# Programming Assignment № 2

Evangelia Panourgia, Athens University of Economics and Business University      27/03/2025

## Part I

### The python script for the first part.

The provided Python script is an asynchronous Kafka producer designed to simulate real-time movie rating events. It uses the `Faker` library to generate synthetic user names and reads actual movie titles from a CSV file using `pandas`. Each generated user rates a randomly selected movie with a random rating between 1 and 10, along with a current timestamp. These events are serialized to JSON and sent to a Kafka topic named `"test"` using the `aiokafka` library. The process runs continuously, with each user producing one rating per minute, making it ideal for streaming applications such as those built with Apache Spark Structured Streaming. Additionally, the script appends a specific user name, "Evangelia Panourgia", to ensure consistent personal data inclusion for testing or filtering.

Listing 1: python-kafka-example

```
1  """
2  Async Kafka Producer: Fake Movie Ratings Generator
3
4  Overview:
5  ---------
6  This script continuously generates **fake movie rating events** and sends ←
       them to a Kafka topic ("test").
7  It simulates users rating random movies with timestamps, mimicking real-←
       time data.
8
9  Components:
10 -----------
11 1. Faker: Generates fake human names.
12 2. Pandas: Loads real movie titles from a CSV file.
13 3. aiokafka: Asynchronous Kafka producer library.
14 4. Kafka: Acts as the message broker for event streaming.
15
16 Use Case:
17 ---------
18 Run this script to simulate live rating data for a Spark Structured ←
       Streaming application to consume and process.
19 """
```

```python
20  import json
21  import asyncio
22  import random
23  from datetime import datetime
24
25  import pandas as pd
26  from faker import Faker
27  from aiokafka import AIOKafkaProducer
28
29  # Load movie titles
30  movies_df = pd.read_csv('data/movies.csv', header=None, names=["↵
      movie_title"])
31  movie_titles = movies_df['movie_title'].dropna().tolist()
32
33  # Generate names
34  fake = Faker()
35  names = [fake.name() for _ in range(10)]
36  names.append("Evangelia Panourgia") # add my own name
37
38  # Kafka settings
39  KAFKA_TOPIC = "test"
40  KAFKA_BOOTSTRAP_SERVERS = "localhost:29092"
41
42  # JSON serializer
43  def serializer(data):
44      return json.dumps(data).encode()
45
46  # Producer logic
47  async def produce():
48      producer = AIOKafkaProducer(
49          bootstrap_servers=KAFKA_BOOTSTRAP_SERVERS,
50          value_serializer=serializer,
51          compression_type="gzip"
52      )
53
54      await producer.start()
55      try:
56          while True:
57              for name in names:
58                  movie = random.choice(movie_titles)
59                  rating = random.randint(1, 10)
60                  timestamp = datetime.now().isoformat()
61
62                  message = {
63                      "name": name,
64                      "movie": movie,
65                      "timestamp": timestamp,
```

```
66                "rating": rating
67            }
68
69            print(f"Sending: {message}")
70            await producer.send(KAFKA_TOPIC, message)
71            await asyncio.sleep(60) # One message per user per minute
72    finally:
73        await producer.stop()
74
75 # Run it (no deprecation warning)
76 if __name__ == "__main__":
77    loop = asyncio.new_event_loop()
78    asyncio.set_event_loop(loop)
79    loop.run_until_complete(produce())
```

The image below depicts the output of kafka prodicer runing in the first terminal.



Figure 1: Terminal 1: Kafka Producer.

An additional modification was applied to the console script located in the `examples` directory, specifically to the file named `console-spark-streaming-example.py`. This adjustment was necessary to tailor the script to the project's requirements and ensure compatibility with the streaming architecture.

The image below depicts the console output helping us to debug the kafka producer real-time data. It runs in another terminal.

Figure 2: Terminal 2: Console for Kafka Producer.

# Part 2

## The pyspark script for the second part.

This PySpark application establishes a real-time data pipeline that reads movie rating events from a Kafka topic, enriches them with metadata from a static Netflix dataset, and writes the processed results into a Cassandra database. The script defines two schemas: one for parsing JSON-encoded rating messages received from Kafka, and another for reading static metadata from a CSV file. It converts timestamps to proper `TimestampType`, and groups data by hourly buckets for time-based aggregations. Ratings are joined with metadata on the movie title to create an enriched DataFrame, which is then written to Cassandra using the `foreachBatch` method in micro-batch mode. This architecture enables structured streaming analysis and time-windowed queries for individual users and movies.

Listing 2: cassandra-spark-streaming-example.py.

```
1  from pyspark.sql import SparkSession
2  from pyspark.sql.types import StructType, StructField, StringType, ↩
       IntegerType, TimestampType
3  from pyspark.sql.functions import from_json, col, to_timestamp, ↩
       date_format
4
5  # Schema for Kafka messages (movie ratings)
6  ratingSchema = StructType([
7      StructField("name", StringType(), False),
8      StructField("movie", StringType(), False),
9      StructField("timestamp", StringType(), False),
10     StructField("rating", IntegerType(), False)
11 ])
12
13 # Schema for Netflix CSV
14 netflixSchema = StructType([
15     StructField("show_id", StringType(), True),
16     StructField("title", StringType(), False),
17     StructField("director", StringType(), True),
```

```python
        StructField("country", StringType(), True),
        StructField("release_year", StringType(), True),
        StructField("rating", StringType(), True),
        StructField("duration", StringType(), True)
])

# Initialize Spark session
spark = (
    SparkSession.builder
    .appName("MovieRatingStreamer")
    .config("spark.jars.packages", "org.apache.spark:spark-sql-kafka-0-10↵
        _2.12:3.5.0,com.datastax.spark:spark-cassandra-connector_2↵
        .12:3.4.1")
    .config("spark.cassandra.connection.host", "localhost")
    .getOrCreate()
)

spark.sparkContext.setLogLevel("ERROR")

# Read Netflix CSV with renamed rating column to avoid conflict with "↵
    rating" in Kafka schema
netflix_df = (
    spark.read.schema(netflixSchema)
        .option("header", True)
        .csv("data/netflix.csv")
        .withColumnRenamed("rating", "rating_category")  # Rename avoids ↵
            conflict
        .cache()
)

# Read streaming data from Kafka
df = (
    spark.readStream.format("kafka")
        .option("kafka.bootstrap.servers", "localhost:29092")
        .option("subscribe", "test")
        .option("startingOffsets", "latest")
        .load()
)

# Parse JSON from Kafka messages
ratings_df = (
    df.selectExpr("CAST(value AS STRING)")
      .select(from_json(col("value"), ratingSchema).alias("data"))
      .select("data.*")
      .withColumn("timestamp", to_timestamp("timestamp"))  # Convert ↵
          string timestamp to TimestampType
      .withColumn("hour_bucket", date_format("timestamp", "yyyy-MM-dd HH↵
```

```python
            :00"))  # Grouping by hour
60  )
61
62  # Join ratings with static Netflix metadata
63  enriched_df = (
64      ratings_df.join(netflix_df, ratings_df.movie == netflix_df.title, "↩
            left")
65                  .drop("title")
66  )
67
68  # Optional: Print schema for debugging (can be removed later)
69  enriched_df.printSchema()
70
71  # Define Cassandra write logic
72  def writeToCassandra(writeDF, _):
73      (
74          writeDF.select(
75              "name", "movie", "timestamp", "rating", "hour_bucket",
76              "show_id", "director", "country", "release_year", "↩
                    rating_category", "duration"
77          )
78          .write
79          .format("org.apache.spark.sql.cassandra")
80          .mode('append')
81          .options(table="movie_ratings", keyspace="netflix_ks")
82          .save()
83      )
84
85  # Write to Cassandra in a streaming loop
86  result = None
87  while result is None:
88      try:
89          result = (
90              enriched_df.writeStream
91                  .foreachBatch(writeToCassandra)
92                  .outputMode("update")
93                  .option("checkpointLocation", "/tmp/checkpoints/movie")
94                  .trigger(processingTime="30 seconds")
95                  .start()
96                  .awaitTermination()
97          )
98      except Exception as e:
99          print(f"Streaming error: {e}")
```

**Details about your Cassandra data model.**

The Cassandra data model was carefully designed to optimize query performance for the specific requirement of aggregating movie ratings by user and hour. The primary key is composed of a composite partition key (`name, hour_bucket`), ensuring that all ratings provided by a specific user within a given hour are stored in the same partition. This design eliminates the need for inefficient filtering operations and allows for fast, targeted queries such as retrieving all movies rated by a user during a specific hour or calculating the average rating and duration of those movies. The clustering columns (`timestamp, movie`) further organize the data chronologically within each partition, enabling efficient range queries and ordered retrieval. Crucially, the fields `rating` and `duration` are stored as integers, allowing native support for aggregation functions like `AVG()`, which would not be possible with text fields. Overall, this schema supports the core analytical tasks required by the application while adhering to Cassandra's best practices for scalability and performance.

Listing 3: Cassandra schema for movie ratings

```
1  CREATE KEYSPACE IF NOT EXISTS netflix_ks
2  WITH replication = {'class': 'SimpleStrategy', 'replication_factor': 1};
3
4  CREATE TABLE IF NOT EXISTS netflix_ks.movie_ratings (
5      name text,
6      hour_bucket text,
7      timestamp timestamp,
8      movie text,
9      rating int,
10     show_id text,
11     director text,
12     country text,
13     release_year int,        -- changed from text : int
14     rating_category text,
15     duration int,            -- changed from text : int
16     PRIMARY KEY ((name, hour_bucket), timestamp, movie)
17 );
```

**A sample of persisted lines (around 50) of your Cassandra table.**

To inspect the data stored in the `movie_ratings` table of the `netflix_ks` keyspace, we can use the `cqlsh` shell interface provided by Apache Cassandra. The following query retrieves approximately 50 rows that were previously persisted by the streaming pipeline:

```
SELECT * FROM netflix_ks.movie_ratings LIMIT 50;
```

Figure 3: Terminal 3: Spark pre-process, enrich with metadata..

This query provides a representative snapshot of the data ingested into Cassandra, including key fields such as `name`, `movie`, `timestamp`, `rating`, and `duration`. Analyzing these records ensures data integrity, validates the correctness of the ingestion pipeline, and confirms that the schema design supports the desired querying capabilities.



Figure 4: A sample output of the Cassandra table with 50 persisted entries.

**Two CQL queries and their results in your database about your own name and a particular hour that generate the average runtime of the movies that you've rated during this hour, and the names of the movies, respectively.**

To support time-based analytics in our Cassandra data model, we designed the schema to partition data by `name` and `hour_bucket`, allowing efficient access to all ratings a specific user made during a given hour. The first query retrieves the titles of movies rated by the user *Evangelia Panourgia* during the hour of `2025-03-27 09:00` using a direct lookup on the partition key. The second query, which includes the `ALLOW FILTERING` clause, calculates the average duration of those rated movies within the same hour. While Cassandra does not natively support aggregations without filtering, our schema minimizes performance issues by narrowing queries to a specific partition, ensuring practical execution even with the filtering clause.

Figure 5: Names of the movies you rated during a particular hour SELECT movie.



Figure 6: Average duration of the movies you rated during the same hour SELECT avg(duration).

## Hosted on GitHub Repository - Deploy Instructions.

The complete implementation of this project is hosted on GitHub and is publicly accessible at `https://github.com/e-panourgia/large-data-kafka-cassandra`. The repository contains all necessary components, including Kafka producers, Spark Structured Streaming jobs, and Cassandra schemas. To execute the project, clone the repository using `git clone https://github.com/e-panourgia/large-data-kafka-cassandra`, then follow the instructions in the README file to set up Docker containers, start Kafka and Cassandra services, and run the streaming application using `spark-submit`. Ensure that you have installed Docker, Docker Compose, and Spark locally or are running within the provided Vagrant virtual machine for a consistent and reproducible development environment.

Important note, this implementation was designed to run on mac m4 pc.

## Acknowledgements

I would like to sincerely thank **Professor Liakos Panayiotis** for his invaluable guidance, support, and dedication throughout the entire teaching process. His insightful lectures and continuous encouragement greatly contributed to the successful completion of this project.