

## 20 | 流处理案例实战：分析纽约市出租车载客信息

2019-06-03 蔡元楠 来自北京

《大规模数据处理实战》



你好，我是蔡元楠。

今天我要与你分享的主题是“流处理案例实战：分析纽约市出租车载客信息”。

在上一讲中，我们结合加州房屋信息的真实数据集，构建了一个基本的预测房价的线性回归模型。通过这个实例，我们不仅学习了处理大数据问题的基本流程，而且还进一步熟练了对 RDD 和 DataFrame API 的使用。

你应该已经发现，上一讲的实例是一个典型的批处理问题，因为处理的数据是静态而有边界的。今天让我们来一起通过实例，更加深入地学习用 Spark 去解决实际的流处理问题。

相信你还记得，在前面的章节中我们介绍过 Spark 两个用于流处理的组件——Spark Streaming 和 Structured Streaming。其中 Spark Streaming 是 Spark 2.0 版本前的流处

理库，在 Spark 2.0 之后，集成了 DataFrame/DataSet API 的 Structured Streaming 成为 Spark 流处理的主力。

今天就让我们一起用 Structured Streaming 对纽约市出租车的载客信息进行处理，建立一个实时流处理的 pipeline，实时输出各个区域内乘客小费的平均数来帮助司机决定要去哪里接单。

## 数据集介绍

今天的数据集是纽约市 2009 ~ 2015 年出租车载客的信息。每一次出行包含了两个事件，一个事件代表出发，另一个事件代表到达。每个事件都有 11 个属性，它的 schema 如下所示：

SCHEMA 01			
1	rideId	Long	这次出行的ID
2	taxiId	Long	出租车的ID
3	driverId	Long	出租车司机的ID
4	isStart	Boolean	如果是出发事件，则为true；否则为false
5	startTime	DateTime	这次出行出发的时间
6	endTime	DateTime	这次出行结束的时间。对于出发事件，这个值默认为"1970-01-01 00:00:00"
7	startLon	Float	出发位置的经度
8	startLat	Float	出发位置的纬度
9	endLon	Float	目的地的经度
10	endLat	Float	目的地的纬度
11	passengerCnt	Short	乘客数量

这部分数据有个不太直观的地方，那就是同一次出行会有两个记录，而且代表出发的事件没有任何意义，因为到达事件已经涵盖了所有必要的信息。现实世界中的数据都是这样复杂，不可能像学校的测试数据一样简单直观，所以处理之前，我们要先对数据进行清洗，只留下必要的信息。

这个数据还包含有另外一部分信息，就是所有出租车的付费信息，它有 8 个属性，schema 如下所示。

SCHEMA 02			
1	rideId	Long	这次出行的ID
2	taxiId	Long	出租车的ID
3	driverId	Long	出租车司机的ID
4	startTime	DateTime	这次出行出发的时间
5	paymentType	String	"CSH"（现金） or "CRD"（信用卡）
6	tip	Float	小费金额
7	tolls	Float	过路费金额
8	totalFare	Float	总付费金额

这个数据集可以从[网上](#)下载到，数据集的规模在 100MB 左右，它只是节选了一部分出租车的载客信息，所以在本机运行就可以了。详细的纽约出租车数据集超过了 500GB，同样在[网上](#)可以下载，感兴趣的同学可以下载来实践一下。


## 流数据输入

你可能要问，这个数据同样是静态、有边界的，为什么要用流处理？

因为我们手里没有实时更新的流数据源。我也没有权限去公开世界上任何一个上线产品的数据流。所以，这里只能将有限的数据经过 Kafka 处理，输出为一个伪流数据，作为我们要构建的 pipeline 的输入。


在模块二中，我们曾经初步了解过 Apache Kafka，知道它是基于 Pub/Sub 模式的流数据处理平台。由于我们的专栏并不涉及 Apache Kafka 的具体内容，所以我在这里就不讲如何把这个数据输入到 Kafka 并输出的细节了。你只要知道，在这个例子中，Consumer 是之后要写的 Spark 流处理程序，这个消息队列有两个 Topic，一个包含出行的地理位置信息，一个包含出行的收费信息。Kafka 会**按照时间顺序**，向这两个 Topic 中发布事件，从而模拟一个实时的流数据源。

相信你还记得，写 Spark 程序的第一步就是创建 SparkSession 对象，并根据输入数据创建对应的 RDD 或者 DataFrame。你可以看下面的代码。

 复制代码

```
1 from pyspark.sql import SparkSession
2
3 spark = SparkSession.builder
4     .appName("Spark Structured Streaming for taxi ride info")
5     .getOrCreate()
6
7 rides = spark
8     .readStream
9     .format("kafka")
10    .option("kafka.bootstrap.servers", "localhost:xxxx") //取决于Kafka的配置
11    .option("subscribe", "taxirides")
12    .option("startingOffsets", "latest")
13    .load()
14    .selectExpr("CAST(value AS STRING)")
15
16 fares = spark
17     .readStream
18     .format("kafka")
19     .option("kafka.bootstrap.servers", "localhost:xxxx")
20     .option("subscribe", "taxifares")
21     .option("startingOffsets", "latest")
22     .load()
23     .selectExpr("CAST(value AS STRING)")
```

在这段代码里，我们创建了两个 Streaming DataFrame，并订阅了对应的 Kafka topic，一个代表出行位置信息，另一个代表收费信息。Kafka 对数据没有做任何修改，所以流中的每一个数据都是一个长 String，属性之间是用逗号分割的。

 复制代码


```
1 417986,END,2013-01-02 00:43:52,2013-01-02 00:39:56,-73.984528,40.745377,-73.9759
```

## 数据清洗

现在，我们要开始做数据清洗了。要想分离出我们需要的位置和付费信息，我们首先要将数据分割成一个个属性，并创建对应的 DataFrame 中的列。为此，我们首先要根据数据类型创建




对应的 schema。

 复制代码

```
1 ridesSchema = StructType([
2     StructField("rideId", LongType()), StructField("isStart", StringType()),
3     StructField("endTime", TimestampType()), StructField("startTime", TimestampType()),
4     StructField("startLon", FloatType()), StructField("startLat", FloatType()),
5     StructField("endLon", FloatType()), StructField("endLat", FloatType()),
6     StructField("passengerCnt", ShortType()), StructField("taxiId", LongType()),
7     StructField("driverId", LongType())])
8
9 faresSchema = StructType([
10    StructField("rideId", LongType()), StructField("taxiId", LongType()),
11    StructField("driverId", LongType()), StructField("startTime", TimestampType()),
12    StructField("paymentType", StringType()), StructField("tip", FloatType()),
13    StructField("tolls", FloatType()), StructField("totalFare", FloatType())])
```

接下来，我们将每个数据都用逗号分割，并加入相应的列。

 复制代码


```
1 def parse_data_from_kafka_message(sdf, schema):
2     from pyspark.sql.functions import split
3     assert sdf.isStreaming == True, "DataFrame doesn't receive streaming data"
4     col = split(sdf['value'], ',')
5     for idx, field in enumerate(schema):
6         sdf = sdf.withColumn(field.name, col.getItem(idx).cast(field.dataType))
7     return sdf.select([field.name for field in schema])
8
9 rides = parse_data_from_kafka_message(rides, ridesSchema)
10 fares = parse_data_from_kafka_message(fares, faresSchema)
```

在上面的代码中，我们定义了函数 `parse_data_from_kafka_message`，用来把 Kafka 发来的 message 根据 schema 拆成对应的属性，转换类型，并加入到 DataFrame 的表中。

正如我们之前提到的，读入的数据包含了一些无用信息。

首先，所有代表出发的事件都已被删除，因为到达事件已经包含了出发事件的所有信息，而且只有到达之后才会付费。

其次，出发地点和目的地在纽约范围外的数据，也可以被删除。因为我们的目标是找出纽约市内小费较高的地点。DataFrame 的 filter 函数可以很容易地做到这些。

 复制代码

```
1 MIN_LON, MAX_LON, MIN_LAT, MAX_LAT = -73.7, -74.05, 41.0, 40.5
2 rides = rides.filter(
3     rides["startLon"].between(MIN_LON, MAX_LON) &
4     rides["startLat"].between(MIN_LAT, MAX_LAT) &
5     rides["endLon"].between(MIN_LON, MAX_LON) &
6     rides["endLat"].between(MIN_LAT, MAX_LAT))
7 rides = rides.filter(rides["isStart"] == "END")
```

上面的代码中首先定义了纽约市的经纬度范围，然后把所有起点和终点在这个范围之外的数据都过滤掉了。最后，把所有代表出发事件的数据也移除掉。

当然，除了前面提到的清洗方案，可能还会有别的可以改进的地方，比如把不重要的信息去掉，例如乘客数量、过路费等，你可以自己思考一下。

## Stream-stream Join

我们的目标是找出小费较高的地理区域，而现在收费信息和地理位置信息还在两个 DataFrame 中，无法放在一起分析。那么要用怎样的方式把它们联合起来呢？

你应该还记得，DataFrame 本质上是把数据当成一张关系型的表。在我们这个例子中，rides 所对应的表的键值（Key）是 rideId，其他列里我们关心的就是起点和终点的位置；fares 所对应的表键值也是 rideId，其他列里我们关心的就是小费信息（tips）。

说到这里，你可能会自然而然地想到，如果可以像关系型数据表一样，根据共同的键值 rideId 把两个表 inner join 起来，就可以同时分析这两部分信息了。但是这里的 DataFrame 其实是两个数据流，Spark 可以把两个流 Join 起来吗？

答案是肯定的。在 Spark 2.3 中，流与流的 Join（Stream-stream join）被正式支持。这样的 Join 难点就在于，在任意一个时刻，流数据都不是完整的，流 A 中后面还没到的数据有可

能要和流 B 中已有的数据 Join 起来再输出。为了解决这个问题，我们就要引入**数据水印** (Watermark) 的概念。

数据水印定义了我们可以对数据延迟的最大容忍限度。

比如说，如果定义水印是 10 分钟，数据 A 的事件时间是 1:00，数据 B 的事件时间是 1:10，由于数据传输发生了延迟，我们在 1:15 才收到了 A 和 B，那么我们将只处理数据 B 并更新结果，A 会被无视。在 Join 操作中，好好利用水印，我们就知道什么时候可以不用再考虑旧数据，什么时候必须把旧数据保留在内存中。不然，我们就必须把所有旧数据一直存在内存里，导致数据不断增大，最终可能会内存泄漏。

在这个例子中，为什么我们做这样的 Join 操作需要水印呢？

这是因为两个数据流并不保证会同时收到同一次出行的数据，因为收费系统需要额外的时间去处理，而且这两个数据流是独立的，每个都有可能产生数据延迟。所以要对时间加水印，以免出现内存中数据无限增长的情况。

那么下一个问题就是，究竟要对哪个时间加水印，出发时间还是到达时间？

前面说过了，我们其实只关心到达时间，所以对 rides 而言，我们只需要对到达时间加水印。但是，在 fares 这个 DataFrame 里并没有到达时间的任何信息，所以我们没法选择，只能对出发时间加水印。因此，我们还需要额外定义一个时间间隔的限制，出发时间和到达时间的间隔要在一定的范围内。具体内容你可以看下面的代码。

 复制代码

```
1 faresWithWatermark = fares
2   .selectExpr("rideId AS rideId_fares", "startTime", "totalFare", "tip")
3   .withWatermark("startTime", "30 minutes")
4
5 ridesWithWatermark = rides
6   .selectExpr("rideId", "endTime", "driverId", "taxiId", "startLon", "startLat", "
7   .withWatermark("endTime", "30 minutes")
8
9 joinDF = faresWithWatermark
10   .join(ridesWithWatermark,
11     expr("""
```

```
12     rideId_fares = rideId AND
13     endTime > startTime AND
14     endTime <= startTime + interval 2 hours
15     """)
```

在这段代码中，我们对 fares 和 rides 分别加了半小时的水印，然后把两个 DataFrame 根据 rideId 和时间间隔的限制 Join 起来。这样，joinDF 就同时包含了地理位置和付费信息。


接下来，就让我们开始计算实时的小费最高区域。

## 计算结果并输出

到现在为止，我们还没有处理地点信息。原生的经纬度信息显然并没有很大用处。我们需要做的是把纽约市分割成几个区域，把数据中所有地点的经纬度信息转化成区域信息，这样司机们才可以知道大概哪个地区的乘客比较可能给高点的小费。

纽约市的区域信息以及坐标可以从网上找到，这部分处理比较容易。每个接收到的数据我们都可以判定它在哪个区域内，然后对 joinDF 增加一个列 “area” 来代表终点的区域。现在，让我们假设 area 已经加到现有的 DataFrame 里。接下来我们需要把得到的信息告诉司机了。

还记得第 16 讲和第 17 讲中提到的滑动窗口操作吗？这是流处理中常见的输出形式，即输出每隔一段时间内，特定时间窗口的特征值。在这个例子中，我们可以每隔 10 分钟，输出过去半小时内每个区域内的平均小费。这样的话，司机可以每隔 10 分钟查看一下数据，决定下一步去哪里接单。这个查询（Query）可以由以下代码产生。

 复制代码

```
1 tips = joinDF
2     .groupBy(
3         window("endTime", "30 minutes", "10 minutes"),
4         "area")
5     .agg(avg("tip"))
```

最后，我们把 tips 这个流式 DataFrame 输出。



```
1 query.writeStream
2   .outputMode("append")
3   .format("console")
4   .option("truncate", False)
5   .start()
6   .awaitTermination()
```

你可能会问，为什么我们不可以把输出结果按小费多少进行排序呢？

这是因为两个流的 inner-join 只支持附加输出模式（Append Mode），而现在 Structured Streaming 不支持在附加模式下进行排序操作。希望将来 Structured Streaming 可以提供这个功能，但是现在，司机们只能扫一眼所有的输出数据来大概判断哪个地方的小费最高了。

## 小结

流处理和批处理都是非常常见的应用场景，而且相较而言流处理更加复杂，对延迟性要求更高。今天我们再次通过一个实例帮助你了解要如何利用 Structured Streaming 对真实数据集进行流处理。Spark 最大的好处之一就是它拥有统一的批流处理框架和 API，希望你在课下要进一步加深对 DataSet/DataFrame 的熟练程度。

## 思考题

今天的主题是“案例实战”，不过我留的是思考题，而不是实践题。因为我不确定你是否会使用 Kafka。如果你的工作中会接触到流数据，那么你可以参考今天这个案例的思路和步骤来解决问题，多加练习以便熟悉 Spark 的使用。如果你还没有接触过流数据，但却想往这方面发展的话，我就真的建议你去学习一下 Kafka，这是个能帮助我们更好地做流处理应用开发和部署的利器。

现在，来说一下今天的思考题吧。

1. 为什么流的 Inner-Join 不支持完全输出模式？
2. 对于 Inner-Join 而言，加水印是否是必须的？Outer-Join 呢？

欢迎你把答案写在留言区，与我和其他同学一起讨论。

如果你觉得有所收获，也欢迎把文章分享给你的朋友。

© 版权归极客邦科技所有，未经许可不得传播售卖。页面已增加防盗追踪，如有侵权极客邦将依法追究其法律责任。

## 精选留言 (13)



**never leave**

2019-06-03

官网上说inner join的watermark是可选的，outer join的watermark是必选的。但是我感觉应该都是必选的吧，就像案例中的inner join一样，如果不是必须的话，旧数据一直保存在内存中，有可能导致内存不够。

作者回复: never leave同学，感谢提问。现阶段不仅Inner-join不支持完全输出模式，任何类型的Join都不支持完全输出模式。因为完全输出模式要求每当有新数据输入时，输出完整的结果表。而对于无边界数据，我们很难把所有历史数据存在内存中。所以，一般Join的都是在某个时间窗口内的流数据，这就是引入watermarking的原因。希望将来Spark可以引入新的机制来支持这一点。

Outer join是要在Inner Join的基础上，把没有匹配的行的对应列设为NULL。但是由于流数据的无边界性，Spark永远无法知道在未来会不会找到匹配的数据。所以为了保证Outer Join的正确性，加水印是必须的。这样Spark的执行引擎只要在水印的有效期内没找到与之匹配的数据，就可以把对应的列设为NULL并输出。

由于Inner Join不需要连接两个表中所有的行，所以在Spark官网的叙述中，水印和事件时间的限制不是必须的。但是如果不加任何限制，流数据会不断被读入内存，这样不安全的。所以，即便是Inner Join，我也推荐你加水印和事件时间的限制。



👍 23



**FengX**

2019-06-03

老师，请问join操作里有riderId了，为什么要加上endTime > startTime AND endTime <= startTime + interval 2 hours?

作者回复: Feng.X, 感谢提问。

这个限制目的在于抛弃任何长于2个小时的出租车出行数据。

对于这个例子来说, 这样一个对事件时间的限制确实不是必须的。加入它其实是为了告诉你, 在基于事件时间来join两个流时, 我们一般不考虑时间跨度过大的情况, 因为它们没有普遍意义, 还会影响数据分析的结果。

举例, 对于一个网页广告, 我们需要知道用户看到一个广告后要多长时间才会去点击它, 从而评估广告的效果。这里显然有两个流: 一个代表用户看到广告的事件, 另一个代表用户点击广告的事件。尽管我们可以通过用户的ID来Join这两个流, 但是我们需要加一个限制, 就是点击广告的时间不能比看到广告的时间晚太久, 否则Join的结果很可能是不准确的。比如, 用户可能在1:00和2:00都看到了广告, 但是只在2:01点击了它, 我们应该把2:00和2:01Join起来, 而不应该Join1:00和2:01, 因为1:00看到的广告并没有促使他点击。

共 2 条评论 >

👍 8



**Poleness**

2019-06-04

请问下, 这里解析kafka的value的时候, 自定义了schema, 但真正生产中很多数据的类型结构是很复杂的, 徒手写schema的方式不一定可行。不知道有没有更优雅的方式?

(看了源码, 如果是json等格式好像可以自动推断, 但是对于kafka, 他的sourceSchema好像是写死的, 不知大家有没有好的建议或者经验?)

共 1 条评论 >

👍 5



**谢志斌**

2020-04-23

老师好, 纽约市出租车第一个数据集链接, 无法访问。



👍 2



**jon**

2019-06-03

不支持完全输出是因为join的只是一个时间窗口内的数据

在这个例子中inner join使用watermark 是必须的, left join watermark不是必须的

作者回复: jon, 感谢提问。

现阶段不仅Inner-join不支持完全输出模式，任何类型的Join都不支持完全输出模式。因为完全输出模式要求每当有新数据输入时，输出完整的结果表。而对于无边界数据，我们很难把所有历史数据存在内存中。所以，一般Join的都是在某个时间窗口内的流数据，这就是引入watermarking的原因。希望将来Spark可以引入新的机制来支持这一点。

Outer join是要在Inner Join的基础上，把没有匹配的行的对应列设为NULL。但是由于流数据的无边界性，Spark永远无法知道在未来会不会找到匹配的数据。所以为了保证Outer Join的正确性，加水印是必须的。这样Spark的执行引擎只要在水印的有效期内没找到与之匹配的数据，就可以把对应的列设为NULL并输出。

由于Inner Join不需要连接两个表中所有的行，所以在Spark官网的叙述中，水印和事件时间的限制不是必须的。但是如果不加任何限制，流数据会不断被读入内存，这样不安全的。所以，即便是Inner Join，我也推荐你加水印和事件时间的限制。



**lhk**

2019-09-17

老师你好，请教个watermark的问题：水印是为了解决数据出现延迟时，流处理程序要等待多久。那超过这个时间的数据就丢弃了吗？程序不会再处理他们了吗？比如水印设置30分钟，那31分钟到来的数据就不管了是吧？



**刘万里**

2019-06-03

老师 您好，最近好久没用spark，有个问题请教一下，现在最新spark是否已经支持cep了



**YX**

2021-10-19

比如说，如果定义水印是 10 分钟，数据 A 的事件时间是 1:00，数据 B 的事件时间是 1:10，由于数据传输发生了延迟，我们在 1:15 才收到了 A 和 B，那么我们将只处理数据 B 并更新结果，A 会被无视。

-----  
这里对水印的表述存在一定的不准确，应该是和具体收到的时间无关，而是「max event time seen by the engine」系统当前最大的event time。





天敌

2021-03-24

老师，数据集下载不了了，能再分享一下吗？



之渊

2020-08-22

java版本demo: 模拟的数据集。 <https://gitee.com/oumin12345/daimademojihe/tree/master/cloudx/bigdata/src/main/java/test/spark/streaming>

还是花了不少时间的。对于初学者来说还是值得写点demo的

共 2 条评论 >



北冥有鱼

2020-05-12

老师，比如A和B表join，且A和B都是实时数据，A需要用到B表的历史全量数据，即通过A，保证能取到B中数据，要怎么处理呢？



都市夜归人

2019-06-12

这部分数据有个不太直观的地方，那就是同一次出行会有两个记录，...  
为何会出现两个记录？用一条记录也能记录出发和到达吧？



Ming

2019-06-03

我猜：

对于inner join来说，用不用watermark只是纯粹的一个性能考量，不影响单条数据的正确性，只影响最终分析的样本大小。

对于outer join来说，用watermark会影响单条数据正确性，所以在逻辑上看应该是不推荐的，除非会有内存泄漏的风险。

我倒是好奇为啥spark把这个特性叫水印

作者回复: Ming，感谢提问。



现阶段不仅Inner-join不支持完全输出模式，任何类型的Join都不支持完全输出模式。因为完全输出模式要求每当有新数据输入时，输出完整的结果表。而对于无边界数据，我们很难把所有历史数据存在内存中。所以，一般Join的都是在某个时间窗口内的流数据，这就是引入watermarking的原因。希望将来Spark可以引入新的机制来支持这一点。

Outer join是要在Inner Join的基础上，把没有匹配的行的对应列设为NULL。但是由于流数据的无边界性，Spark永远无法知道在未来会不会找到匹配的数据。所以为了保证Outer Join的正确性，加水印是必须的。这样Spark的执行引擎只要在水印的有效期内没找到与之匹配的数据，就可以把对应的列设为NULL并输出。

由于Inner Join不需要连接两个表中所有的行，所以在Spark官网的叙述中，水印和事件时间的限制不是必须的。但是如果不加任何限制，流数据会不断被读入内存，这样不安全的。所以，即便是Inner Join，我也推荐你加水印和事件时间的限制。

共 2 条评论 >

