

MACHINE LEARNING - SCREENSHOTS OF RESULTS

Student ID: 700734361

Student Name: Kiran Kumar Kongari

Subject Code: CS5710

CRN: 12675

In []:

```
# Importing our libraries
# ---

import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
%matplotlib inline

# to display all columns
pd.set_option('display.max_columns', None)

# to display the entire contents of a cell
pd.set_option('display.max_colwidth', None)

/usr/local/lib/python3.6/dist-packages/statsmodels/tools/_testing.py:19: FutureWarning: pandas.util.testing is deprecated. Use the functions in the public API at pandas.testing instead.
import pandas.util.testing as tm
```

In []:

```
# Dataset url = https://bit.ly/SpotifySongsDS
# ---
df = pd.read_csv('https://bit.ly/SpotifySongsDS')
```

In []:

```
# Checking the first 5 rows of data
# ---
df.head()
```

Out[3]:

	track_id	track_name	track_artist	track_popularity	track_album_id	track_album_name	track_album_release_date	playlist
0	I Don't Care (with Justin Bieber) - Loud Luxury Remix	Ed Sheeran	66	2oCs0DGtRo98Gh5ZSi2Cx	I Don't Care (with Justin Bieber) [Loud Luxury Remix]	2019-06-14	Pop	

In []:

```
df = pd.read_csv('https://bit.ly/SpotifySongsDS')
```

In []:

```
# Checking the first 5 rows of data
# ---
df.head()
```

Out[3]:

	track_id	track_name	track_artist	track_popularity	track_album_id	track_album_name	track_album_release_date	playlist
0	I Don't Care (with Justin Bieber) - Loud Luxury Remix	Ed Sheeran	66	2oCs0DGtRo98Gh5ZSi2Cx	I Don't Care (with Justin Bieber) [Loud Luxury Remix]	2019-06-14	Pop	
1	Memories - Dillon Francis Remix	Maroon 5	67	63rPSO264uRjW1X5E6cWv6	Memories (Dillon Francis Remix)	2019-12-13	Pop	
2	All the Time (Don Diablo Remix)	Zara Larsson	70	1HoSmj2eLcsrR0vE9gThr4	All the Time (Don Diablo Remix)	2019-07-05	Pop	
3	Call You Mine - The Chainsmokers	Keanu Silva	60	1ngYsOeffyKKuGOVchbsk6	Call You Mine - The Remixes	2019-07-19	Pop	
4	Someone You Loved (Future Humans Remix)	Lewis Capaldi	69	7m7vv9wlQ4i0LFuJlE2zsQ	Someone You Loved (Future Humans Remix)	2019-03-05	Pop	

In []:

```
# Checking the last 5 rows of data
# ---
df.tail()
```

Out[4]:

	track_id	track_name	track_artist	track_popularity	track_album_id	track_album_name	track_album_release_date	playlist
0	I Don't Care (with Justin Bieber) - Loud Luxury Remix	Ed Sheeran	66	2oCs0DGtRo98Gh5ZSi2Cx	I Don't Care (with Justin Bieber) [Loud Luxury Remix]	2019-06-14	Pop	
1	Memories - Dillon Francis Remix	Maroon 5	67	63rPSO264uRjW1X5E6cWv6	Memories (Dillon Francis Remix)	2019-12-13	Pop	
2	All the Time (Don Diablo Remix)	Zara Larsson	70	1HoSmj2eLcsrR0vE9gThr4	All the Time (Don Diablo Remix)	2019-07-05	Pop	
3	Call You Mine - The Remixes	Keanu Silva	60	1ngYsOeffyKKuGOVchbsk6	Call You Mine - The Remixes	2019-07-19	Pop	
4	Someone You Loved (Future Humans Remix)	Lewis Capaldi	69	7m7vv9wlQ4i0LFuJlE2zsQ	Someone You Loved (Future Humans Remix)	2019-03-05	Pop	

In []: # Sample 10 rows of data
df.sample(10)

	track_id	track_name	track_artist	track_popularity	track_album_id	track_album_name	track_album_release_date
32832	29zWqhca3zt5NsckZqDf6c	Typhoon - Original Mix	Julian Calor	27	0X3mUOm6MhxR7PzxG95rAo	Typhoon/Storm	2014-03-03
31686	5CMzNnyno6YvgjgYGpyVGH	Ghetto for You	Brame & Hamo	37	5LHBYza65k4scoy1ltP0	The Parish Rumors EP	2015-04-06
5584	6keactpF0I9b9NzwWDDb9	Forward Backward	Max Keeble	40	24SZva5P5KIM7XHINk4cv	Titled Gully	2019-12-10
25394	59zN1921ag2pVv80Oz3sWz	Love Riddim	Rotimi	62	0FtaeZat9X3t0cDiGzrWM	Love Riddim	2019-04-18
4622	78fHrTXB1WRS7Ub24V	Lucky	Chelsea Cutler	73	6esPAZhJXKGhByoXAHSEG	Lucky	2019-11-01
13218	7cdnq45E9aP2XDSHg5vd7	Cherry Bomb	The Runaways	65	5DVNCzpvDrSEIfIU7hm8ey	The Runaways	1976
25538	7roq9n77lul6o1Dg14Kg	It's Love	Jill Scott	34	1X26qkSSJHmYl6kUhkPN	Who Is Jill Scott? (Words And Sounds Volume 1)	2000-07-18
10864	2JfrRptu8DCHaykjUhwD	Outra Vez	Mäolee	51	6OS4M4H2eS2ZYjl0yriIu	Outra Vez	2019-09-26
31104	2thYAVPCmf4QxN7R69nPu	Wasted - Ummet Ozcan Remix	Tiesto	39	7HVxqv7UYQqpFyBvVmRL	Wasted (Remixes)	2014-06-06
5257	2aiAzLlkHcdJOnkLipgHUE	Mortel	Fishbach	33	5PFwIT6HoxgzsRsRQEduSw	A ta merci (deluxe)	2018-02-02
5161	4Ea02XR7gVaHQ8DW7nu1Y	Paris	Magic Man	53	1UBdVrk0aaya4VX1r0t3Un	Before the Waves	2014-07-07

In []: # Checking number of rows and columns
df.shape

In []: # Checking number of rows and columns
df.shape

Out[5]: (32833, 23)

In []: # Checking datatypes

df.dtypes

track_id	object
track_name	object
track_artist	object
track_popularity	int64
track_album_id	object
track_album_name	object
track_album_release_date	object
playlist_name	object
playlist_id	object
playlist_genre	object
playlist_subgenre	object
danceability	float64
energy	float64
key	int64
loudness	float64
mode	int64
speechiness	float64
acousticness	float64
instrumentalness	float64
liveness	float64
valence	float64
tempo	float64
duration_ms	int64
dtype:	object

In []: # Checking datatypes

df.dtypes

The screenshot shows a Jupyter Notebook window titled "Spotify_song_popularity_prediction (1) Last Checkpoint 11/10/2022 (autosaved)". The notebook interface includes a toolbar with File, Edit, View, Insert, Cell, Kernel, Widgets, Help, and a Python 3 (ipykernel) button. A status bar at the bottom shows "Not Trusted".

In the code cell, the command `df.describe()` is run, and the output is displayed as a table:

	track_popularity	danceability	energy	key	loudness	mode	speechiness	acousticness	instrumentalness	liveness
count	32833.000000	32833.000000	32833.000000	32833.000000	32833.000000	32833.000000	32833.000000	32833.000000	32833.000000	32833.000000
mean	42.477081	0.654850	0.699619	5.374471	-6.719499	0.565711	0.107068	0.175334	0.084747	0.190171
std	24.984074	0.145085	0.180910	3.611657	2.988436	0.499571	0.101314	0.219633	0.224230	0.15431
min	0.000000	0.000000	0.000175	0.000000	-46.448000	0.000000	0.000000	0.000000	0.000000	0.000000
25%	24.000000	0.563000	0.581000	2.000000	-8.171000	0.000000	0.041000	0.015100	0.000000	0.09270
50%	45.000000	0.672000	0.721000	6.000000	-6.166000	1.000000	0.062500	0.080400	0.000016	0.12700
75%	62.000000	0.761000	0.840000	9.000000	-4.645000	1.000000	0.132000	0.255000	0.004830	0.24800
max	100.000000	0.983000	1.000000	11.000000	1.275000	1.000000	0.918000	0.994000	0.996000	1.000000

From the dataset provided we are able to make the following observations:

- The data set has 32833 records and various used datatypes including object, float and integers.

3. External Data Source Validation

From our data I can confirm that the information about tracks in the dataset is valid because the data is evidenced by records from the Spotify mobile app.

4. Data Preparation

The screenshot shows a Jupyter Notebook window titled "Spotify_song_popularity_prediction (1) Last Checkpoint 11/10/2022 (autosaved)". The notebook interface includes a toolbar with File, Edit, View, Insert, Cell, Kernel, Widgets, Help, and a Python 3 (ipykernel) button. A status bar at the bottom shows "Not Trusted".

In the code cell, the command `# Checking datatypes and missing entries of all the variables` is run, and the output is displayed as a table:

	track_id	track_name	track_artist	track_popularity	track_album_id	track_album_name	track_album_release_date	playlist_name	playlist_id	playlist_genre	playlist_subgenre	danceability	energy	key	loudness	mode	speechiness	acousticness	instrumentalness	liveness	valence	tempo	duration_ms
Out[9]:	0	5	5	0	0	5	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

From our observation we observe that some columns contain some missing values but that's very negligible.

From our observation we observe that some columns contain some missing values but that's very negligible.

```
In [ ]: # missing values by percentage
missing_pct = df.isna().mean().round(4)*100
missing_pct
```

```
Out[10]: track_id      0.00
track_name     0.02
track_artist    0.02
track_popularity 0.00
track_album_id   0.00
track_album_name 0.02
track_album_release_date 0.00
playlist_name   0.00
playlist_id     0.00
playlist_genre   0.00
playlist_subgenre 0.00
danceability     0.00
energy          0.00
key             0.00
loudness        0.00
mode            0.00
speechiness     0.00
acousticness    0.00
instrumentalness 0.00
liveness         0.00
valence          0.00
tempo            0.00
duration_ms     0.00
dtype: float64
```

```
In [ ]: # Checking how many duplicate rows are there in the data
# ---
```

```
# ---
```

We observe the following from our dataset:

- There are no duplicated rows in the dataset

```
In [ ]: # Checking if any of the columns are all null
# ---
# df.isnull().all()
```

```
Out[12]: track_id      False
track_name     False
track_artist    False
track_popularity  False
track_album_id   False
track_album_name  False
track_album_release_date  False
playlist_name   False
playlist_id     False
playlist_genre   False
playlist_subgenre  False
danceability     False
energy          False
key             False
loudness        False
mode            False
speechiness     False
acousticness    False
instrumentalness  False
```

Downloads/ Spotify_song_popularity_prediction.ipynb

jupyter Spotify_song_popularity_prediction (1) Last Checkpoint: 11/10/2022 (autosaved)

File Edit View Insert Cell Kernel Widgets Help

Not Trusted | Python 3 (ipykernel) O

liveness False
valence False
tempo False
duration_ms False
dtype: bool

- There are no all null columns in the dataset.

In []: # Checking if any of the rows are all null

sum(df.isnull().all(axis = 1))

Out[13]: 0

- There are no all null rows in the dataset.

In []: df['track_id'].value_counts()

Out[14]: 78KLCZ1jbUBVqRiz2FVLVw 10
14505L36385F30LBhew4 9
3eeekarcy7kvN4y5t5ZfZltw 9
6oJ6le65835EqWmRNXuJy 8
0\$1zqH55qcw8pgymFoQ 8
..
6vFsxu9TUJjWzMtualgyT 1
7ADrwGz12XJxxlw5c9oMy 1
4HzI9Nb5FE5n9V4cvn6il 1
7lVN7XK13CHfvXLI8damb 1
5zNfetnJFcsGMwK067xlck 1
Name: track_id, Length: 28356, dtype: int64

40°F Moderate air

Search

7:09 PM 12/5/2022

Downloads/ Spotify_song_popularity_prediction.ipynb

jupyter Spotify_song_popularity_prediction (1) Last Checkpoint: 11/10/2022 (autosaved)

File Edit View Insert Cell Kernel Widgets Help

Not Trusted | Python 3 (ipykernel) O

checking unique values for our columns

In []: columns = df.columns

for col in columns:
 print("Variable:", col)
 print("Number of unique values:", df[col].nunique())
 print(df[col].unique())
 print()

Variable: track_id
Number of unique values: 28356
['6f807x0ima9a1j3Vpb7WN' '0r7CvbZTMZgbTCYdfa2P31'
'1z1Hg7VbAhDlEm0e791' ... '1mMqPP3QiyfUhvsdnwEo'
'2m69hnhfq10g61GtxUhgX' '29Zwqhc3zt5Nsck2qf6c']

Variable: track_name
Number of unique values: 23449
['I Don't Care (with Justin Bieber) - Loud Luxury Remix'
'Memories - Dillon Francis Remix' 'All the Time - Don Diablo Remix' ...
'Sweet Surrender - Radio Edit' 'Only For You - Maor Levi Remix'
'Typhoon - Original Mix']

Variable: track_artist
Number of unique values: 10692
['Ed Sheeran' 'Maroon 5' 'Zara Larsson' ...
'Ferry Corsten feat. Jenny Wahlinström' 'Tegan and Sara' 'Mat Zo']

Variable: track_popularity
Number of unique values: 101
[66 67 70 60 69 62 68 58 63 65 35 64 8 30 56 55 59 87
83 61 57 53 34 74 46 52 45 51 48 39 71 81 85 76 75 73
79 16 84 37 36 72 77 82 41 5 24 18 47 80 96 86 19 27
43 91 42 21 54 9 23 40 2 11 50 10 49 6 25 95 1 4
28 29 38 17 31 26 32 12 44 7 0 13 78 88 94 22 98 97
90 93 92 3 33 20 14 15 89 100 99]

40°F Moderate air

Search

7:09 PM 12/5/2022

Cleaning the data

Renaming columns

```
In [ ]: df.columns = df.columns.str.lower()
df.head()
```

Out[16]:

	track_id	track_name	track_artist	track_popularity	track_album_id	track_album_name	track_album_release_date	playlist
0	6f807x0ima9a1j3VPbc7VN	I Don't Care (with Justin Bieber) - Loud Luxury Remix	Ed Sheeran	68	2oCs0DGTSrO98Gh5ZSi2Cx	I Don't Care (with Justin Bieber) [Loud Luxury Remix]	2019-06-14	Pop
1	0r7CvbZTWZgbTCYdfaP31	Memories - Dillon Francis Remix	Maroon 5	67	63rPSO264uRjW1X5E6cWv6	Memories (Dillon Francis Remix)	2019-12-13	Pop
2	1z1Hg7Vb0AhDiEmnDE79l	All the Time - Don Diablo Remix	Zara Larsson	70	1HoSmj2eLcsrR0vE9gThr4	All the Time (Don Diablo Remix)	2019-07-05	Pop
3	75FpbthrwQmzHIBJLuGdC7	Call You Mine - Keanu Silva Remix	The Chainsmokers	60	1nqYsOefIyKKuGOvchbsk6	Call You Mine - The Remixes	2019-07-19	Pop
4	1e8PAfcKUYoKbxPhrHqw4x	Someone You Loved - Future Humans Remix	Lewis Capaldi	69	7m7vv9wQ4i0LFuJIe2zsQ	Someone You Loved (Future Humans Remix)	2019-03-05	Pop

Dropping rows with missing values

Dropping rows with missing values

```
In [ ]: df.dropna(inplace = True)
df.isna().sum()
```

Out[17]:

```
track_id      0
track_name     0
track_artist    0
track_popularity  0
track_album_id   0
track_album_name  0
track_album_release_date 0
playlist_name    0
playlist_id      0
playlist_genre    0
playlist_subgenre  0
danceability     0
energy          0
key             0
loudness        0
mode            0
speechiness      0
acousticness     0
instrumentalness 0
liveness         0
valence          0
tempo           0
duration_ms      0
dtype: int64
```

Dropping irrelevant variables

```
In [ ]: df.drop(columns = ['track_album_id','playlist_id','playlist_name','track_album_name','track_name'], inplace = True)
```

Downloads/ Spotify_song_popularity_prediction.ipynb

jupyter Spotify_song_popularity_prediction (1) Last Checkpoint: 11/10/2022 (autosaved)

File Edit View Insert Cell Kernel Widgets Help

In []: df.drop(columns = ['track_album_id','playlist_id','playlist_name','track_album_name','track_name'], inplace = True)
df.sample(2)

Out[18]:

	track_id	track_artist	track_popularity	track_album_release_date	playlist_genre	playlist_subgenre	danceability	energy	key	lou
21907	3WRlaWsws011vHfd9uzPJG	Future	69	2019-01-18	r&b	urban contemporary	0.734	0.683	1	-
943	1kSEroHstaZPoBmouyCbHU	Noah Neiman	31	2019-10-18	pop	dance pop	0.573	0.672	1	-

In []: df.info()

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 32828 entries, 0 to 32832
Data columns (total 18 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   track_id         32828 non-null   object 
 1   track_artist     32828 non-null   object 
 2   track_popularity 32828 non-null   int64  
 3   track_album_release_date 32828 non-null   object 
 4   playlist_genre   32828 non-null   object 
 5   playlist_subgenre 32828 non-null   object 
 6   danceability     32828 non-null   float64
 7   energy           32828 non-null   float64
 8   key              32828 non-null   int64  
 9   loudness         32828 non-null   float64
 10  mode             32828 non-null   int64  
 11  speechiness      32828 non-null   float64
 12  acousticness    32828 non-null   float64
 13  instrumentalness 32828 non-null   float64
 14  liveliness       32828 non-null   float64
 15  valence          32828 non-null   float64
 16  tempo            32828 non-null   float64
```

40°F AQI 51

Search

ENG US

7:10 PM 12/5/2022

Downloads/ Spotify_song_popularity_prediction.ipynb

jupyter Spotify_song_popularity_prediction (1) Last Checkpoint: 11/10/2022 (autosaved)

File Edit View Insert Cell Kernel Widgets Help

In []: df_clean = df.copy()
df_clean.sample(20)

Out[20]:

	track_id	track_artist	track_popularity	track_album_release_date	playlist_genre	playlist_subgenre	danceability	energy	key	lou
27009	3yTR8aMxpwaPLEsjHmv0Ew	Two Figures	28	2018-06-25	edm	electro house	0.901	0.964	4	-
29967	5TR2ewQlnmctowBoPTCuL	Licho	40	2019-12-16	edm	pop edm	0.708	0.588	5	-
25715	0z9PmxrJUVCLc3e7fESJL	Mustard	30	2016-03-14	r&b	neo soul	0.784	0.798	4	-
11576	3Fdhg0.bh3DXNtGLh5pFu	Miami Sound Machine	56	1985-08-13	rock	album rock	0.789	0.655	2	-
32314	0FNtBIL5Hcn7pLVTp6TG	Quintino	44	2016-01-01	edm	progressive electro house	0.567	0.772	5	-
20777	4T4CZ3Q4z3GG7lU720YhM	Lunay	70	2019-02-21	latin	latin hip hop	0.778	0.731	5	-
30590	2ea8kp4WKy1GdqHoYaz9Nu	Simon Field	34	2017-09-15	edm	pop edm	0.628	0.787	0	-
15076	42lgso7JdwvXnLygpKBu	Alveole	17	2019-10-04	rock	hard rock	0.521	0.924	7	-
17034	3qaYuV5XZAwI0vGE91HwS	Mike Perry	48	2016-10-21	latin	tropical	0.561	0.800	8	-
32197	1bSO9nzT5h55OrJ07BbHL2	DJ MEG	38	2019-11-29	edm	progressive electro house	0.644	0.936	9	-
13979	6Volbz0vhCyz7OIEoRYDw	Weezer	73	1994-05-10	rock	permanent wave	0.634	0.551	3	-
2170	00INx0OcTjS3MKHcb80HY	Jax Jones	14	2016-12-09	pop	post-teen pop	0.876	0.669	11	-
27178	5NK8imKaxqZlgzbSqeZ	Bob Gordon	44	2019-11-18	edm	electro house	0.800	0.832	1	-
29412	2H98dHwAMJf0bg0gcoc09N	Bassjackers	0	2015-03-02	edm	big room	0.692	0.968	1	-
24134	22sTmBjk27SLWmtsdyL14u	Mary J. Blige	2	1992-01-01	r&b	new jack swing	0.760	0.793	0	-
9997	2985JpMNKG1QT5i4RSjehz	Giaime	63	2019-11-08	rap	trap	0.758	0.804	11	-
12996	0FMhMAFPgLgSEjRROS0Vnj	ZZ Top	55	2004-06-08	rock	classic rock	0.552	0.651	4	-

40°F AQI 51

Search

ENG US

7:10 PM 12/5/2022

Downloads/ Spotify_song_popularity_prediction.ipynb

jupyter Spotify_song_popularity_prediction (1) Last Checkpoint: 11/10/2022 (autosaved)

File Edit View Insert Cell Kernel Widgets Help

Not Trusted Python 3 (ipykernel) O

In []: `# installing the pandas profiling package from
pip install https://github.com/pandas-profiling/pandas-profiling/archive/master.zip`

```
Collecting https://github.com/pandas-profiling/pandas-profiling/archive/master.zip  
  Using cached https://github.com/pandas-profiling/pandas-profiling/archive/master.zip  
Requirement already satisfied (use --upgrade to upgrade): pandas-profiling==2.9.0 from https://github.com/pandas-profiling/pandas-profiling/archive/master.zip in /usr/local/lib/python3.6/dist-packages  
Requirement already satisfied: joblib in /usr/local/lib/python3.6/dist-packages (from pandas-profiling==2.9.0) (0.16.0)  
Requirement already satisfied: scipy>=1.4.1 in /usr/local/lib/python3.6/dist-packages (from pandas-profiling==2.9.0) (1.4.1)  
Requirement already satisfied: pandas!=1.0.0,>=1.0.1,>=1.0.2,>=1.1.0,>=0.25.3 in /usr/local/lib/python3.6/dist-packages (from pandas-profiling==2.9.0) (1.0.5)  
Requirement already satisfied: matplotlib>=3.2.0 in /usr/local/lib/python3.6/dist-packages (from pandas-profiling==2.9.0) (3.2.2)  
Requirement already satisfied: confuse>=1.0.0 in /usr/local/lib/python3.6/dist-packages (from pandas-profiling==2.9.0) (1.3.0)  
Requirement already satisfied: jinja2>=2.11.1 in /usr/local/lib/python3.6/dist-packages (from pandas-profiling==2.9.0) (2.11.2)  
Requirement already satisfied: visions[type_image_path]==0.5.0 in /usr/local/lib/python3.6/dist-packages (from pandas-profiling==2.9.0) (0.5.0)  
Requirement already satisfied: numpy>=1.16.0 in /usr/local/lib/python3.6/dist-packages (from pandas-profiling==2.9.0) (1.18.5)
```

40°F AQI 51

ENG US 7:10 PM 12/5/2022

Downloads/ Spotify_song_popularity_prediction.ipynb

jupyter Spotify_song_popularity_prediction (1) Last Checkpoint: 11/10/2022 (autosaved)

File Edit View Insert Cell Kernel Widgets Help

Not Trusted Python 3 (ipykernel) O

```
Created wheel for pandas-profiling: filename=pandas-profiling-2.9.0-py2.py3-none-any.whl size=258952 sha256=949780d48f9d2029fd3892d0392b4ec649d798df9e943f74ec33783ced9b195b1  
Stored in directory: /tmp/pip-ephem-wheel-cache-ofic756a/wheels/56/c2/dd/8d945b0443c35df7df62fa9e9ae105a2d8b286302b92e  
0109  
Successfully built pandas-profiling
```

Type Markdown and LaTeX: α^2

In []: `import pandas_profiling`

```
df_clean.profile_report()  
HBox(children=(FloatProgress(value=0.0, description='Summarize dataset', max=33.0, style=ProgressStyle(description_wi...  
  
HBox(children=(FloatProgress(value=0.0, description='Generate report structure', max=1.0, style=ProgressStyle(...  
  
HBox(children=(FloatProgress(value=0.0, description='Render HTML', max=1.0, style=ProgressStyle(description_wi...
```

Out[22]:

FEATURE ENGINEERING

FOMATING YEAR OF RELEASE VARIABLE

In []: `df_clean['track_album_release_date'] = pd.to_datetime(df_clean['track_album_release_date']).dt.year`

Creating a new column : Frequency to check on whether the song appears in more than one subgenres.

40°F AQI 51

ENG US 7:10 PM 12/5/2022

In []:

```
counts = df_clean.groupby('playlist_genre')[“track_id”].value_counts().rename(“Frequency”).reset_index()
df_clean = df_clean.merge(counts)
df_clean.sample(100)
```

Out[24]:

	track_id	track_artist	track_popularity	track_album_release_date	playlist_genre	playlist_subgenre	danceability	energy	key	lou
26011	2gpEHrW0XBpeAlGORiBWr5	Kendra Morris	32	2013	r&b	neo soul	0.499	0.523	5	
5984	5fu3apTO0JuTERm53CGXf	Lil Mosey	41	2019	rap	hip hop	0.697	0.685	7	
12284	3yuJVbCLjxcoUT6BzoiWxg	Queen	31	1976	rock	album rock	0.439	0.561	2	
10710	4ohD3ipfysmcWex1AyddI	Cynthia Luz	67	2018	rap	trap	0.517	0.460	8	
26729	6AIb4vBOGWM4FJ65yNWxu	James Brown	51	1991	r&b	neo soul	0.646	0.350	8	
...
31654	05jSQ5KtlbYQThxoijnP	Rebeat	19	2014	edm	progressive electro house	0.569	0.960	5	
2921	17Fd6Yb7msBmKG0LoWfI	Mike Posner	17	2016	pop	electropop	0.663	0.713	7	
32159	3Gb01FC1a8NRiDytBcoix	Roberto Rosso	4	2017	edm	progressive electro house	0.568	0.883	5	
16660	0UnJphaDGplP59XRMaE	Juan Luis Guerra 4 40	52	2019	latin	tropical	0.720	0.878	5	
30912	6LoQHlo74tOzQ8EsLEkhgF	Selena Gomez	69	2013	edm	pop edm	0.546	0.787	7	

100 rows × 19 columns

In []:

```
df_clean[‘Frequency’].value_counts(ascending=False)
```

Out[25]:

1	28236
2	3716
3	884
4	72

Name: Frequency, dtype: int64

In []:

```
df_clean[‘track_popularity’].value_counts()
```

Out[26]:

0	2698
1	575
2	541
3	514
4	514
...	...
97	22
98	15
99	7
100	4

Name: track_popularity, Length: 101, dtype: int64

creating a new variable `Likely_popular` based on the recorded popularity of the tracks

'`Likely_popular`'= 1 if it has popularity of +50 and '`Likely_popular`' = 0 if its less than 50.

```
In [ ]: # Create a new column called Likely_popular where the value is 1 if popularity is greater than 50 and 0 if not.
df_clean['likely_popular'] = np.where(df_clean['track_popularity'] >= 50, 1, 0)
df_clean.sample(20)
```

	track_id	track_artist	track_popularity	track_album_release_date	playlist_genre	playlist_subgenre	danceability	energy	key	temp
17150	5REoTLUGvmtKUIMvdKwKboUn	Iceleak	48	2019	latin	tropical	0.722	0.724	8	116.0
9251	3YrExZjAZrchrO72vG5gK9	BONES	49	2019	rap	gangster rap	0.552	0.598	6	116.0
10500	4JvI2SVUJP3vDsH6oJDuL	Khea	51	2018	rap	trap	0.748	0.742	6	116.0
1285	0lDpsyFvt1KaxNBYWz2aaq	John Maus	54	2011	pop	dance pop	0.573	0.639	0	116.0
20642	3LUUsUqjlpMuXnNjHzDsf	Promise Circle	11	2007	latin	latin hip hop	0.750	0.911	1	116.0
2171	1RwmmiVLAIPmxAqKVifwg	Bleachers	53	2014	pop	post-teen pop	0.456	0.891	9	116.0
24320	2vAbKDdgA6cEBHU9zQiaB8g	LL Cool J	0	2009	r&b	new jack swing	0.608	0.665	1	116.0
29969	6P6YKV0Z2XSYOMhlyWpOUWu	Timo Odv	1	2015	edm	pop edm	0.733	0.696	10	116.0
9674	67XyC1GdSkKb9sbRRkUDXS	NLE Choppa	80	2019	rap	trap	0.894	0.511	2	116.0
15440	39AKmo6kwifRcyOSzbSu00T	Bruce Springsteen	1	2009	rock	hard rock	0.529	0.877	11	116.0
24917	1krKp0OXeCH6Si5SXIBtu5	Janet Jackson	48	1986	r&b	new jack swing	0.661	0.711	0	116.0
6663	2XXvYN4zZdmxJMiizuWPB	TiMT	2	2018	rap	hip hop	0.567	0.204	8	116.0

encoding categorical variables

```
In [ ]: from sklearn.preprocessing import LabelEncoder

le_track_artist = LabelEncoder().fit(df_clean['track_artist'].unique())
```

Downloads/ Spotify_song_popularity_prediction.ipynb

jupyter Spotify_song_popularity_prediction (1) Last Checkpoint: 11/10/2022 (autosaved)

File Edit View Insert Cell Kernel Widgets Help

Not Trusted | Python 3 (ipykernel) O

In []: `#dropping 'playlist_genre' and 'playlist_subgenre' they do have certain single tracks repeated more than once in observation`

df_clean.drop(columns=['track_album_release_date'], inplace=True)

df_clean.sample(10)

Out[29]:

	track_id	track_artist	track_popularity	playlist_genre	playlist_subgenre	danceability	energy	key	loudness	mode	speechiness	acousticness	ins
3399	12617	5315	30	2	5	0.709	0.434	11	-12.089	1	0.0668	0.5190	
2421	20262	3148	19	2	17	0.768	0.567	2	-7.046	0	0.0312	0.6600	
6528	11193	33	63	4	8	0.848	0.699	5	-5.425	0	0.1570	0.0401	
21957	23082	1135	76	3	9	0.868	0.646	1	-4.674	0	0.2880	0.3030	
19375	25895	1709	62	1	19	0.583	0.823	4	-4.663	0	0.2050	0.1250	
11643	17966	6363	0	5	0	0.340	0.943	9	-8.970	1	0.0851	0.0285	
16697	9801	7077	20	1	22	0.812	0.783	5	-5.598	0	0.0528	0.0478	
30367	24079	360	47	0	16	0.701	0.942	5	-2.253	1	0.2410	0.1590	
21503	18036	8203	47	3	23	0.565	0.905	6	-4.396	1	0.1490	0.1080	
26263	5884	53	45	3	13	0.513	0.379	1	-13.587	0	0.1710	0.2590	

In []:

In []: `# visualizing the distribution of outliers -`

df_clean.to_csv('spotify.csv')

df_clean=pd.read_csv('spotify.csv')

DATA MODELING

In []: `# split into features (X) and target (Y)`

Downloads/ Spotify_song_popularity_prediction.ipynb

jupyter Spotify_song_popularity_prediction (1) Last Checkpoint: 11/10/2022 (autosaved)

File Edit View Insert Cell Kernel Widgets Help

Not Trusted | Python 3 (ipykernel) O

In []: `# visualizing the distribution of outliers -`

df_clean.to_csv('spotify.csv')

df_clean=pd.read_csv('spotify.csv')

DATA MODELING

In []: `# split into features (X) and target (Y)`

X = df_clean.drop(['likely_popular','track_id','track_popularity'], axis=1)

y = df_clean.likely_popular

print(X.shape)

print(y.shape)

(32828, 18)

(32828,)

In []: `# splitting into 75-25 training and test sets`

from sklearn.model_selection import train_test_split as tts

X_train, X_test, y_train, y_test = tts(X, y, test_size = 0.3, random_state = 0)

In []: `# scaling our features`

from sklearn.preprocessing import StandardScaler

sc = StandardScaler()

X_train = sc.fit_transform(X_train)

X_test= sc.transform(X_test)

In []: `#fetching regressors`

from sklearn.linear_model import LogisticRegression

from sklearn.tree import DecisionTreeClassifier

```
random_forest_classifier.fit(X_train, y_train)
random_forest_y_prediction = random_forest_classifier.predict(X_test)

# Ada Boosting
from sklearn.ensemble import AdaBoostClassifier
ada_boost_classifier = AdaBoostClassifier()
ada_boost_classifier.fit(X_train, y_train)
ada_boost_y_prediction = ada_boost_classifier.predict(X_test)

# Gradient Boosting
from sklearn.ensemble import GradientBoostingClassifier
g_boost_classifier = GradientBoostingClassifier()
g_boost_classifier.fit(X_train, y_train)
g_boost_y_prediction = g_boost_classifier.predict(X_test)

# XG Boosting
from xgboost import XGBClassifier
xg_boost_classifier = XGBClassifier()
xg_boost_classifier.fit(X_train, y_train)
xg_boost_y_prediction = xg_boost_classifier.predict(X_test)

# Logistic Regression
log = LogisticRegression(solver = 'lbfgs', max_iter = 200)
log.fit(X_train, y_train)
log_pred = log.predict(X_test)

# Evaluation of the model
print('Logistic Regression Performance: ')
print('RMSE: ', np.sqrt(metrics.mean_squared_error(y_test, log_pred)))
print('R2: ', metrics.r2_score(y_test, log_pred))

Logistic Regression Performance:
RMSE: 0.592965255615975
R2: -0.4455959385923156
```

```
# Support Vector Machine
svm = SVC(kernel = 'rbf')
svm.fit(X_train, y_train)
svm_pred = svm.predict(X_test)

# Evaluation of the model
print('Logistic Regression Performance: ')
print('RMSE: ', np.sqrt(metrics.mean_squared_error(y_test, svm_pred)))
print('R2: ', metrics.r2_score(y_test, svm_pred))

Logistic Regression Performance:
RMSE: 0.5798946410340153
R2: -0.38256244545704576

# XG Boosting
from xgboost import XGBClassifier
xgb = XGBClassifier()
xgb.fit(X_train, y_train)
xgb_pred = xgb.predict(X_test)

# Evaluation of the model
print('Logistic Regression Performance: ')
print('RMSE: ', np.sqrt(metrics.mean_squared_error(y_test, xgb_pred)))
print('R2: ', metrics.r2_score(y_test, xgb_pred))

Logistic Regression Performance:
RMSE: 0.5552068489685973
R2: -0.2673489083356253

# Random Forest
rf = RandomForestClassifier()
rf.fit(X_train, y_train)
```

Downloads/ Spotify_song_popularity_prediction.ipynb

jupyter Spotify_song_popularity_prediction (1) Last Checkpoint: 11/10/2022 (autosaved)

File Edit View Insert Cell Kernel Widgets Help

xgb_pred = xgb.predict(X_test)

```
# Evaluation of the model
print('Logistic Regression Performance: ')
print('RMSE: ', np.sqrt(metrics.mean_squared_error(y_test, log_pred)))
print('R2: ', metrics.r2_score(y_test, log_pred))

Logistic Regression Performance:
RMSE: 0.5552068489685973
R2: -0.2673489083356253
```

In []:

```
# Random Forest
rf = RandomForestClassifier()
rf.fit(X_train, y_train)
rf_pred = rf.predict(X_test)

# Evaluation of the model
print('Random Forest Performance: ')
print('RMSE: ', np.sqrt(metrics.mean_squared_error(y_test, rf_pred)))
print('R2: ', metrics.r2_score(y_test, rf_pred))

Random Forest Performance:
RMSE: 0.5177330562311706
R2: -0.10204252898750021
```

In []:

```
# evaluating the classification reports and confusion matrices of each classifier
from sklearn.metrics import classification_report, confusion_matrix

# Logistic Regression
print("Logistic Regression Confusion Matrix:\n", confusion_matrix(y_test, log_pred))
print("Classification Report:\n", classification_report(y_test, log_pred))
print()

# Support Vector Machine
```

40°F Mostly cloudy

Search

Logout

Python 3 (ipykernel) O

Not Trusted

File Edit View Insert Cell Kernel Widgets Help

7:11 PM 12/5/2022

Downloads/ Spotify_song_popularity_prediction.ipynb

jupyter Spotify_song_popularity_prediction (1) Last Checkpoint: 11/10/2022 (autosaved)

File Edit View Insert Cell Kernel Widgets Help

```
# evaluating the classification reports and confusion matrices of each classifier
from sklearn.metrics import classification_report, confusion_matrix

# Logistic Regression
print("Logistic Regression Confusion Matrix:\n", confusion_matrix(y_test, log_pred))
print("Classification Report:\n", classification_report(y_test, log_pred))
print()

# Support Vector Machine
print("Support Vector Machine Confusion Matrix:\n", confusion_matrix(y_test, svm_pred))
print("Classification Report:\n", classification_report(y_test, svm_pred))
print()

# Random Forest
print("Random Forest Confusion Matrix:\n", confusion_matrix(y_test, rf_pred))
print("Classification Report:\n", classification_report(y_test, rf_pred))
print()

# XG Boosting
print("XG Boosting Confusion Matrix:\n", confusion_matrix(y_test, xgb_pred))
print("Classification Report:\n", classification_report(y_test, xgb_pred))

# classification report for K-Nearest Neighbors Classifier
print("K-Nearest Neighbors classification report:", confusion_matrix(y_test, knn_y_prediction))
print("Classification Report:\n", classification_report(y_test, knn_y_prediction))
print()
#bagging
print("bagging classification report:", confusion_matrix(y_test, bagging_y_prediction))
print("Classification Report:\n", classification_report(y_test, bagging_y_prediction))
print()

# ada boosting
print("ada_boost classification report:", confusion_matrix(y_test, ada_boost_y_prediction))
print("Classification Report:\n", classification_report(y_test, ada_boost_y_prediction))
```

40°F Mostly cloudy

Search

Logout

Python 3 (ipykernel) O

Not Trusted

File Edit View Insert Cell Kernel Widgets Help

7:11 PM 12/5/2022

Downloads/ Spotify_song_popularity_prediction.ipynb

jupyter Spotify_song_popularity_prediction (1) Last Checkpoint: 11/10/2022 (autosaved)

File Edit View Insert Cell Kernel Widgets Help

Not Trusted | Python 3 (ipykernel) O

```
print("bagging classification report:", confusion_matrix(y_test, bagging_y_prediction))
print("Classification Report:\n", classification_report(y_test, bagging_y_prediction))
print()

# ada boosting
print("ada_boost classification report:", confusion_matrix(y_test, ada_boost_y_prediction))
print("Classification Report:\n", classification_report(y_test, ada_boost_y_prediction))
print()

Logistic Regression Confusion Matrix:
[[1788 2326]
 [1137 4598]]
Classification Report:
precision recall f1-score support
0 0.61 0.43 0.51 4114
1 0.66 0.80 0.73 5735

accuracy macro avg weighted avg
0.64 0.62 0.64 9849

Support Vector Machine Confusion Matrix:
[[1672 2442]
 [ 870 4865]]
Classification Report:
precision recall f1-score support
0 0.66 0.41 0.50 4114
1 0.67 0.85 0.75 5735

accuracy macro avg weighted avg
0.66 0.63 0.64 9849
```

40°F Mostly cloudy

Search

ENG US

7:11 PM 12/5/2022

Downloads/ Spotify_song_popularity_prediction.ipynb

jupyter Spotify_song_popularity_prediction (1) Last Checkpoint: 11/10/2022 (autosaved)

File Edit View Insert Cell Kernel Widgets Help

Not Trusted | Python 3 (ipykernel) O

```
Random Forest Confusion Matrix:
[[2585 1529]
 [1111 4624]]
Classification Report:
precision recall f1-score support
0 0.70 0.63 0.66 4114
1 0.75 0.81 0.78 5735

accuracy macro avg weighted avg
0.73 0.72 0.73 9849

XG Boosting Confusion Matrix:
[[2128 1986]
 [1050 4685]]
Classification Report:
precision recall f1-score support
0 0.67 0.52 0.58 4114
1 0.70 0.82 0.76 5735

accuracy macro avg weighted avg
0.69 0.67 0.69 9849

K-Nearest Neighbors classification report: [[2076 2038]
 [1628 4107]]
Classification Report:
precision recall f1-score support
0 0.56 0.50 0.53 4114
1 0.67 0.72 0.69 5735

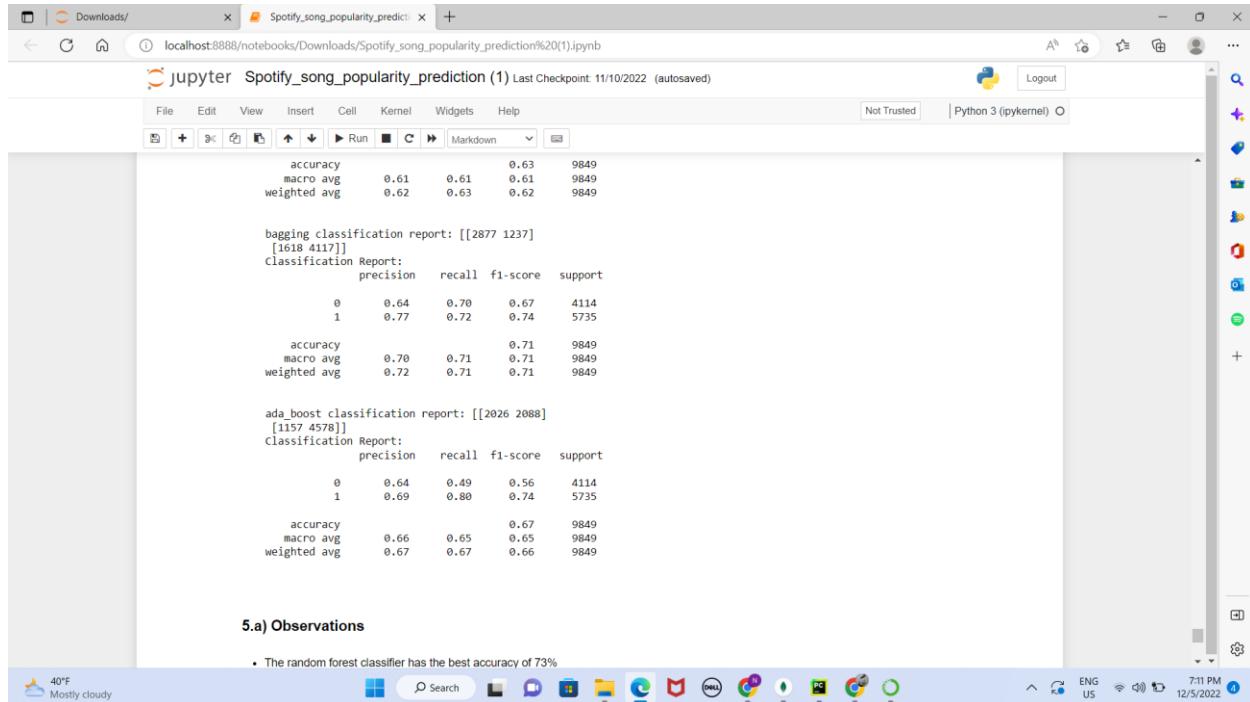
accuracy macro avg weighted avg
0.63 0.61 0.63 9849
```

40°F Mostly cloudy

Search

ENG US

7:11 PM 12/5/2022



accuracy 0.63 9849
macro avg 0.61 0.61 0.61 9849
weighted avg 0.62 0.63 0.62 9849

bagging classification report: [[2877 1237]
[1618 4117]]
Classification Report:
precision recall f1-score support
0 0.64 0.70 0.67 4114
1 0.77 0.72 0.74 5735

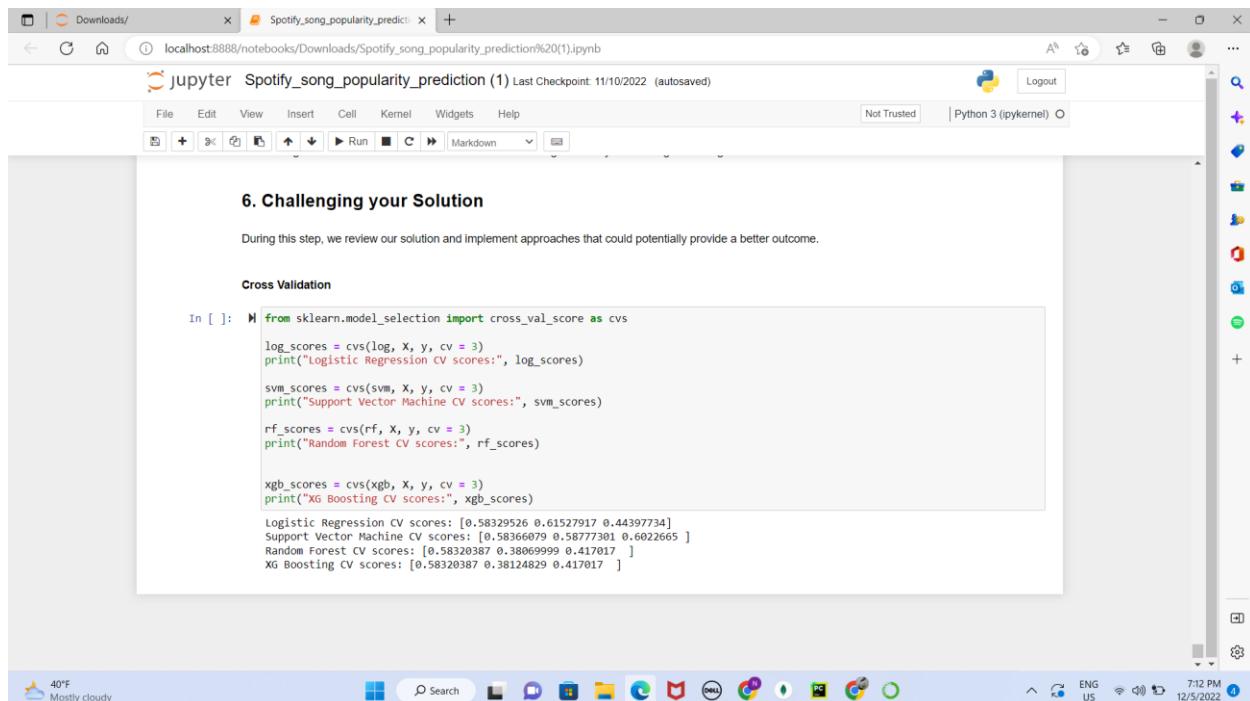
accuracy 0.71 9849
macro avg 0.70 0.71 0.71 9849
weighted avg 0.72 0.71 0.71 9849

ada_boost classification report: [[2026 2088]
[1157 4578]]
Classification Report:
precision recall f1-score support
0 0.64 0.49 0.56 4114
1 0.69 0.80 0.74 5735

accuracy 0.67 9849
macro avg 0.66 0.65 0.65 9849
weighted avg 0.67 0.67 0.66 9849

5.a) Observations

- The random forest classifier has the best accuracy of 73%



Cross Validation

```
In [ ]: from sklearn.model_selection import cross_val_score as cvs
log_scores = cvs(log, X, y, cv = 3)
print("Logistic Regression CV scores:", log_scores)

svm_scores = cvs(svm, X, y, cv = 3)
print("Support Vector Machine CV scores:", svm_scores)

rf_scores = cvs(rf, X, y, cv = 3)
print("Random Forest CV scores:", rf_scores)

xgb_scores = cvs(xgb, X, y, cv = 3)
print("XG Boosting CV scores:", xgb_scores)
```

Logistic Regression CV scores: [0.58329526 0.61527917 0.44397734]
Support Vector Machine CV scores: [0.58366079 0.58777301 0.6022665]
Random Forest CV scores: [0.58320387 0.38069999 0.417017]
XG Boosting CV scores: [0.58320387 0.38124829 0.417017]

GITHUB Link:

https://github.com/KiranKumarKongari/ML-Final_Project