Large Scale Data Management
Athens University of Economics and Business
M.Sc. in Data Science

Programming Project II

Theofanis Nitsos – p3352325

## Part I

The part_1.py file is the script for the 1st part of the project that produces messages to the Kafka topic.

A custom JSON encoder is used to overcome some issues that came up during the execution of this script. Then, a loop limited to 1000 counts is used to generate 1000 songs, person name and current timestamp combinations with a 1 second delay. 10 random names are generated using the Faker library and the list is appended with the name 'Mr. Fanis Nitsos', while the song names are extracted from the provided csv file.

## Part II

The part_2.py file is the script for the 2nd part of the project. This script effectively sets up a continuous stream processing job using Spark, consuming data from Kafka, processing it, and persisting the results into a Cassandra database.

In detail this script receives JSON information from Kafka (defined in Part I) in the schema of person name, timestamp and song name. It initializes a Spark session and reads the data from the provided CSV. The information from the CSV is joined with the information from Kafka based on the song name and are streamed to Cassandra with a 30 second interval.

## Cassandra

The Cassandra table is created using the variable person name (personname) as partition key and the timestamp (listenattime) as clustering key to facilitate aggregations for a particular person and hour. The commands for Cassandra are included in the cassandra.txt file.

```
CREATE TABLE spotify.fanis_records_2 (
    personname TEXT,
    listenattime TIMESTAMP,
    name TEXT,
    song TEXT,
```

```
    artists TEXT,
    duration_ms BIGINT,
    album_name TEXT,
    album_release_date DATE,
    danceability FLOAT,
    energy FLOAT,
    key INT,
    loudness FLOAT,
    mode INT,
    speechiness FLOAT,
    acousticness FLOAT,
    instrumentalness FLOAT,
    liveness FLOAT,
    valence FLOAT,
    tempo FLOAT,
    PRIMARY KEY (personname, listenattime)
);
```

Subsequently 2 queries are executed, using Cassandra's features to our advantage; filtering the data based on the partition and clustering keys.

1st query calculates the average danceability for the user Mr. Fanis Nitsos between two timestamps

select avg(danceability) from spotify.fanis_records_2 where personname='Mr. Fanis Nitsos' and listenattime >= '2024-03-28 21:20:00' and listenattime <= '2024-03-28 21:25:00';


2nd query retrieves the songs player for the user Mr. Fanis Nitsos between two timestamps

select name from spotify.fanis_records_2 where personname='Mr. Fanis Nitsos' and listenattime >= '2024-03-28 21:20:00' and listenattime <= '2024-03-28 21:25:00';

```
(1 rows)
cqlsh> select * from spotify.fanis_records_2 limit 50;

 personname     | listenattime                  | acousticness | album_name                                                           | album_release_date | artists                          | danceability | dura
tion_ms | energy | instrumentalness | key | liveness | loudness | mode | name                                                       | song                                            | speechiness | tempo        | valence
```

Figure 1 Output of the 1st Cassandra Query & 50 persisted lines of the Cassandra table

```
(1 rows)
cqlsh> select name from spotify.fanis_records_2 where personna

 name
------------------------------------------------------
                                      WEST AFRICA TIME
                                                EXTEND
                                Beverly Hills Freestyle
                                     Different Pattern
                                              Estrella
                                            My Brother
                                          Petit génie
                                        Yüreğim Ağlar
                                                 Drama
                                                 Laços
                                            BETTER NOW
                                                 Turné
                                                Anders
                                                  24/7
                                                Daheya
                                          Mia Na' Mas
                                                  2:30
                                            Allegria...
                                            Ciao Bella
                               Niet Genoeg (feat. Idaly)
                                         Diamond Days
                                          Wacht Op Mij
                                Nu tändas tusen juleljus
                          See You Again (feat. Kali Uchis)
                                            Ted a tady
                               Hjertesorg Betaler Regninga
                                           Mexri Telous
                                              Mentirosa
                                            Friesenjung
                                              Memories
                                                  Mala
                                              ROCKSTAR
                                  Cold Heart – PNAU Remix
                                      DIS-MOI (feat. SDM)
                     I Remember Everything (feat. Kacey Musgraves)
                                          Reyah El Hayah
                                                   Dum
                                          Arjan Vailly
                               Banan Melon Kiwi & Citron
                                     Seni Dert Etmeler
                                      Viimeisiä sanoja
                                             Algo De Mi
                                                你的世界
                                        Para Sa Streets
                                   lovely (with Khalid)
                                              Името ти
```

```
            I Remember Everything (feat. Kacey Musgraves)
                                      Reyah El Hayah
                                               Dum
                                      Arjan Vailly
                           Banan Melon Kiwi & Citron
                                 Seni Dert Etmeler
                                  Viimeisiä sanoja
                                         Algo De Mi
                                            你的世界
                                    Para Sa Streets
                               lovely (with Khalid)
                                          Името ти
                                          Ma Jolie
                       Il Salto (feat. Massimo Pericolo)
                                         Wunderbar
                                     O helga natt
                                            Pullup
                                         mielipide
                          DIMMI CHE NON È UN ADDIO
                       Alt vi har er nu (Artigeardit, Lamin)
                                         na Circus
                                           Periodt
                                       Sulle tehty
                                     Prywatny bal
                                        Aiga qarap
                              ดาวหางอ่ำ ลาลย์
                                           Destiny
                                              Özür
                            Je te laisserai des mots
        Yes or No (Feat. HUH YUNJIN of LE SSERAFIM, Crush)
                                       Scary Movie
                                   Tanrım Reva Mı
                                           Madison
                                  Coração de Gelo
                                       ENA TSIGARO
                                             TOXIC
                             Er Her (Artigeardit, KESI)
                            Pujaanku (feat. Aisyah Aziz)
                          VIDÍM JAK SA NA MŇA POZERÁŠ
                                              GEME
                                              Cast
                                         Bráðna
                             Jaga Jaga (with Babyboy AV)
                                           Navidad
                                           BILLETS
                                      Roll With It
                                     väljateenitud
                           (It Goes Like) Nanana – Edit
                                          На зape

(83 rows)
cqlsh>
```

*Output of the 2nd Cassandra query*