

# Managing ML Data and Models for Spotify Music Recommendation

## Final Report

**Name:** Kehan Bi

**USC ID:** 2608579504

**Course:** DSCI 551, Tue, Fall 2021

**Date:** December 1, 2021

### I. About Me:

My background is Mathematics and minor in Business in Undergraduate school. This is one-person group and I did not have computer science background or UI background.

### II. Project Introduction and Motivation:

Music is an indispensable part of life. Music accompanies people on the way to and get off work, time before going to bed, and even an alarm to wake up. I love listening to music, but I often feel troubled because there is usually not enough music in my playlist. The music software recommended me music that I might like every day. Everyone has a specific type of music they like. So how does Spotify recommend music to us? The goal of this project is to recommend few artists that users may like based on genre of music which given by users.

### III. Parts and Functions in project:

(The details will be explained later)

#### 1. Cloud Database

#### 2. Data preparation

- a) Clean and transform data
- b) Extract features and metadata from raw data
- c) Upload raw data, metadata, and features to Firebase

#### 3. Interface for user to upload file and upload metadata

- a) User can upload .csv or .json files
- b) User can customize the special split symbols for their .csv files
- c) User can type in the description of dataset, the input from users will storage in dataset.
- d) Metadata and features are auto extraction by app when user click “upload” to upload the file.
- e) Metadata, features, description, and dataset upload to Firebase at same time.

#### 4. Interface for user to view the dataset from Firebase:

- a) User can select the datasets, explore the description, metadata, features, first 10 lines, last 10 lines, and random 10 lines of datasets.
- b) The new dataset which uploaded by users show in the list with their metadata. Users able to explore them.

## 5. Machine Learning Method - LightFM:

- a) Set up ML model
- b) Set the parameter

## 6. Show recommended result

- a) User can type in the genre of music they like

## IV. Cloud Database:

The data is storage in Firebase Realtime Database:

The screenshot shows the Firebase Realtime Database interface. On the left, there's a sidebar with various project settings like Authentication, Firestore Database, and Machine Learning. The main area is titled "Realtime Database" and has tabs for Data, Rules, Backups, and Usage. A banner at the top right says "Protect your Realtime Database resources from abuse, such as billing fraud or phishing" and "Configure App Check". Below this, a message indicates "Read-only and non-realtime mode activated in the data viewer to improve browser performance. Select a key with fewer records to edit or view in realtime". The data viewer shows a single node under "final-project-551-8f1e5-default-rtdb/Spotify/991". The data is as follows:

```

{
  "991": {
    "acousticness": 0.99,
    "artist_id": 2,
    "artist_name": "Leopold Stokowski",
    "danceability": 0.43,
    "duration_ms": 114480,
    "energy": 0.116,
    "genre": "Movie",
    "instrumentalness": 0.908,
    "key": "G",
    "liveness": 0.121,
    "loudness": -24.73,
    "mode": "Major",
    "popularity": 8
  }
}

```

Database location: United States (us-central1)

## V. Datasets:

- SpotifyFeature.csv download from Kaggle
- Features: Genre; Artist\_name; Track\_name; Track\_id; Popularity; Acousticness; Danceability; Duration\_ms; Energy; Instrumentalness; Key; Liveness; Loudness; Mode; Speechiness; Tempo; Time\_signature; Valence; Artist\_id
- Genre, artist\_name, popularity, and artist\_id will be used in my model
- I choose the data with the Artist who has music records more than 50 times. Because the filtered artists are more popular and make the train dataset in more general condition. The result of recommendation will be more reliability.
- The dataset is too large to upload to and request from Firebase, I selected top 10000 rows for upload and use in model.

### Data preparation:

Link of Firebase: <https://final-project-551-8f1e5-default-rtdb.firebaseio.com/.json>

#### 1. Upload dataset to Firebase:

I give each “artist\_name” an “actist\_id”, because integer type of data is easier to run in machine learning model.

```

1 #upload dataset to Firebase Realtime database
2 spotify = pd.read_csv("data/SpotifyFeatures.csv", error_bad_lines=False, sep = ',')
3 #use data which the artist occurs more than 50 times
4 spotify = spotify.groupby('artist_name').filter(lambda x : len(x)>=50)
5 df_artist = pd.DataFrame(spotify["artist_name"].unique())
6 df_artist = df_artist.reset_index()
7 df_artist = df_artist.rename(columns={'index':'artist_id', 0:'artist_name'})
8 df_artist
9 spotify = pd.merge(spotify , df_artist, how='inner', on='artist_name')
10
11 jsondata = spotify.head(100000).to_json(orient = "records")
12 url = "https://final-project-551-8f1e5-default-rtdb.firebaseio.com/Spotify.json"
13 response = requests.put(url, jsondata)
14
15 print("-----Data Uploaded To Firebase Successful-----")

```

-----Data Uploaded To Firebase Successful-----

- Extract metadata and features from dataset and upload to Firebase: (The users could view data and metadata through interface, I will explain it later).

This part gets the file size, file type, number of rows and columns, and how many features and what they are.

```

1 #upload metadata to firebase
2 url2 = "https://final-project-551-8f1e5-default-rtdb.firebaseio.com/Spotify/Metadata.json"
3 path = "data/SpotifyFeatures.csv"
4
5 fsize = os.path.getsize(path)
6 fsize = fsize/float(1024*1024)
7 fsize = str(round(fsize,2)) + ' MB'
8 ext = os.path.splitext(path)[-1]
9 rows_num = spotify.shape[0]
10 columns_num = spotify.shape[1]
11 features_num = len(spotify.columns.values)
12 features = spotify.columns.values
13
14 meta_dic = {"File_size": [fsize], "File_type": [ext], "Number of Rows": [rows_num],
15             "Number of columns": [columns_num], "Number of features": [features_num], "Features": [features],
16             "Description": ["Spotify music information such as artist, genre of music, track name, popularity, etc"]}
17
18 describe = pd.DataFrame(meta_dic)
19 describe = describe.to_json(orient = "records")
20 des = requests.put(url2, describe)

```

The metadata shows in Firebase as following:

The screenshot shows the Firebase Realtime Database interface. At the top, there's a URL bar with the address <https://final-project-551-8f1e5-default-rtdb.firebaseio.com/>. Below the URL bar, a message indicates "Read-only and non-realtime mode activated in the data viewer to improve browser performance". It also says "Select a key with fewer records to edit or view in realtime".

The main content area displays a hierarchical tree structure. The root node has several child nodes labeled 99994, 99995, 99996, 99997, 99998, and 99999. Underneath the root node, there is a single child node named "Metadata". The "Metadata" node has one child node, "0", which contains the following data:

- Description: "Spotify music information such as artist, genre..."
- Features
- File\_size: "32.15MB"
- File\_type: ".csv"
- Number of columns: 19
- Number of features: 19
- Number of Rows: 120265

At the bottom of the tree view, there is a button labeled "99001 more values...".

At the very bottom of the screen, there is a footer note: "Database location: United States (us-central1)".

The screenshot shows the Firebase Realtime Database Data Viewer interface. The URL in the address bar is <https://final-project-551-8f1e5-default.firebaseio.com/>. A message at the top indicates "Read-only and non-realtime mode activated in the data viewer to improve browser performance. Select a key with fewer records to edit or view in realtime". The database structure under the key "99999" is as follows:

```

 99999
  +-- Metadata
  |    +-- 0
  |         +-- Description: "Spotify music information such as artist, genre..."
  |         +-- Features
  |             +-- 0: "genre"
  |             +-- 1: "artist_name"
  |             +-- 10: "key"
  |             +-- 11: "liveness"
  |             +-- 12: "loudness"
  |             +-- 13: "mode"
  |             +-- 14: "speechiness"
  |             +-- 15: "tempo"
  |             +-- 16: "time_signature"
  |             +-- 17: "valence"
  |             +-- 18: "artist_id"
  |
  +-- 2
  
```

Database location: United States (us-central1)

## VI. Interface for user to upload file and upload metadata:

1. This part allows the users upload .csv and .json files to Firebase.
2. When user click on the upper “Upload” button, user can choose file from local. When user click on the lower “Upload” button, the file will be upload to Firebase, and the metadata, features, and description will be upload as well.
3. I set the split symbol of csv file as “;” by default. User can customize the special split symbols for their .csv files, i.e “\t” in film1.csv from our HW5.
4. User allow to type in the description of their dataset. The content will be upload with metadata to Firebase by click lower “Upload” button.
5. Metadata and features of their files are auto extraction by my app when user click “Upload” to upload the file.
6. Metadata, features, description, and dataset upload to Firebase at same time.

Code shows as follow:

```

1 # user can upload file here
2 def on_button_upload_file(b):
3     with output1:
4         output1.clear_output()
5         uploaded_file = uploader.value
6         print("Successful upload file", list(uploaded_file.keys())[0])
7         #print(uploaded_file)
8         metadata = list(uploaded_file.values())[0]['metadata']
9         ext = os.path.splitext(metadata['name'])[-1]
10        print("Metadata:\nFile type:", ext)
11        fsize = metadata['size']
12        fsize = fsize/float(1024*1024)
13        print("File size:", str(round(fsize,2)), 'MB')
14        file_name = os.path.splitext(metadata['name'])[0]
15        #print(file_name)
16        content = list(uploaded_file.values())[0]['content']
17        describe = str(text_file_description.value)
18        if ext == '.csv': #csv file
19            if len(str(text_split_symbol.value))==0:
20                split = ";"
21            else:
22                split = str(text_split_symbol.value)
23            upload_csv = pd.read_csv(io.BytesIO(content), sep = split)
24            rows_num = upload_csv.shape[0]
25            columns_num = upload_csv.shape[1]
26            features_num = len(upload_csv.columns)
27            features = upload_csv.columns.values
28            print(features)
29            print(type(features))
30            print("There are {} rows and {} columns in dataset.".format(rows_num, columns_num))
31            print("There are {} features in dataset".format(features_num))
32            print("The features and their datatypes:")
33            for i in range(len(upload_csv.columns)):
34                print(" - " + upload_csv.columns[i] + " - " + str(upload_csv.dtypes[i]))
35
36            #upload dataset to firebase
37            url_upload_csv = "https://final-project-551-8fle5-default-firebase.com/" \
38                            + str(file_name) + ".json"
39            csv_to_json = upload_csv.to_json(orient = "records")
40            #print(csv_to_json)
41            response = requests.put(url_upload_csv, csv_to_json)
42            print("{} has been uploaded to Firebase.".format(list(uploaded_file.keys())[0]))
43
44            #upload metadata
45            url_upload_csv_metadata = "https://final-project-551-8fle5-default-firebase.com/" \
46                            + str(file_name) + "/Metadata.json"
47            meta_dic_csv = {"File_size": [fsize], "File_type": [ext], "Number of Rows": [rows_num],
48                            "Number of columns": [columns_num], "Number of features": [features_num],
49                            "Features": [features], "Description": [describe]}
50            describe_csv = pd.DataFrame(meta_dic_csv)
51            describe_csv = describe_csv.to_json(orient = "records")
52            des = requests.put(url_upload_csv_metadata, describe_csv)
53        elif ext == '.json': #json file
54            json_content = json.loads(str(content, 'utf-8'))
55            rows_num_json = len(json_content)
56            json_feature = list(json_content[0].keys())
57            features_num_json = len(json_feature)
58
59            print("There are {} rows and {} columns in dataset.".format(rows_num_json, features_num_json))
60            print("There are {} features in dataset".format(features_num_json))
61            print("The features and their datatypes:")
62            for key, value in json_content[0].items():
63                print(" - ", key, " - ", str(type(value)))
64
65            #upload dataset to firebase
66            upload_json = str(content, 'utf-8')
67            url_upload_json = "https://final-project-551-8fle5-default-firebase.com/" \
68                            + str(file_name) + ".json"
69            response_json = requests.put(url_upload_json, upload_json)
70            print("{} has been uploaded to Firebase.".format(list(uploaded_file.keys())[0]))
71
72            #upload metadata to firebase
73            url_upload_json_metadata = "https://final-project-551-8fle5-default-firebase.com/" \
74                            + str(file_name) + "/Metadata.json"
75            meta_dic_json = {"File_size": [fsize], "File_type": [ext], "Number of Rows": [rows_num_json],
76                            "Number of columns": [features_num_json], "Number of features": [features_num_json],
77                            "Features": [json_feature], "Description": [describe]}
78            describe_json = pd.DataFrame(meta_dic_json)
79            describe_json = describe_json.to_json(orient = "records")
80            des = requests.put(url_upload_json_metadata, describe_json)
81        else:
82            print("The file type is uncorrect, please upload other file.")

```

```

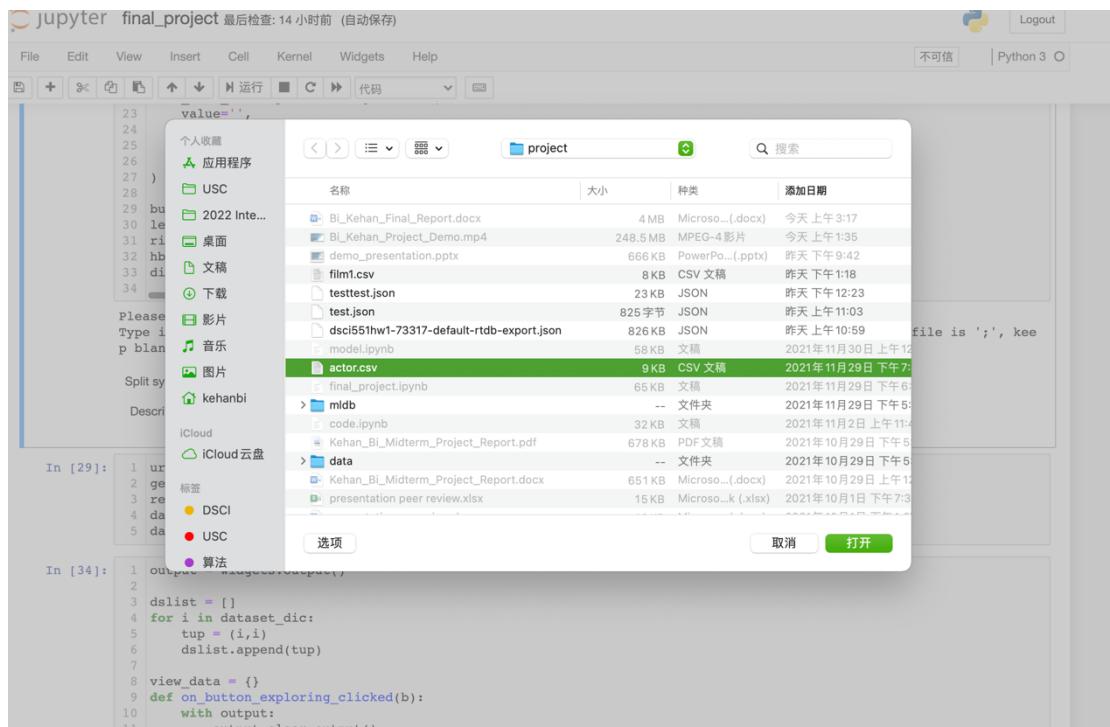
1 print("Please upload .csv or .json files here.")
2 print("Type in the special split symbol of .csv files here, e.g. ',';'\t'; etc. If the split symbol of your file is ','; keep blank." )
3 output1 = widgets.Output()
4 uploader = widgets.FileUpload(
5     accept=".csv,.json", # Accepted file extension e.g. '.txt', '.pdf', 'image/*', 'image/*,.pdf'
6     multiple=True # True to accept multiple files upload else False
7 )
8 button = widgets.Button(
9     description='Upload',
10    disabled=False,
11    button_style='', # 'success', 'info', 'warning', 'danger' or ''
12    tooltip='Click me!',
13    icon='check' # (FontAwesome names without the `fa-` prefix)
14 )
15 text_split_symbol = widgets.Text(
16     value='',
17     placeholder='Type in the special split symbol',
18     description='Split symbol:',
19     disabled=False
20 )
21
22 text_file_description = widgets.Text(
23     value='',
24     placeholder='Describe your file in short sentence',
25     description='Description:',
26     disabled=False
27 )
28
29 button.on_click(on_button_upload_file)
30 left_box = widgets.VBox([text_split_symbol, text_file_description])
31 right_box = widgets.VBox([uploader, button])
32 hbox = widgets.HBox([left_box, right_box])
33 display(hbox, output1)
34

```

Please upload .csv or .json files here.  
Type in the special split symbol of .csv files here, e.g. ',';'\t'; etc. If the split symbol of your file is ','; keep blank.

Split symbol:	<input type="text" value="Type in the special split symbol"/>	<input type="button" value="Upload (0)"/>
Description:	<input type="text" value="Describe your file in short sentence"/>	<input type="button" value="Upload"/>

- When click upper “Upload” button:



- Type in description and click upload:

Please upload .csv or .json files here.  
Type in the special split symbol of .csv files here, e.g. ','; '\t'; etc. If the split symbol of your file is ';', keep blank.

Split symbol:	Type in the special split symbol	<input type="button" value="Upload (1)"/>
Description:	<input type="button" value="▼ Upload"/>	

```
Successful upload file actor.csv
Metadata:
File type: .csv
File size: 0.01 MB
['actor_id' 'first_name' 'last_name' 'last_update']
<class 'numpy.ndarray'>
There are 200 rows and 4 columns in dataset.
There are 4 features in dataset
The features and their datatypes:
|-actor_id - int64
|-first_name - object
|-last_name - object
|-last_update - object
actor.csv has been uploaded to Firebase.
```

- The data of actor.csv shows in Firebase:

The screenshot shows the Firebase Realtime Database Data viewer. At the top, there's a header with 'final project 551' and 'Go to docs'. Below the header, it says 'Realtime Database' and has tabs for 'Data', 'Rules', 'Backups', and 'Usage'. A note below the tabs says 'Protect your Realtime Database resources from abuse, such as billing fraud or phishing' and 'Configure App Check'. The main area shows a tree structure under 'final-project-551-8f1e5-default-rtdb'. The 'actor' node is expanded, showing 200 items numbered 0 to 199. Each item has fields: 'actor\_id', 'first\_name', 'last\_name', and 'last\_update'. The first item (0) has values: 'actor\_id: 1', 'first\_name: "PENELOPE"', 'last\_name: "GUINNESS"', and 'last\_update: "2006-02-15 04:34:33"'. There are also collapsed sections for '98', '99', and 'Spotify'.

- The metadata shows in Firebase as well:

The screenshot shows the Firebase Realtime Database Data viewer. At the top, it says 'Read-only and non-realtime mode activated in the data viewer to improve browser performance' and 'Select a key with fewer records to edit or view in realtime'. The main area shows a tree structure under 'final-project-551-8f1e5-default-rtdb'. The 'Metadata' node is expanded, showing a single item '0' which contains the following details: 'Description: "This is actor.csv"', 'Features' (with indices 0 to 3: 'actor\_id', 'first\_name', 'last\_name', 'last\_update'), 'File\_size: 0.0082483292', 'File\_type: ".csv"', 'Number of columns: 4', 'Number of features: 4', and 'Number of Rows: 200'. There are also collapsed sections for '98', '99', and 'Spotify'.

## VII. Interface for user to view the dataset from Firebase:

- a) User can select the datasets, explore the description, metadata, features, first 10 lines, last 10 lines, and random 10 lines of datasets.
- b) The new dataset which uploaded by users show in the list with their metadata. Users able to explore them.

### 1. Request data from Firebase:

```

1 url_get_data = "https://final-project-551-8f1e5-default-rtdb.firebaseio.com/.json"
2 get_data_response = requests.get(url_get_data)
3 reply = get_data_response.json()
4 dataset_dic = reply.keys()
5 dataset_dic = list(dataset_dic)

```

### 2. Set up definition of user interface, I use ipywidget package.

```

1 output = widgets.Output()
2
3 dslist = []
4 for i in dataset_dic:
5     tup = (i,i)
6     dslist.append(tup)
7
8 view_data = {}
9 def on_button_exploring_clicked(b):
10     with output:
11         output.clear_output()
12         view_data = reply[dropdown_ds.value]
13         meta = pd.DataFrame(view_data["Metadata"])
14         #print(meta)
15         df2 = view_data.copy()
16         del df2['Metadata']
17         line_data = pd.DataFrame(df2)
18         if b.description == 'View Dataset':
19             print("Dataset: {}".format(dropdown_ds.value))
20             print("Description:", view_data["Metadata"][0]['Description'])
21         elif b.description == 'Metadata':
22             print(meta.T)
23             print("\nFeatures are: ")
24             for i in list(meta["Features"])[0]:
25                 print("-", i)
26                 #, list(meta["Features"])[0])
27                 #print("Dataset has {} rows and {} columns\nNumber of features = {}\nFeatures are: {}".format(valdic[0]
28         elif b.description == 'First 10 lines':
29             print(line_data.head(10).T)
30             #print(valdic[0][1].head(10))
31         elif b.description == 'Last 10 lines':
32             print(line_data.tail(10).T)
33             #print(valdic[0][1].tail(10))
34         else:
35             print(line_data.sample(n=2).T)

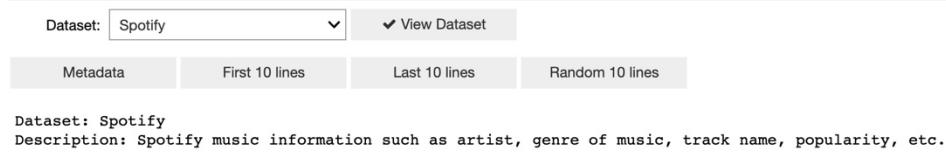
```

### 3. Interface show as follow:

```

1 dropdown_ds = widgets.Dropdown(
2     options = dslist,
3     value = "Spotify",
4     description='Dataset:',
5 )
6 button_exploring = widgets.Button(
7     description='View Dataset',
8     disabled=False,
9     button_style='',
10    tooltip='Click me2',
11    icon='check'
12 )
13 button_metadata = widgets.Button(
14     description='Metadata',
15     disabled=False,
16     button_style='',
17 )
18 button_head10 = widgets.Button(
19     description='First 10 lines',
20     disabled=False,
21     button_style='',
22 )
23 button_tail10 = widgets.Button(
24     description='Last 10 lines',
25     disabled=False,
26     button_style='',
27 )
28 button_rand10 = widgets.Button(
29     description='Random 10 lines',
30     disabled=False,
31     button_style='',
32 )
33 radiobuttons = widgets.RadioButtons(
34     options=['First 10 lines', 'Last 10 lines', 'Random 10 lines'],
35     description='Preview of data:',
36     disabled=False
37 )
38
39 button_exploring.on_click(on_button_exploring_clicked)
40 button_metadata.on_click(on_button_exploring_clicked)
41 button_head10.on_click(on_button_exploring_clicked)
42 button_tail10.on_click(on_button_exploring_clicked)
43 button_rand10.on_click(on_button_exploring_clicked)
44
45 left_box = widgets.VBox([dropdown_ds])
46 right_box = widgets.VBox([button_exploring])
47 hbox1 = widgets.HBox([left_box, right_box])
48 hbox2 = widgets.HBox([button_metadata, button_head10,button_tail10,button_rand10])
49 display(hbox1,hbox2,output)

```



### 4. The drop-down box shows the dataset storage in Firebase. “Actor”, “film1”, and “test” is test file that I upload in previous interface:

```

45 left_box = widgets.VBox([dropdown_ds])
46 right_box = widgets.VBox([button_exploring])
47 hbox1 = widgets.HBox([left_box, right_box])
48 hbox2 = widgets.HBox([button_metadata, button_head10,button_tail10,button_rand10])
49 display(hbox1,hbox2,output)

```

Dataset:  Spotify

actor	film1	Last 10 lines	Random 10 lines
Metadata	test		

```

32]: 1 def on_button_ml_model(b):
2     with output2:
3         output2.clear_output()
4         if b.description == 'Upload':
5             rec_list = sample_recommendation_user(model = model,

```

5. When click “View Dataset”, the dataset name and description will be show in the output:

Dataset: Spotify  View Dataset

Metadata First 10 lines Last 10 lines Random 10 lines

Dataset: Spotify  
Description: Spotify music information such as artist, genre of music, track name, popularity, etc.

---

6. When click “Metadata”, the metadata and features will be show as follow:

Dataset: Spotify  View Dataset

Metadata First 10 lines Last 10 lines Random 10 lines

Description	Spotify music information such as artist, genr...
Features	[genre, artist_name, track_name, track_id, pop...
File_size	32.15MB
File_type	.csv
Number of Rows	120265
Number of columns	19
Number of features	19
 Features are:	
- genre	
- artist_name	
- track_name	
- track_id	
- popularity	
- acousticness	
- danceability	
- duration_ms	
- energy	
- instrumentalness	
- key	
- liveness	
- loudness	
- mode	
- speechiness	
- tempo	
- time_signature	
- valence	
- artist_id	

7. When click on “First 10 lines” button, user can preview the dataset. “Last 10 lines” and “Random 10 lines” are similar:

Dataset: Spotify  View Dataset

Metadata First 10 lines Last 10 lines Random 10 lines

```

acousticness artist_id      artist_name danceability duration_ms   energy \
0          0.611           0 Henri Salvador       0.389     99373    0.91
1          0.703           0 Henri Salvador       0.24      152427   0.326
2          0.749           0 Henri Salvador       0.578     160627   0.0948
3          0.381           0 Henri Salvador       0.451     194360   0.491
4          0.689           0 Henri Salvador       0.704     161773   0.804
...
...
...
...
99995     0.249           930 The Cat Empire      0.672     313053   0.603
99996     0.189           930 The Cat Empire      0.524     225467   0.765
99997     0.019           930 The Cat Empire      0.419     195547   0.575
99998     0.249           930 The Cat Empire      0.672     313053   0.603
99999     0.00727         930 The Cat Empire      0.613     258573   0.762

genre instrumentalness key liveness
0 Movie          0 C#  0.346
1 Movie          0 C#  0.0985
2 Movie          0 C#  0.107
3 Movie          0 D   0.152
4 Movie          0.0422 C   0.18
...
...
...
...
99995 Jazz        0.00277 D   0.127
99996 Ska         0 F   0.153
99997 Ska        4.83e-06 A#  0.158
99998 Ska        0.00277 D   0.127
99999 Ska         0 B   0.188

```

[100000 rows x 10 columns]

### VIII. Machine Learning Method - LightFM:

- LightFM is a Python implementation of a number of popular recommendation algorithms for both implicit and explicit feedback.
- The model learns embeddings (latent representations in a high-dimensional space) for users and items in a way that encodes user preferences over items. When multiplied together, these representations produce scores for every item for a given user; items scored highly are more likely to be interesting to the user.
- When user give a genre of music, the model will give recommended artists base on the genre, artists, and popularity. The model will give each music record score, and sort the score list until the model get the correct one. The top n of list will be output as results.
- Code:
  1. Get data from Firebase:

```
In [11]: 1 url_model = "https://final-project-551-8f1e5-default-rtdb.firebaseio.com/Spotify.json"
2 data_model = requests.get(url_model)
3 reply_model = data_model.json()
4 del reply_model['Metadata']
5 spotify_firebase = pd.DataFrame(reply_model)
6 spotify_firebase = spotify_firebase.T
```

2. Create\_interaction\_matrix(): Group the dataset by genre, artist\_id, and rating as popularity. I set the parameters for user\_col = “genre”, item\_col = “artist\_id”, and rating\_col = “popularity”.

Create\_user\_dict(); create\_item\_dict(): Generate the dictionary by “genre” and “artist\_id” separately.

runMF(): Main model of lightFM, run lightFM model and fit it on my data.

Parameters:

- n\_components – the dimensionality of the feature latent embeddings.
- loss – one of ('logistic', 'bpr', 'warp', 'warp-kos'): the loss function.  
WARP: Weighted Approximate-Rank Pairwise loss. Maximises the rank of positive examples by repeatedly sampling negative examples until rank violating one is found.
- k – for k-OS training, the k-th positive example will be selected from the n positive examples sampled for every genre.

```
1 def create_interaction_matrix(df,user_col, item_col, rating_col, norm= False, threshold = None):
2     interactions = df.groupby([user_col, item_col])[rating_col].sum().unstack().reset_index().fillna(0) \
3         .set_index(user_col)
4     if norm:
5         interactions = interactions.applymap(lambda x: 1 if x > threshold else 0)
6     return interactions
7
8 def create_user_dict(interactions):
9     user_id = list(interactions.index)
10    user_dict = {}
11    counter = 0
12    for i in user_id:
13        user_dict[i] = counter
14        counter += 1
15    return user_dict
16
17 def create_item_dict(df,id_col,name_col):
18    item_dict = {}
19    for i in range(df.shape[0]):
20        item_dict[(df.loc[i,id_col])] = df.loc[i,name_col]
21    return item_dict
22
23 def runMF(interactions, n_components=30, loss='warp', k=15, epoch=30,n_jobs = 4):
24     model = LightFM(no_components= n_components, loss=loss,k=k)
25     model.fit(x,epochs=epoch,num_threads = n_jobs)
26     return model
```

```

def sample_recommendation_user(model, interactions, actor_id, user_dict,
                               item_dict, threshold = 0, nrec_items = 10, show = True):
    n_users, n_items = interactions.shape
    user_x = user_dict[actor_id]
    scores = pd.Series(model.predict(user_x, np.arange(n_items)))
    scores.index = interactions.columns
    scores = list(pd.Series(scores.sort_values(ascending=False)).index)

    known_items = list(pd.Series(interactions.loc[actor_id, :][interactions.loc[actor_id, :] > threshold].index) \
                       .sort_values(ascending=False))

    scores = [x for x in scores if x not in known_items]
    return_score_list = scores[0:nrec_items]
    known_items = list(pd.Series(known_items).apply(lambda x: item_dict[x]))
    scores = list(pd.Series(return_score_list).apply(lambda x: item_dict[x]))
    if show == True:
        print("\n Recommended:")
        counter = 1
        for i in scores:
            print(str(counter) + ' - ' + i)
            counter+=1
    return return_score_list

```

### 3. Run all model by set up the parameters:

```

1 interactions = create_interaction_matrix(df = spotify_firebase, user_col = "genre", item_col = 'artist_id',
                                           rating_col = 'popularity', norm= False, threshold = None)
2
3 interactions

```

	artist_id	0	1	2	3	4	5	6	7	8	9	...	1144	1145	1146	1147	1148	1149	1150	1151	1152
	genre																				
Alternative	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
Anime	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
Blues	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
Children's Music	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
Children's Music	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
Classical	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
Comedy	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
Country	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
Dance	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
Electronic	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
Folk	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
Hip-Hop	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
Indie	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
Jazz	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
Movie	2205.0	266.0	713.0	9730.0	165.0	1495.0	1.0	601.0	996.0	1123.0	...	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
Opera	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
Pop	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0

```

1 user_dict = create_user_dict(interactions=interactions)
2 artists_dict = create_item_dict(df = df_artist, id_col = 'artist_id', name_col = 'artist_name')
3 x = sparse.csr_matrix(interactions.values)
4 train, test = lightfm.cross_validation.random_train_test_split(x, test_percentage=0.2, random_state=None)
5

```

```

1 %time
2 model = runMF(interactions = train,
                 n_components = 30,
                 loss = 'warp',
                 k = 15,
                 epoch = 30,
                 n_jobs = 4)

```

CPU times: user 2 µs, sys: 1 µs, total: 3 µs  
Wall time: 5.96 µs

```

1 train_auc = auc_score(model, train, num_threads=4).mean()
2 print('Train AUC: %s' % train_auc)
3 test_auc = auc_score(model, test, train_interactions=train, num_threads=4).mean()
4 print('Test AUC: %s' % test_auc)
5 train_precision = precision_at_k(model, train, k=10).mean()
6 test_precision = precision_at_k(model, test, k=10, train_interactions=train).mean()
7 print('train Precision %.2f, test Precision %.2f.' % (train_precision, test_precision))

```

Train AUC: 0.9879274  
Test AUC: 0.99997014  
train Precision 0.78, test Precision 0.92.

#### 4. Interface that let user type in the genre of music:

```

1 def on_button_ml_model(b):
2     with output2:
3         output2.clear_output()
4         if b.description == 'Upload':
5             rec_list = sample_recommendation_user(model = model,
6                                                 interactions = interactions,
7                                                 actor_id = text_song_type.value,
8                                                 user_dict = user_dict,
9                                                 item_dict = artists_dict,
10                                                threshold = 0,
11                                                nrec_items = 10,
12                                                show = True)
13
14
15
16
17
18
19
20
21
22
23

```

Type:

#### 5. The model output the recommended artists: There are two examples:

Type:

Recommended:  
1- El Alfa  
2- Keyshia Cole  
3- Reik  
4- Papa Roach  
5- J Balvin  
6- Godsmack  
7- Creed  
8- The Lacs  
9- Ari Lennox  
10- Akon

Type:

Recommended:  
1- Old Town School of Folk Music  
2- Tom's Music Box  
3- Yann Tiersen  
4- Jóhann Jóhannsson  
5- Umberto Giordano  
6- Junkie XL  
7- Gaetano Donizetti  
8- Shigeru Umebayashi  
9- Masakatsu Takagi  
10- Dora The Explorer

## IX. Learning Experiences

What are the lessons, challenges, experiences, and skills gained from the project?

After this project, I fully understand how to use the cloud databases and

what it can do for us to make the value. I learned how to set up the user interface, this is really helpful because the programs aim to help the people who do not understand the coding. I also learned how to structure a whole project. The logical relationship between modules are really important for a project. Need to keep mind very clearly.