

# 埋点日志数据处理

## 埋点日志

埋点就是在应用中特定的流程收集一些信息，用来跟踪应用使用的状况，后续用来进一步优化产品或是提供运营的数据支撑。比如用于作为用于实现个性化推荐的数据支撑。

埋点方式主流的无非两种方式：

- 自行研发：在研发的产品中注入代码进行统计，并搭建起相应的后台查询和处理
- 第三方平台：第三方统计工具，如友盟、百度移动等

埋点日志主要分为类型：

- 曝光日志：商品被展示到页面被称为曝光，曝光日志也就是指商品一旦被展示出来，则记录一条曝光日志
  - 曝光时间
  - 曝光场景
  - 用户唯一标识
  - 商品ID
  - 商品类别ID
- 点击流日志：用户浏览、收藏、加购物车、购买、评论、搜索等行为记录日志
  - 被曝光时间：对应曝光日志的曝光时间(浏览)
  - 被曝光场景：对应曝光日志的曝光场景(浏览)
  - 用户唯一标识
  - 行为时间
  - 行为类型
  - 商品ID
  - 商品类别ID
  - 停留时长(浏览)
  - 评分(评论)
  - 搜索词(搜索)

## 埋点日志意义

用户行为偏好分析

- 利用点击流日志分析个体/群体用户的行为特征，预测出用户行为的偏好

统计指标分析

- 点击率：顾名思义被点击的概率，计算公式通常是：点击次数/曝光次数。如某商品共曝光或展示了100次，曝光后总共被点击了10次，那么点击率则是10%
- 跳出率：用户访问一个页面后，之后没有再没有其他操作，称为跳出，计算公式常用的是：访问一次就退出的访问量/总的访问量
  - 整体(整个板块/应用)跳出率
  - 单页面的跳出率。如某页面共计有100个用户访问，但其中有10个用户访问当前页面后就再也没有其他访问了，那么当前页面的跳出率是10%

- 转化率：电商中的转化率计算：商品订单成交量/商品访问量。如某商品累计访问量是100个，最终提交订单的只有5个，那么该商品转化率就是5%

注意：跳出率和转化率中的访问量指的是独立用户访问量。独立用户访问量：首先独立用户访问量并不等价于来访问的总用户个数。比如某用户A在1月1日访问了商品1，1月2日又访问了商品1，那么这里商品1的访问量应算作2次，而不是1次。

```
In [1]: 1 import logging#log: 记录
2
3 def get_logger(logger_name, path, level):
4
5     # 创建logger
6     logger = logging.getLogger(logger_name)
7     # level: OFF、FATAL、ERROR、WARN、INFO、DEBUG、ALL或者自己定义的级别
8     logger.setLevel(level)
9
10    # 创建formatter
11    # %(asctime)s: 打印日志的时间
12    # %(message)s: 打印日志信息
13    fmt = '%(asctime)s: %(message)s'
14    datefmt = '%Y/%m/%d %H:%M:%S'
15    formatter = logging.Formatter(fmt, datefmt)
16
17    # 创建handler
18    # FileHandler: writes formatted logging records to disk files
19    handler = logging.FileHandler(path)
20    handler.setLevel(level)
21
22    # 添加handler和formatter 到 logger
23    handler.setFormatter(formatter)
24    logger.addHandler(handler)
25
26    return logger
```

```
In [14]: 1 import time
2 import logging
3
4 exposure_logger = get_logger('exposure', '/root/workspace/3.rs_project/project2/meiduo/
5                               logging.DEBUG)
6 # 曝光日志
7 exposure_timesteamp = time.time()
8 exposure_loc = 'detail'
9 uid = 1
10 sku_id = 1
11 cate_id = 1
12
```

```
In [15]: 1 exposure_logger.info("exposure_timesteamp<%d> exposure_loc<%s> uid<%d> sku_id<%d> cate
2                               %(exposure_timesteamp, exposure_loc, uid, sku_id, cate_id))
```

```
In [16]: 1 # 每运行上条指令一次，就会在log中记录一条信息
          2 # 显示最后五条信息
          3 !cat /root/workspace/3.rs_project/project2/meiduoSourceCode/logs/exposure.log | tail
```

2020/12/23 21:16:43:exposure\_timesteamp<1608729319> exposure\_loc<detail> uid<1> sku\_id<1> cate\_id<1>  
 2020/12/23 21:21:21:exposure\_timesteamp<1608729679> exposure\_loc<detail> uid<1> sku\_id<1> cate\_id<1>  
 2020/12/23 21:21:21: exposure\_timesteamp<1608729679> exposure\_loc<detail> uid<1> sku\_id<1> cate\_id<1>  
 2020/12/24 11:09:20: exposure\_timesteamp<1608779359> exposure\_loc<detail> uid<1> sku\_id<1> cate\_id<1>  
 2020/12/24 11:09:20: exposure\_timesteamp<1608779359> exposure\_loc<detail> uid<1> sku\_id<1> cate\_id<1>

```
In [17]: 1 import time
          2 import logging
          3 click_trace_logger = get_logger('click_trace', '/root/workspace/3.rs_project/project2/
          4                                     logging.DEBUG)
```

```
In [18]: 1 # 点击流日志
          2 exposure_timesteamp = exposure_timesteamp
          3 exposure_loc = exposure_loc
          4 timesteamp = time.time()
          5 behavior = 'pv' # pv fav cart buy
          6 uid = 1
          7 sku_id = 1
          8 cate_id = 1
          9 stay_time = 60
          10 # 假设某点击流日志记录格式如下:
          11 click_trace_logger.info("exposure_timesteamp<%d> exposure_loc<%s> timesteamp<%d> behav
          12                                     %(exposure_timesteamp, exposure_loc, timesteamp, behavior, uid
          13
```

In [19]: 1 !cat /root/workspace/3.rs\_project/project2/meiduoSourceCode/logs/click\_trace.log

```
2018/11/30 03:20:24: exposure_timesteamp<1543519181> exposure_loc<detail> timesteamp<
1543519224> behavior<pv> uid<1> sku_id<1> cate_id<1> stay_time<60>
2018/11/30 11:54:45: exposure_timesteamp<1543519181> exposure_loc<detail> timesteamp<
1543550085> behavior<pv> uid<1> sku_id<1> cate_id<1> stay_time<60>
2018/11/30 11:54:45: exposure_timesteamp<1543519181> exposure_loc<detail> timesteamp<
1543550085> behavior<pv> uid<1> sku_id<1> cate_id<1> stay_time<60>
2018/11/30 13:17:16: exposure_timesteamp<1543519181> exposure_loc<detail> timesteamp<
1543555036> behavior<pv> uid<1> sku_id<1> cate_id<1> stay_time<60>
2018/11/30 13:17:16: exposure_timesteamp<1543519181> exposure_loc<detail> timesteamp<
1543555036> behavior<pv> uid<1> sku_id<1> cate_id<1> stay_time<60>
2018/11/30 13:17:16: exposure_timesteamp<1543519181> exposure_loc<detail> timesteamp<
1543555036> behavior<pv> uid<1> sku_id<1> cate_id<1> stay_time<60>
2018/11/30 13:18:16: exposure_timesteamp<1543555093> exposure_loc<detail> timesteamp<
1543555096> behavior<pv> uid<1> sku_id<1> cate_id<1> stay_time<60>
2018/11/30 13:18:16: exposure_timesteamp<1543555093> exposure_loc<detail> timesteamp<
1543555096> behavior<pv> uid<1> sku_id<1> cate_id<1> stay_time<60>
2018/11/30 13:18:16: exposure_timesteamp<1543555093> exposure_loc<detail> timesteamp<
1543555096> behavior<pv> uid<1> sku_id<1> cate_id<1> stay_time<60>
2018/11/30 13:18:16: exposure_timesteamp<1543555093> exposure_loc<detail> timesteamp<
1543555096> behavior<pv> uid<1> sku_id<1> cate_id<1> stay_time<60>
```

### flume采集日志

In [20]: 1 !hadoop fs -ls /meiduo\_mall/logs

```
Found 2 items
drwxr-xr-x - root supergroup          0 2020-12-24 08:55 /meiduo_mall/logs/click-trac
e
drwxr-xr-x - root supergroup          0 2020-12-23 22:31 /meiduo_mall/logs/click_trac
e
```

```
In [21]: 1 import re
2 s = '2018/12/01 02:35:13: exposure_timesteamp<1543601846> exposure_loc<detail> timeste
3
4 match = re.search("\
5 exposure_timesteamp<(P<exposure_timesteamp>.*?)> \
6 exposure_loc<(P<exposure_loc>.*?)> \
7 timesteamp<(P<timesteamp>.*?)> \
8 behavior<(P<behavior>.*?)> \
9 uid<(P<uid>.*?)> \
10 sku_id<(P<sku_id>.*?)> \
11 cate_id<(P<cate_id>.*?)> \
12 stay_time<(P<stay_time>.*?)>", s)
13
14 result = []
15 if match:
16     result.append(("exposure_timesteamp", match.group("exposure_timesteamp")))
17     result.append(("exposure_loc", match.group("exposure_loc")))
18     result.append(("timesteamp", match.group("timesteamp")))
19     result.append(("behavior", match.group("behavior")))
20     result.append(("uid", match.group("uid")))
21     result.append(("sku_id", match.group("sku_id")))
22     result.append(("cate_id", match.group("cate_id")))
23     result.append(("stay_time", match.group("stay_time")))
24 result
```

```
Out[21]: [('exposure_timesteamp', '1543601846'),
('exposure_loc', 'detail'),
('timesteamp', '1543602913'),
('behavior', 'pv'),
('uid', '1'),
('sku_id', '1'),
('cate_id', '1'),
('stay_time', '60')]
```

```
In [ ]:
```

```
1
```