

Stream the data stored on the GCS bucket into Kafka by breaking the data into batches of 10 records that are written to Kafka separated by a sleep time of 10 seconds until 100 records are written. Use Spark Streaming to read from Kafka every 5 seconds and emit the count of rows seen in the last 10 seconds.

```
me19b190@me19b190-ssh8:~$ sudo apt update && sudo apt install -y tmux vim openjdk-8-jdk wget curl git python3-pip
Hit:1 http://us-central1.gce.archive.ubuntu.com/ubuntu focal InRelease
Hit:2 http://us-central1.gce.archive.ubuntu.com/ubuntu focal-updates InRelease
Hit:3 http://us-central1.gce.archive.ubuntu.com/ubuntu focal-backports InRelease
Hit:4 http://security.ubuntu.com/ubuntu focal-security InRelease
Reading package lists... Done
Building dependency tree
Reading state information... Done
16 packages can be upgraded. Run 'apt list --upgradable' to see them.
Reading package lists... Done
Building dependency tree
Reading state information... Done
curl is already the newest version (7.68.0-lubuntu2.18).
git is already the newest version (1:2.25.1-lubuntu3.10).
tmux is already the newest version (3.0a-2ubuntu0.4).
vim is already the newest version (2:8.1.2269-lubuntu5.12).
wget is already the newest version (1.20.3-lubuntu2).
openjdk-8-jdk is already the newest version (8u362-ga-0ubuntu1-20.04.1).
python3-pip is already the newest version (20.0.2-5ubuntu1.8).
The following packages were automatically installed and are no longer required:
  libatasmart4 libblockdev-fs2 libblockdev-loop2 libblockdev-part-err2 libblockdev-part2 libblockdev-swap2 libblockdev-utils2 libblockdev2 libbmbm-glib4
  libbmbm-proxy libbmm-glib0 libnumal libparted-fs-resize0 libqmi-glib5 libqmi-proxy libudisks2-0 libxmlb2 usb-modeswitch usb-modeswitch-data
Use 'sudo apt autoremove' to remove them.
0 upgraded, 0 newly installed, 0 to remove and 16 not upgraded.
me19b190@me19b190-ssh8:~$

me19b190@me19b190-ssh8:~$ pip install kafka-python pyspark pyspark[sql]
Requirement already satisfied: kafka-python in ./local/lib/python3.8/site-packages (2.0.2)
Requirement already satisfied: pyspark in ./local/lib/python3.8/site-packages (3.3.2)
Requirement already satisfied: py4j==0.10.9.5 in ./local/lib/python3.8/site-packages (from pyspark) (0.10.9.5)

me19b190@me19b190-ssh8:~$ wget https://dlcdn.apache.org/kafka/3.4.0/kafka_2.13-3.4.0.tgz # Get Kafka Sources
--2023-03-31 22:39:16-- https://dlcdn.apache.org/kafka/3.4.0/kafka_2.13-3.4.0.tgz
Resolving dlcdn.apache.org (dlcdn.apache.org)... 151.101.2.132, 2a04:4e42::644
Connecting to dlcdn.apache.org (dlcdn.apache.org)[151.101.2.132]:443... connected.
HTTP request sent, awaiting response... 200 OK
Length: 106290956 (101M) [application/x-gzip]
Saving to: 'kafka_2.13-3.4.0.tgz.1'

kafka_2.13-3.4.0.tgz.1      100%[=====>] 101.37M  183MB/s   in 0.6s

2023-03-31 22:39:16 (183 MB/s) - 'kafka_2.13-3.4.0.tgz.1' saved [106290956/106290956]
```

[illegible]

```
me19b190me19b190-asm8:~/kafka_2.13-3.4.0$ bin/kafka-topics.sh --create --topic me19b190-asm8-topic --bootstrap-server localhost:9092
Created topic me19b190-asm8-topic.
me19b190me19b190-asm8:~/kafka_2.13-3.4.0$ bin/kafka-topics.sh --list --bootstrap-server localhost:9092
me19b190-asm8-topic
test-topic
me19b190me19b190-asm8:~/kafka_2.13-3.4.0$ bin/kafka-topics.sh --describe --topic me19b190-asm8-topic --bootstrap-server localhost:9092
Topic: me19b190-asm8-topic TopicId: 55x3W61mS0-cmZ530Xdyrrw PartitionCount: 1 ReplicationFactor: 1 Configs:
      Topic: me19b190-asm8-topic Partition: 0 Leader: 0 Replicas: 0 Isr: 0
me19b190me19b190-asm8:~/kafka_2.13-3.4.0$
```

Step 4: Before proceeding to run our producer and consumer files it is important to check if the created topic is working properly (are we able to write/ broadcast etc.), for which we open a new terminal session and split the screen into two panes and run each of the “test producer” and “test consumer” in the two created panes. Kafka topic is the created topic “me19b190-ass8-topic”.

```
me19b190@me19b190-ass8:~/kafka 2.13-3.4.0$ bin/kafka-console-producer.sh --topi
c me19b190-ass8-topic --bootstrap-server localhost:9092
>I am Gangadhar
>This is a Trail run
>Exiting...
>

me19b190@me19b190-ass8:~/kafka 2.13-3.4.0$ bin/kafka-console-consumer.sh --topi
c me19b190-ass8-topic --from-beginning --bootstrap-server localhost:9092
I am Gangadhar
This is a Trail run
Exiting...
```

Thus the created topic is running successfully.

Step 5: Upload a “.csv” to the bucket. A new bucket named “me19b190-ass8” is created and we upload the “Customers.csv” data set which is downloaded from Kaggle. It has 8 columns and 2001 rows. The following figure illustrates the dataset.

CustomerID	Gender	Age	Annual Income (\$)	Spending Score (1-100)	Profession	Work Experience	Family Size
1	Male	19	15000	39	Healthcare	1	4
2	Male	21	35000	81	Engineer	3	3
3	Female	20	86000	6	Engineer	1	1
4	Female	23	59000	77	Lawyer	0	2
5	Female	31	38000	40	Entertainment	2	6
6	Female	22	58000	76	Artist	0	2
7	Female	35	31000	6	Healthcare	1	3
8	Female	23	84000	94	Healthcare	1	3
9	Male	64	97000	3	Engineer	0	3
10	Female	30	98000	72	Artist	1	4

only showing top 10 rows

Step 6: The following kafka_producer.py file reads the “.csv” file from the bucket and pushes it to the kafka topic “me19b190-ass8-topic”. The producer terminates once it pushes 100 records in batches of 10, which are separated by a sleep time of 10 seconds.

```
GNU nano 4.8      kafka_producer.py      Modified

spark = SparkSession.builder.appName("pysparkdf").getOrCreate() #create a spark session

KAFKA_TOPIC_NAME_CONS = "me19b190-ass8-topic" #topic name
KAFKA_BOOTSTRAP_SERVERS_CONS = 'localhost:9092' #local host

if __name__ == "__main__":
    storage_client = storage.Client()
    bucket = storage_client.get_bucket('me19b190-bucket-ass8')
    blob = bucket.blob('Customers.csv')
    records_str = blob.download_as_string().decode('utf-8').split('\n')
    records_str = list(csv.reader(records_str, delimiter=","))
    data = records_str[1:-1] #construct the dataframe
    header = records_str[0] #header of the dataframe
    df = spark.createDataFrame(data, header)

    num_records_terminate = 100 #termination number of records = 100
    print("Kafka Producer Application Started ... ")
    #create the producer object
    kafka_producer_obj = KafkaProducer(bootstrap_servers=KAFKA_BOOTSTRAP_SERVERS_CONS,value_serializer=lambda x: dumps(x).encode('utf-8'))
    df = df.withColumn("CustomerID", df['CustomerID'].cast("Integer"))
    i = 1
    while i <= num_records_terminate+1:
        #collect batches of 10 records each and push them to topic
        print("Current records: ")
        records = (df.filter((df.CustomerID >= i - 10) & (df.CustomerID < i)))
        records.show()
        for row in records.rdd.collect():
            kafka_producer_obj.send(KAFKA_TOPIC_NAME_CONS, row)
        #increment i
        i = i+10
        time.sleep(10) #sleep time of the producer, send data after a wait time of 10 seconds

    print("Kafka Producer Application Completed. ")

Get Help  Write Out  Where Is  Cut Text  Justify  Cur Pos  M-U Undo  M-M Mark Text  M-] To Bracket  M-; Previous
Exit      Read File  Replace  Paste Text  To Spell  Go To Line  M-R Redo  M-C Copy Text  M-_ Where Was  M-N Next
Ctrl-Ove Ctrl-Unde Ctrl-Undo Ctrl-Undo Ctrl-Undo Ctrl-Undo Ctrl-Undo Ctrl-Undo Ctrl-Undo Ctrl-Undo
```

Step 7: The below given file is the code for consumer. It reads the data from kafka every 5 seconds and emits the count of rows seen in the last 10 seconds.

```

GNU nano 4.8                                kafka_streaming_json_demo.py                                Modified
import time

kafka_topic_name = "me19b190-ass8-topic"
kafka_bootstrap_servers = 'localhost:9092'

if __name__ == "__main__":
    #create a spark session
    spark = SparkSession\
        .builder\
        .appName("ME19B190-Ass8-consumer")\
        .getOrCreate()
    spark.sparkContext.setLogLevel("ERROR")

    # Create DataFrame from stream of input lines from the localhost 9092
    records = spark\
        .readStream\
        .format('kafka')\
        .option('kafka.bootstrap.servers', kafka_bootstrap_servers)\
        .option('subscribe', kafka_topic_name)\
        .option("startingOffsets", 'latest')\
        .load()

    #count the number of rows in the read input with the desired window size
    windowedCounts = records.groupby(
        window(records.timestamp, '10 seconds', '5 seconds')).count()

    query = windowedCounts\
        .writeStream\
        .outputMode('complete')\
        .option('truncate', 'false')\
        .format('console')\
        .start()

    query.awaitTermination()

```

Step 8: Open a new terminal session and split the terminal into two panes and run the producer and consumer file separately in each of the panes using the following commands

Command to execute the consumer file: *"spark-submit --packages org.apache.spark:spark-sql-kafka-0-10_2.12:3.3.2 kafka_streaming_json_demo.py"*

```

me19b190@me19b190-ass8:~/kafka_2.13-3.4.0/ME19B190-ass8-producer_consumer_files
$ spark-submit --packages org.apache.spark:spark-sql-kafka-0-10_2.12:3.3.2 kafk
a_streaming_json_demo.py

```

Command to execute the producer file: *"python3 kafka_producer.py"*

Step 9: The output of the producer(left) and consumer(right) is shown below. The producer prints the current records (separated by a sleep time of 10 seconds) it is pushing to the kafka and the consumer reads the data from kafka every 5 seconds and emits the count of rows seen in the last 10 seconds.

CustomerID	Gender	Age	Annual Income (\$)	Spending Score (1-100)	Profession	Work Experience	Family Size
9	Male	64	97000	3	Engineer		
10	Female	30	98000	72	Artist		
11	Male	67	7000	14	Engineer		
12	Female	35	93000	99	Healthcare		
13	Female	58	80000	15	Executive		
14	Female	24	91000	77	Lawyer		
15	Male	37	19000	13	Doctor		
16	Male	22	51000	79	Healthcare		
17	Female	35	29000	35	Homemaker		
18	Male	20	89000	66	Healthcare		
19	Male	52	20000	29	Entertainment		
20	Female	35	62000	98	Artist		

Batch	Window	Count
Batch: 0		
	(2023-03-31 22:24:25, 2023-03-31 22:24:35)	10
	(2023-03-31 22:24:20, 2023-03-31 22:24:30)	10
Batch: 1		
	(2023-03-31 22:24:25, 2023-03-31 22:24:35)	10
	(2023-03-31 22:24:20, 2023-03-31 22:24:30)	10
Batch: 2		
	(2023-03-31 22:24:25, 2023-03-31 22:24:35)	10
	(2023-03-31 22:24:20, 2023-03-31 22:24:30)	10
	(2023-03-31 22:24:35, 2023-03-31 22:24:45)	10
	(2023-03-31 22:24:30, 2023-03-31 22:24:40)	10

	49 Female 29 0 1	78000	42 Healthcare		(2023-03-31 22:55:00, 2023-03-31 22:55:10) 10		
	50 Female 31 1	25000	42 Engineer		(2023-03-31 22:55:10, 2023-03-31 22:55:20) 10		
						(2023-03-31 22:54:55, 2023-03-31 22:55:05) 10	
						(2023-03-31 22:55:15, 2023-03-31 22:55:25) 10	
						Batch: 3	
						window	count
	51 Female 49 1	88000	52 Artist		(2023-03-31 22:55:00, 2023-03-31 22:55:10) 10		
	52 Male 33 1	97000	60 Artist		(2023-03-31 22:55:10, 2023-03-31 22:55:20) 10		
	53 Female 31 1	74000	54 Artist		(2023-03-31 22:54:55, 2023-03-31 22:55:05) 10		
	54 Male 59 0	68000	60 Entertainment		(2023-03-31 22:55:25, 2023-03-31 22:55:35) 10		
	55 Female 50 1	18000	45 Marketing		(2023-03-31 22:55:15, 2023-03-31 22:55:25) 10		
	56 Male 47 14	95000	41 Artist		(2023-03-31 22:55:20, 2023-03-31 22:55:30) 10		
	57 Female 51 9	71000	50 Healthcare		Batch: 4		
	58 Male 69 8	8000	46 Doctor		window	count	
	59 Female 27 1	57000	51 Artist				
	60 Male 53 1	89000	46 Lawyer		(2023-03-31 22:55:00, 2023-03-31 22:55:10) 10		
						(2023-03-31 22:55:10, 2023-03-31 22:55:20) 10	
						(2023-03-31 22:55:30, 2023-03-31 22:55:40) 10	
						(2023-03-31 22:55:35, 2023-03-31 22:55:45) 10	
						(2023-03-31 22:54:55, 2023-03-31 22:55:05) 10	
						(2023-03-31 22:55:25, 2023-03-31 22:55:35) 10	
						(2023-03-31 22:55:15, 2023-03-31 22:55:25) 10	
						(2023-03-31 22:55:20, 2023-03-31 22:55:30) 10	
	69 Male 19 0	81000	59 Artist		(2023-03-31 22:55:20, 2023-03-31 22:55:30) 10		
	70 Female 32 2	97000	47 Homemaker		Batch: 4		
						window	count
						(2023-03-31 22:55:00, 2023-03-31 22:55:10) 10	
						(2023-03-31 22:55:10, 2023-03-31 22:55:20) 10	
						(2023-03-31 22:55:30, 2023-03-31 22:55:40) 10	
						(2023-03-31 22:55:35, 2023-03-31 22:55:45) 10	
	71 Male 70 8	91000	55 Healthcare		(2023-03-31 22:54:55, 2023-03-31 22:55:05) 10		
	72 Female 47 8	9000	42 Entertainment		(2023-03-31 22:55:25, 2023-03-31 22:55:35) 10		
	73 Female 60 0	44000	49 Artist		(2023-03-31 22:55:15, 2023-03-31 22:55:25) 10		
	74 Female 60 8	10000	56 Healthcare		(2023-03-31 22:55:20, 2023-03-31 22:55:30) 10		
	75 Male 59 0	89000	47 Artist		Batch: 5		
	76 Male 26 0	49000	54 Homemaker		window	count	
	77 Female 45 0	67000	53 Doctor				
	78 Male 40 4	99000	48 Artist		(2023-03-31 22:55:00, 2023-03-31 22:55:10) 10		
	79 Female 23 2	97000	52 Engineer		(2023-03-31 22:55:10, 2023-03-31 22:55:20) 10		
	80 Female 49 1	98000	42 Doctor		(2023-03-31 22:55:30, 2023-03-31 22:55:40) 10		
						(2023-03-31 22:55:45, 2023-03-31 22:55:55) 10	
						(2023-03-31 22:55:20, 2023-03-31 22:55:30) 10	
	89 Female 34 7	76000	60 Entertainment		window	count	
	90 Female 50 0	56000	46 Doctor				
						(2023-03-31 22:55:00, 2023-03-31 22:55:10) 10	
						(2023-03-31 22:55:10, 2023-03-31 22:55:20) 10	
						(2023-03-31 22:55:30, 2023-03-31 22:55:40) 10	
						(2023-03-31 22:55:35, 2023-03-31 22:55:45) 10	
						(2023-03-31 22:54:55, 2023-03-31 22:55:05) 10	
						(2023-03-31 22:55:25, 2023-03-31 22:55:35) 10	
						(2023-03-31 22:55:15, 2023-03-31 22:55:25) 10	
						(2023-03-31 22:55:40, 2023-03-31 22:55:50) 10	
						(2023-03-31 22:55:45, 2023-03-31 22:55:55) 10	
						(2023-03-31 22:55:20, 2023-03-31 22:55:30) 10	
						Batch: 6	
	91 Female 68 0	46000	55 Entertainment		window	count	
	92 Male 18 1	36000	41 Artist				
	93 Male 48 1	88000	49 Artist		(2023-03-31 22:56:05, 2023-03-31 22:56:15) 10		
	94 Female 40 1	31000	40 Entertainment		(2023-03-31 22:55:00, 2023-03-31 22:55:10) 10		
	95 Female 32 0	24000	42 Artist		(2023-03-31 22:56:00, 2023-03-31 22:56:10) 10		
	96 Male 24 10	80000	52 Artist		(2023-03-31 22:55:10, 2023-03-31 22:55:20) 10		
	97 Female 47 0	2000	47 Artist		(2023-03-31 22:55:30, 2023-03-31 22:55:40) 10		
	98 Female 27 0	67000	50 Artist		(2023-03-31 22:55:35, 2023-03-31 22:55:45) 10		
	99 Male 48 1	58000	42 Doctor		(2023-03-31 22:56:55, 2023-03-31 22:56:05) 10		
	100 Male 20 3	80000	49 Engineer		(2023-03-31 22:54:55, 2023-03-31 22:55:05) 10		
						(2023-03-31 22:55:25, 2023-03-31 22:55:35) 10	
						(2023-03-31 22:55:15, 2023-03-31 22:55:25) 10	
						(2023-03-31 22:55:40, 2023-03-31 22:55:50) 10	
						(2023-03-31 22:55:45, 2023-03-31 22:55:55) 10	
						(2023-03-31 22:55:20, 2023-03-31 22:55:30) 10	
						(2023-03-31 22:56:10, 2023-03-31 22:56:20) 10	

Step 10: Clean up the logs using the following command “*rm -rf /tmp/kafka-logs /tmp/zookeeper /tmp/kraft-combined-logs*”

```
me19b190@me19b190-ass8:~/kafka_2.13-3.4.0/ME19B190-ass8-producer_consumer_files$
rm -rf /tmp/kafka-logs /tmp/zookeeper /tmp/kraft-combined-logs
me19b190@me19b190-ass8:~/kafka_2.13-3.4.0/ME19B190-ass8-producer_consumer_files$
```