

# Assignment 9.1

By Kurt Stoneburner

<https://sparkbyexamples.com/pyspark-tutorial/> (<https://sparkbyexamples.com/pyspark-tutorial/>)

watermarks are discussed here: <https://spark.apache.org/docs/latest/structured-streaming-programming-guide.html> (<https://spark.apache.org/docs/latest/structured-streaming-programming-guide.html>)

<https://sparkbyexamples.com/pyspark/pyspark-exception-java-gateway-process-exited-before-sending-the-driver-its-port-number/> (<https://sparkbyexamples.com/pyspark/pyspark-exception-java-gateway-process-exited-before-sending-the-driver-its-port-number/>)

<https://quabr.com/58723314/pyspark-failed-to-find-data-source-kafka> (<https://quabr.com/58723314/pyspark-failed-to-find-data-source-kafka>)

[https://search.maven.org/search?q=g:org.apache.spark%20AND%20a:spark-streaming-kafka-0-8-assembly\\_2.11](https://search.maven.org/search?q=g:org.apache.spark%20AND%20a:spark-streaming-kafka-0-8-assembly_2.11) ([https://search.maven.org/search?q=g:org.apache.spark%20AND%20a:spark-streaming-kafka-0-8-assembly\\_2.11](https://search.maven.org/search?q=g:org.apache.spark%20AND%20a:spark-streaming-kafka-0-8-assembly_2.11))

<https://www.rittmanmead.com/blog/2017/01/getting-started-with-spark-streaming-with-python-and-kafka/> (<https://www.rittmanmead.com/blog/2017/01/getting-started-with-spark-streaming-with-python-and-kafka/>)

<https://www.analyticsvidhya.com/blog/2021/06/setting-up-real-time-structured-streaming-with-spark-and-kafka-on-windows-os/> (<https://www.analyticsvidhya.com/blog/2021/06/setting-up-real-time-structured-streaming-with-spark-and-kafka-on-windows-os/>)

<https://spark.apache.org/docs/latest/api/python/reference/api/pyspark.SparkConf.html> (<https://spark.apache.org/docs/latest/api/python/reference/api/pyspark.SparkConf.html>)

[https://stackoverflow.com/questions/54227744/pyspark-2-x-programmatically-adding-maven-jar-coordinates-to-spark?\\_\\_cpo=aHR0cHM6Ly9zdGFja292ZXJmbG93LmNvbQ](https://stackoverflow.com/questions/54227744/pyspark-2-x-programmatically-adding-maven-jar-coordinates-to-spark?__cpo=aHR0cHM6Ly9zdGFja292ZXJmbG93LmNvbQ) ([https://stackoverflow.com/questions/54227744/pyspark-2-x-programmatically-adding-maven-jar-coordinates-to-spark?\\_\\_cpo=aHR0cHM6Ly9zdGFja292ZXJmbG93LmNvbQ](https://stackoverflow.com/questions/54227744/pyspark-2-x-programmatically-adding-maven-jar-coordinates-to-spark?__cpo=aHR0cHM6Ly9zdGFja292ZXJmbG93LmNvbQ))

<https://duckduckgo.com/?t=ffab&q=SparkSession+readstream+Failed+to+find+data+source%3A+kafka&ia=web> (<https://duckduckgo.com/?t=ffab&q=SparkSession+readstream+Failed+to+find+data+source%3A+kafka&ia=web>)

reintsall pyspark: <https://sparkbyexamples.com/pyspark-tutorial/#pyspark-installation> (<https://sparkbyexamples.com/pyspark-tutorial/#pyspark-installation>)

conda install -c conda-forge pyspark

<https://kontext.tech/column/spark/298/get-the-current-spark-context-settingsconfigurations>

(<https://kontext.tech/column/spark/298/get-the-current-spark-context-settingsconfigurations>)

Find Spark add\_packages resolved the local Spark Issue. [https://github.com/minrk/findspark](https://github.com/minrk/findspark/pull/11) /pull/11 (<https://github.com/minrk/findspark/pull/11>)

In [ ]:

In [1]:

```
1 import os
2 import time
3 import shutil
4 import json
5 from pathlib import Path
6
7 import pandas as pd
8
9 from kafka import KafkaProducer, KafkaAdminClient
10 from kafka.admin.new_topic import NewTopic
11 from kafka.errors import TopicAlreadyExistsError
12
13 from pyspark.sql import SparkSession
14 from pyspark.streaming import StreamingContext
15 from pyspark import SparkConf
16 from pyspark.sql.functions import window, from_json, col
17 from pyspark.sql.types import StringType, TimestampType, DoubleType,
18 from pyspark.sql.functions import udf
19 from IPython.display import clear_output
20 import findspark
21 findspark.init()
22
23
24 current_dir = Path(os.getcwd()).absolute()
25 checkpoint_dir = current_dir.joinpath('checkpoints')
26 locations_checkpoint_dir = checkpoint_dir.joinpath('locations')
27 accelerations_checkpoint_dir = checkpoint_dir.joinpath('accelerations')
28
29 if locations_checkpoint_dir.exists():
30     shutil.rmtree(locations_checkpoint_dir)
31
32 if accelerations_checkpoint_dir.exists():
33     shutil.rmtree(accelerations_checkpoint_dir)
34
35 locations_checkpoint_dir.mkdir(parents=True, exist_ok=True)
```

## Configuration Parameters

**TODO:** Change the configuration parameters to the appropriate values for your setup.

In [2]:

```
1 config = dict(
2     bootstrap_servers=['kafka.kafka.svc.cluster.local:9092'],
3     first_name='Kurt',
```

```

4     last_name='Stoneburner'
5 )
6
7 config['client_id'] = '{}{}'.format(
8     config['last_name'],
9     config['first_name']
10 )
11 config['topic_prefix'] = '{}{}'.format(
12     config['last_name'],
13     config['first_name']
14 )
15
16 config['locations_topic'] = '{}-locations'.format(config['topic_prefix'])
17 config['accelerations_topic'] = '{}-accelerations'.format(config['topic_prefix'])
18 config['simple_topic'] = '{}-simple'.format(config['topic_prefix'])
19

```

```

Out[2]: {'bootstrap_servers': ['kafka.kafka.svc.cluster.local:9092'],
'first_name': 'Kurt',
'last_name': 'Stoneburner',
'client_id': 'StoneburnerKurt',
'topic_prefix': 'StoneburnerKurt',
'locations_topic': 'StoneburnerKurt-locations',
'accelerations_topic': 'StoneburnerKurt-accelerations',
'simple_topic': 'StoneburnerKurt-simple'}

```

## Create Topic Utility Function

The `create_kafka_topic` helps create a Kafka topic based on your configuration settings.

For instance, if your first name is *John* and your last name is *Doe*,

`create_kafka_topic('locations')` will create a topic with the name `DoeJohn-locations`. The function will not create the topic if it already exists.

```

In [3]: 1 def create_kafka_topic(topic_name, config=config, num_partitions=1, replication_factor=1):
2     bootstrap_servers = config['bootstrap_servers']
3     client_id = config['client_id']
4     topic_prefix = config['topic_prefix']
5     name = '{}-{}'.format(topic_prefix, topic_name)
6
7     admin_client = KafkaAdminClient(
8         bootstrap_servers=bootstrap_servers,
9         client_id=client_id
10    )
11
12    topic = NewTopic(
13        name=name,
14        num_partitions=num_partitions,
15        replication_factor=replication_factor
16    )
17
18    topic_list = [topic]
19    try:
20        admin_client.create_topics(new_topics=topic_list)
21        print('Created topic "{}"'.format(name))

```

```

22     except TopicAlreadyExistsError as e:
23         print('Topic "{}" already exists'.format(name))
24 for topic in ['locations', 'accelerations']:
Topic "StoneburnerKurt-locations" already exists
Topic "StoneburnerKurt-accelerations" already exists

```

```

In [4]: 1
2  #!/*** Close Spark if already running. Guarantees Spark is loaded with
3  #!/*** Prevents some Ipython/Notebook related issues.
4
5  spark = SparkSession.builder\
6          .appName("Assignment09")\
7          .getOrCreate()
8
9  df_locations = spark\
10     .readStream.format("kafka")\
11     .option("kafka.bootstrap.servers", config['bootstrap_servers'][0])\
12     .option("subscribe", config['locations_topic'])\
13     .load()
14
15

```

**TODO:** Create a data frame called `df_accelerations` that reads from the accelerations topic you published to in assignment 8. In order to read data from this topic, make sure that you are running the notebook you created in assignment 8 that publishes acceleration and location data to the `LastNameFirstname-simple` topic.

```

In [5]: 1 df_accelerations = spark\
2         .readStream.format("kafka")\
3         .option("kafka.bootstrap.servers", config['bootstrap_servers'][0])\
4         .option("subscribe", config['accelerations_topic'])\
5         .load()
6

```

**TODO:** Create two streaming queries, `ds_locations` and `ds_accelerations` that publish to the `LastNameFirstname-simple` topic. See <http://spark.apache.org/docs/latest/structured-streaming-programming-guide.html#starting-streaming-queries> (<http://spark.apache.org/docs/latest/structured-streaming-programming-guide.html#starting-streaming-queries>) and <http://spark.apache.org/docs/latest/structured-streaming-kafka-integration.html> (<http://spark.apache.org/docs/latest/structured-streaming-kafka-integration.html>) for more information.

```

In [ ]: 1 ds_locations = df_locations.writeStream \
2         .format("kafka") \
3         .option("checkpointLocation", locations_checkpoint_dir) \
4         .option("kafka.bootstrap.servers", config['bootstrap_servers'][0]) \
5         .option("topic", config['simple_topic']) \
6         .start()
7
8
9  ds_accelerations = df_locations.writeStream \
10     .format("kafka") \
11     .option("checkpointLocation", locations_checkpoint_dir) \

```

```
12 .option("kafka.bootstrap.servers", config['bootstrap_servers'][0]) \
13 .option("topic", config['simple_topic']) \
14 .start()
15 try:
16     ds_locations.awaitTermination()
17     ds_accelerations.awaitTermination()
18 except KeyboardInterrupt:
```

In [ ]:

```
19 print("DONE!")
```