17               sku_id=match.group("sku_id"),
18               cate_id=match.group("cate_id"))
19
20     return row
21 exposure.rdd.map(map).toDF().show()
```

cate_id	exposure_loc	exposure_timesteamp	sku_id	uid
1	detail	1543519181	1	1
1	detail	1543555093	1	1
1	detail	1543555093	1	1
1	detail	1543600079	1	1
1	detail	1543600079	1	1
1	detail	1543600373	1	1
1	detail	1543600666	1	1
1	detail	1608729679	1	1
1	detail	1608779359	1	1
1	detail	1608779359	1	1
1	detail	1543519181	1	1
1	detail	1543555093	1	1
1	detail	1543555093	1	1
1	detail	1543600079	1	1
1	detail	1543600079	1	1
1	detail	1543600373	1	1

曝光日志和点击流日志对齐

利用点击流日志中行为是"pv"的数据同时 cate_id|exposure_loc|exposure_timesteamp|sku_id|uid
——对应的数据，最终就能得出：所有曝光的商品中，哪些商品被用户浏览了，哪些没有被浏览