

Documentation of Python script: <https://github.com/adaltas/ece-spark/tree/master/structured-streaming> (<https://github.com/adaltas/ece-spark/tree/master/structured-streaming>)

Spark Streaming slides link: <http://bit.ly/2DjKjxL> (<http://bit.ly/2DjKjxL>)

Took 0 sec. Last updated by gauthier at March 18 2020, 1:25:18 PM.

FINISHED

1. SSH to the edge server
2. Get the code from HDFS

```
hdfs dfs -get /learning/code/spark .
```

3. Cd to the code folder: `cd spark/structured-streaming`
4. Run the Python script (choose an unused port number):

```
python3 stream_taxi_data.py edge-1.au.adaltas.cloud 11111 fares
```

Took 0 sec. Last updated by gauthier at March 18 2020, 4:23:30 PM.

```
%pyspark
# Create a socket readStream
fares_raw = spark \
  .readStream \
  .format("socket") \
  .option("host", "edge-1.au.adaltas.cloud") \
  .option("port", 11333) \
  .load()
```

FINISHED

Took 0 sec. Last updated by gauthier at March 18 2020, 2:03:17 PM.

```
%pyspark
from pyspark.sql.functions import explode
from pyspark.sql.functions import split
from pyspark.sql.functions import window

# Parse the socket message "manually"
fares = fares_raw.select(
    split(fares_raw.value, ',')[0].alias('ride_id').cast('int'),
    split(fares_raw.value, ',')[1].alias('taxi_id').cast('int'),
    split(fares_raw.value, ',')[2].alias('driver_id').cast('int'),
    split(fares_raw.value, ',')[3].alias('start_time').cast('timestamp'),
    split(fares_raw.value, ',')[4].alias('payment_type'),
    split(fares_raw.value, ',')[5].alias('tip').cast('float'),
    split(fares_raw.value, ',')[6].alias('tolls').cast('float'),
    split(fares_raw.value, ',')[7].alias('total_fare').cast('float')
)
```

FINISHED

Took 0 sec. Last updated by gauthier at March 18 2020, 2:03:19 PM.

```
%pyspark
# Write all fares events to an in-memory table named "fares"
fares_query = fares \
    .writeStream \
    .outputMode("append") \
    .format("memory") \
    .queryName("fares") \
    .start()
```

FINISHED

Took 0 sec. Last updated by gauthier at March 18 2020, 1:59:51 PM.

```
%pyspark
# Pretty print the result table
z.show(spark.table("fares"))
```

FINISHED

settings ▾

| ride_id ▾ | taxi_id ▾ | driver_id ▾ | start_time ▾ | payment_type ▾ | tip | ≡ |
|-----------|------------|-------------|-----------------------|----------------|------|---|
| 1 | 2013000001 | 2013000001 | 2020-03-18 12:43:39.0 | CSH | 0 | |
| 2 | 2013000002 | 2013000002 | 2020-03-18 12:43:39.0 | CSH | 0 | |
| 3 | 2013000003 | 2013000003 | 2020-03-18 12:43:39.0 | CRD | 2.2 | |
| 4 | 2013000004 | 2013000004 | 2020-03-18 12:43:39.0 | CRD | 1.7 | |
| 5 | 2013000005 | 2013000005 | 2020-03-18 12:43:39.0 | CRD | 4.65 | |
| 6 | 2013000006 | 2013000006 | 2020-03-18 12:43:39.0 | CSH | 0 | |

dsti-a19/ref/lab3

x

FINISHED

Took 0 sec. Last updated by gauthier at March 18 2020, 2:01:12 PM.

FINISHED

FINISHED

FINISHED

SPARK JOB (<http://wrk-1.au.adaltas.cloud:39901/jobs/job?id=86>) FINISHED



 settings 

FINISHED

Took 1 sec. Last updated by gauthier at March 18 2020, 2:26:20 PM.

FINISHED

FINISHED

Took 0 sec. Last updated by gauthier at March 18 2020, 2:34:30 PM.

READY

READY

FINISHED

FINISHED

Took 0 sec. Last updated by gauthier at March 18 2020, 2:38:26 PM.

FINISHED

| ride_id | is_start | start_time | end_time | start_lon | start_lat | end_lon | |
|---------|----------|-----------------------|-----------------------|------------|-----------|------------|--|
| 65 | END | 2020-03-18 13:37:48.0 | 2020-03-18 13:37:50.0 | -73.99221 | 40.725124 | -73.991646 | |
| 137 | END | 2020-03-18 13:38:39.0 | 2020-03-18 13:38:41.0 | 0 | 0 | 0 | |
| 77 | END | 2020-03-18 13:37:59.0 | 2020-03-18 13:38:01.0 | -73.9701 | 40.768005 | -73.96977 | |
| 94 | END | 2020-03-18 13:38:19.0 | 2020-03-18 13:38:21.0 | -74.005165 | 40.72053 | -74.00393 | |

| | | | | | | |
|-------------|-----|-----------------------|-----------------------|-----------|-----------|-----------|
| 70 | END | 2020-03-18 13:57:55.0 | 2020-03-18 13:57:55.0 | -73.97544 | 40.749657 | -73.97733 |
| <div></div> | | | | | | |

dsti-a19/ref/lab3

Took 0 sec. Last updated by gauthier at March 18 2020, 2:40:37 PM.

```
%pyspark
full_rides_query.stop()
```

FINISHED

Took 1 sec. Last updated by gauthier at March 18 2020, 2:41:18 PM.

```
%pyspark
# Expected lab code for 1 KPI
from pyspark.sql.functions import explode
from pyspark.sql.functions import split
from pyspark.sql.functions import window

# Create a socket readStream
fares_raw = spark \
    .readStream \
    .format("socket") \
    .option("host", "edge-1.au.adaltas.cloud") \
    .option("port", 11333) \
    .load()

# Parse the socket message "manually"
fares = fares_raw.select(
    split(fares_raw.value, ',')[0].alias('ride_id').cast('int'),
    split(fares_raw.value, ',')[1].alias('taxi_id').cast('int'),
    split(fares_raw.value, ',')[2].alias('driver_id').cast('int'),
    split(fares_raw.value, ',')[3].alias('start_time').cast('timestamp'),
    split(fares_raw.value, ',')[4].alias('payment_type'),
    split(fares_raw.value, ',')[5].alias('tip').cast('float'),
    split(fares_raw.value, ',')[6].alias('tolls').cast('float'),
    split(fares_raw.value, ',')[7].alias('total_fare').cast('float')
)

fares_count = fares \
    .withWatermark('start_time', '5 minutes') \
    .groupBy(window(fares.start_time, '2 minutes', '2 minutes'), fares.payment_type) \
    .agg({'ride_id': 'count', 'total_fare': 'mean', 'tip': 'mean'})

# Start writting the stream to an in-memory table
fares_count_query = fares_count \
    .writeStream \
    .trigger(processingTime='30 seconds') \
    .outputMode("complete") \
    .format("memory") \
    .queryName("fares_count") \
    .start()
```

FINISHED

Took 0 sec. Last updated by gauthier at March 18 2020, 2:58:05 PM.



Took 0 sec. Last updated by gauthier at March 18 2020, 3:02:43 PM.

```
%pyspark
fares_count_query.stop()
```

FINISHED

Took 0 sec. Last updated by gauthier at March 18 2020, 3:02:38 PM.

FINISHED

Homework

Display using a graph:

- 1. The percentage for each type of payment per window
- 2. The mean ride duration

Took 0 sec. Last updated by ggunther at March 18 2020, 2:52:56 PM.

dsti-a19/ref/lab3

%md

READY