

BillBoard Hits Prediction-Project_Workflow_Details_EDA

Libraries Import

```
In [1]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

from IPython.display import IFrame

import warnings
warnings.filterwarnings('ignore')

from sklearn.preprocessing import StandardScaler

pd.options.plotting.backend = "matplotlib"
#pd.options.plotting.backend = "fplot"

import scipy.stats as stats
from statsmodels.stats import weights as stats
```

Initial Data Analysis and Visualization

Bill Board Data

```
In [2]: df_bb = pd.read_excel('C:/Users/Anjum/Downloads/Technocolabs Materials/Predicting-Billboard-Hits-Using-Spotify-
df_bb.head(2)
```

	Track	Artist	SpotifyID	danceability	energy	key	mode	speechiness	acousticness	instrumentalness	liveness	valence
0	Lucid Dreams	WRLD	285p8ltuF7W6TeW8hRDR	0.511	0.566	6	0	0.20	0.349	0.0	0.340	0.218
1	Better Now	Post Malone	7d6xSM1jpdT8t8C8u5tK	0.680	0.578	10	1	0.04	0.331	0.0	0.135	0.341

```
In [4]: df_bb.shape
Out[4]: (9329, 15)
```

```
In [5]: df_bb.shape
Out[5]: (9329, 15)
```

Dropping duplicate data based on SpotifyID

```
In [6]: df_bb.drop_duplicates(subset='SpotifyID',keep=False,inplace=True)
In [7]: df_bb.shape
Out[7]: (9258, 15)
```

MSD Data

```
In [8]: df_msd = pd.read_excel('C:/Users/Anjum/Downloads/Technocolabs Materials/Predicting-Billboard-Hits-Using-Spotify-
In [9]: df_msd.head(2)
```

		0	1	2	3	4	5	danceability	energy	key	loudness	mode	speechiness	acousticness	instrumentalness	liveness	valence	
0	Sonora Santanera	Numero 1 De La Sonora Santanera	Amor De Cabaret	De 2016	TRAAADZ128F93148C2E	SleyXbf7CVCVQ8DqWfKfANP		0.699	0.519	8	-6.422	1		0.099	0.519	8	-6.422	1
1	JennyAnyKind	I Need You	Young Boy Blues	2000	TRAAAVO128F93133DA	SuN6vY2NLG6GKoywHKZA		0.458	0.570	5	-9.159	0		0.458	0.570	5	-9.159	0

```
In [10]: df_msd.shape
Out[10]: (5603, 18)
```

Renaming columns of MSD Dataframe

```
In [11]: df_msd.rename(columns={0:'Artist',1:'Album',2:'Track',3:'Year',4:'TrackID',5:'SpotifyID'},inplace=True)
In [12]: df_msd.head()
```

	Artist	Album	Track	Year	TrackID	SpotifyID	danceability	energy	key	loudness	mode	speechiness	acousticness	instrumentalness	liveness	valence
0	Sonora Santanera	Numero 1 De La Sonora Santanera	Amor De Cabaret	2016	TRAAADZ128F93148C2E	SleyXbf7CVCVQ8DqWfKfANP	0.699	0.519	8	-6.422	1	0.099	0.519	8	-6.422	1
1	JennyAnyKind	I Need You	Young Boy Blues	2000	TRAAAVO128F93133DA	SuN6vY2NLG6GKoywHKZA	0.458	0.570	5	-9.159	0	0.458	0.570	5	-9.159	0
2		Casual Fear	Idon't Mean To	1994	TRAAAAM128F429D538	01TR6aAKA2c1320gnC0u	0.751	0.549	6	-10.508	0					
3	Jeff And Shri Easter	Ordinary Day	The Moon And I (Ordinary Day Album Version)	2006	TRAAAO128F1481E7F	76Rvc8BfY5HqQzQqJh9AN	0.456	0.472	5	-8.328	1					
4	Tweetfriendly	Cin & Phonic	Drop of Rain	2003	TRAAAAC128E0786D96	6dfpQZQZ04z3Kc6N8V93	0.498	0.300	7	-10.716	1					

Dropping duplicate data based on SpotifyID

```
In [13]: df_msd.drop(columns=['Album','Year','TrackID'],axis=1,inplace=True)
In [14]: df_msd.shape
Out[14]: (5603, 15)
```

```
In [15]: df_msd.drop_duplicates(subset='SpotifyID',keep=False,inplace=True)
In [16]: df_msd.shape
Out[16]: (5522, 15)
```

```
In [17]: df_bb.shape,df_msd.shape
Out[17]: (9258, 15), (5522, 15)
```

```
In [18]: df_bb.head(2)
Out[18]:
```

	Track	Artist	SpotifyID	danceability	energy	key	mode	speechiness	acousticness	instrumentalness	liveness	valence
0	Lucid Dreams	WRLD	285p8ltuF7W6TeW8hRDR	0.511	0.566	6	0	0.20	0.349	0.0	0.340	0.218
1	Better Now	Post Malone	7d6xSM1jpdT8t8C8u5tK	0.680	0.578	10	1	0.04	0.331	0.0	0.135	0.341

```
In [19]: df_msd.head(2)
Out[19]:
```

	Artist	Track	SpotifyID	danceability	energy	key	loudness	mode	speechiness	acousticness	instrumentalness	liveness	valence
0	Sonora Santanera	Amor De Cabaret	SleyXbf7CVCVQ8DqWfKfANP	0.699	0.519	8	-6.422	1	0.099	0.519	8	-6.422	1
1	JennyAnyKind	Young Boy Blues	SuN6vY2NLG6GKoywHKZA	0.458	0.570	5	-9.159	0	0.458	0.570	5	-9.159	0

List of unique spotifyID that are present in bill board data and are hit.

```
In [20]: spotify_id_hit = list(df_bb['SpotifyID'].unique())
spotify_id_hit[0:5]
```

```
['285p8ltuF7W6TeW8hRDR',
'1R1d8tM1jpdT8t8C8u5tK',
'76Rvc8BfY5HqQzQqJh9AN',
'2a1ML1CQ3DGFmKkPnL2D9n',
'21UXuYVDBRvQ8w8eqP7012']
```

Total number of hit tracks

Below are the ids present in only billboard and refers to hit songs

```
In [21]: len(spotify_id_hit)
Out[21]: 9258
```

```
In [22]: df_hit = df_bb[df_bb['SpotifyID'].isin(df_bb['SpotifyID'])]
df_hit.shape
Out[22]: (9258, 15)
```

Below are the ids not present in billboard and refers to non-hit songs

```
In [23]: df_non_hit = df_msd[~df_msd['SpotifyID'].isin(df_bb['SpotifyID'])]
df_non_hit.shape
Out[23]: (5424, 15)
```

```
In [24]: df_hit.head(2)
Out[24]:
```

	Track	Artist	SpotifyID	danceability	energy	key	mode	speechiness	acousticness	instrumentalness	liveness	valence
0	Lucid Dreams	WRLD	285p8ltuF7W6TeW8hRDR	0.511	0.566	6	0	0.20	0.349	0.0	0.340	0.218
1	Better Now	Post Malone	7d6xSM1jpdT8t8C8u5tK	0.680	0.578	10	1	0.04	0.331	0.0	0.135	0.341

```
In [25]: df_non_hit.head(2)
Out[25]:
```

	Artist	Track	SpotifyID	danceability	energy	key	loudness	mode	speechiness	acousticness	instrumentalness	liveness	valence
0	Sonora Santanera	Amor De Cabaret	SleyXbf7CVCVQ8DqWfKfANP	0.699	0.519	8	-6.422	1	0.099	0.519	8	-6.422	1
1	JennyAnyKind	Young Boy Blues	SuN6vY2NLG6GKoywHKZA	0.458	0.570	5	-9.159	0	0.458	0.570	5	-9.159	0

We have to set the sequence of the dataframe same for both.

```
In [26]: df_hit.columns
Out[26]: Index(['Track', 'Artist', 'SpotifyID', 'danceability', 'energy', 'key', 'mode', 'speechiness', 'acousticness', 'instrumentalness', 'liveness', 'valence', 'tempo', 'duration_ms', 'loudness'],
dtype='object')
```

Rearrange columns of df_non_hit as per df_hit dataframe

```
In [27]: df_non_hit = df_non_hit[['Track', 'Artist', 'SpotifyID', 'danceability', 'energy', 'key', 'mode', 'speechiness', 'acousticness', 'instrumentalness', 'liveness', 'valence', 'tempo', 'duration_ms', 'loudness']]
df_non_hit.shape
Out[27]: (5424, 15)
```

The sequence of the columns are now same for both.

```
In [28]: df_hit.head(2)
Out[28]:
```

	Track	Artist	SpotifyID	danceability	energy	key	mode	speechiness	acousticness	instrumentalness	liveness	valence
0	Lucid Dreams	WRLD	285p8ltuF7W6TeW8hRDR	0.511	0.566	6	0	0.20	0.349	0.0	0.340	0.218
1	Better Now	Post Malone	7d6xSM1jpdT8t8C8u5tK	0.680	0.578	10	1	0.04	0.331	0.0	0.135	0.341

```
In [29]: df_non_hit.head(2)
Out[29]:
```

	Track	Artist	SpotifyID	danceability	energy	key	mode	speechiness	acousticness	instrumentalness	liveness	valence
0	Amor De Cabaret	Sonora Santanera	SleyXbf7CVCVQ8DqWfKfANP	0.699	0.519	8	1	0.099	0.519	8	-6.422	1
1	Young Boy Blues	JennyAnyKind	SuN6vY2NLG6GKoywHKZA	0.458	0.570	5	0	0.458	0.570	5	-9.159	0

Appending column 'is hit' indicator as 0 for df_non_hit and 1 for df_hit

```
In [30]: df_hit['is hit']=1
df_non_hit['is hit']=0
```

New column is hit added successfully for both the dataset and ready to be merged together.

```
In [31]: df_hit.head(2)
Out[31]:
```

	Track	Artist	SpotifyID	danceability	energy	key	mode	speechiness	acousticness	instrumentalness	liveness	valence
0	Lucid Dreams	WRLD	285p8ltuF7W6TeW8hRDR	0.511	0.566	6	0	0.20	0.349	0.0	0.340	0.218
1	Better Now	Post Malone	7d6xSM1jpdT8t8C8u5tK	0.680	0.578	10	1	0.04	0.331	0.0	0.135	0.341

```
In [32]: df_non_hit.head(2)
Out[32]:
```

	Track	Artist	SpotifyID	danceability	energy	key	mode	speechiness	acousticness	instrumentalness	liveness	valence
0	Amor De Cabaret	Sonora Santanera	SleyXbf7CVCVQ8DqWfKfANP	0.699	0.519	8	1	0.099	0.519	8	-6.422	1
1	Young Boy Blues	JennyAnyKind	SuN6vY2NLG6GKoywHKZA	0.458	0.570	5	0	0.458	0.570	5	-9.159	0

```
In [33]: df_merge = pd.concat([df_hit,df_non_hit])
In [34]: df_merge.shape
Out[34]: (14682, 16)
```

Shuffling the dataframe and exporting the data for pre-processing.

```
In [35]: df_merge = df_merge.sample(frac=1).reset_index(drop=True)
So, this is the data we can use for further pre-processing.
```

```
In [37]: df_merge.head()
Out[37]:
```

	Track	Artist	SpotifyID	danceability	energy	key	mode	speechiness	acousticness	instrumentalness	liveness	valence
0	Firedance	Mythic Rhythms Band	2yScBrgZv6WnaUfJLCL	0.637	0.876	8	1	0.0415	0.1060	0.313000	0.136	0.0
1	Stop The World	Extreme	4Vghb9YLD9BudOpMKWOK	0.452	0.724	8	1	0.0435	0.0038	0.001120	0.687	0.0
2	Nights	Frank Ocean	7eqoqGKwvQaWNNh9u0Ez	0.466	0.548	5	0	0.1180	0.4200	0.000001	0.113	0.0
3	Suds In The Bucket	Sara Evans	6NhpldYofuNnNBgOztc	0.511	0.903	10	1	0.0397	0.0520	0.000529	0.304	0.0
4	Shrivatsa	William Coultier	0LkxMfW5u8z1vqGCH6G	0.358	0.285	9	0	0.0381	0.9410	0.893000	0.121	0.0

EDA and Visualization

```
In [38]: df_merge.info()
Out[38]:
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 14682 entries, 0 to 14681
Data columns (total 16 columns):
#   Column                Non-Null Count  Dtype  ---
0   Track                 14682 non-null object
1   Artist                14682 non-null object
2   SpotifyID             14682 non-null object
3   danceability           14682 non-null float64
4   energy                 14682 non-null float64
5   key                   14682 non-null int64
6   mode                  14682 non-null int64
7   speechiness           14682 non-null float64
8   acousticness          14682 non-null float64
9   instrumentalness       14682 non-null float64
10  liveness               14682 non-null float64
11  valence                14682 non-null float64
12  tempo                 14682 non-null float64
13  duration_ms            14682 non-null int64
14  loudness               14682 non-null float64
15  is_hit                 14682 non-null int64
dtypes: float64(9), int64(4), object(3)
memory usage: 1.8+ MB
```

Setting spotifyID as index as it is unique for all records.

```
In [39]: df_merge.set_index('SpotifyID',inplace=True)
In [40]: df_merge.head()
Out[40]:
```

	Track	Artist	danceability	energy	key	mode	speechiness	acousticness	instrumentalness	liveness	valence	
	2yScBrgZv6WnaUfJLCL	Firedance	Mythic Rhythms Band	0.637	0.876	8	1	0.0415	0.1060	0.313000	0.136	0.0
	4Vghb9YLD9BudOpMKWOK	Stop The World	Extreme	0.452	0.724	8	1	0.0435	0.0038	0.001120	0.687	0.0
	7eqoqGKwvQaWNNh9u0Ez	Nights	Frank Ocean	0.466	0.548	5	0	0.1180	0.4200	0.000001	0.113	0.0
	6NhpldYofuNnNBgOztc	Suds In The Bucket	Sara Evans	0.511	0.903	10	1	0.0397	0.0520	0.000529	0.304	0.0
	0LkxMfW5u8z1vqGCH6G	Shrivatsa	William Coultier	0.358	0.285	9	0	0.0381	0.9410	0.893000	0.121	0.0

Missing values Treatment

```
In [41]: df_merge.isna().sum()
Out[41]:
Track                1
Artist               2
danceability          0
energy                0
key                  0
mode                 0
speechiness           0
acousticness          0
instrumentalness      0
liveness              0
valence               0
tempo                 0
duration_ms           0
loudness              0
is_hit                0
dtype: int64
```

Dropping the records where Track and Artist are null.

```
In [42]: df_merge = df_merge.dropna(subset=['Track','Artist'],axis=0)
Now, No missing values left to be treated.
```

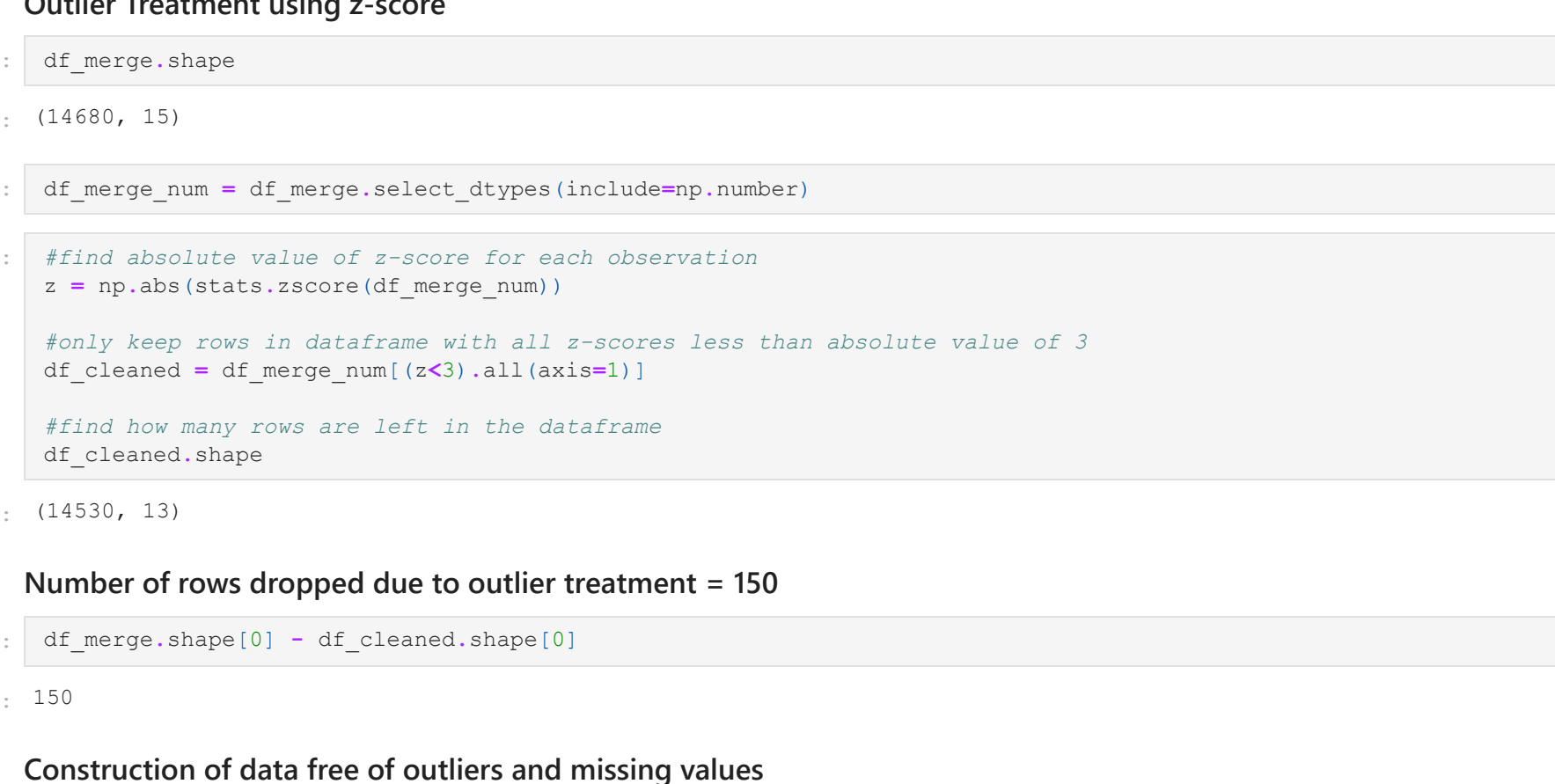
```
In [43]: df_merge.info()
Out[43]:
<class 'pandas.core.frame.DataFrame'>
Int64Index: 14680 entries, 2yScBrgZv6WnaUfJLCL to 17yQoQo5RkdA7iukL0okv
Data columns (total 15 columns):
#   Column                Non-Null Count  Dtype  ---
0   Track                 14680 non-null object
1   Artist                14680 non-null object
2   energy                14680 non-null float64
3   key                   14680 non-null int64
4   mode                  14680 non-null int64
5   speechiness           14680 non-null float64
6   acousticness          14680 non-null float64
7   instrumentalness       14680 non-null float64
8   liveness               14680 non-null float64
9   valence                14680 non-null float64
10  tempo                 14680 non-null float64
11  duration_ms            14680 non-null int64
12  loudness               14680 non-null float64
13  is_hit                 14680 non-null int64
dtypes: float64(9), int64(4), object(2)
memory usage: 1.8+ MB
```

```
In [44]: df_numeric = df_merge.select_dtypes(include=np.number)
In [45]: df_numeric.columns
Out[45]: Index(['danceability', 'energy', 'key', 'mode', 'speechiness', 'acousticness', 'instrumentalness', 'liveness', 'valence', 'tempo', 'duration_ms', 'loudness'],
dtype='object')
```

```
In [46]: ss = StandardScaler()
In [47]: df_numeric_ss = ss.fit_transform(df_numeric)
In [48]: df_numeric_ss = pd.DataFrame(df_numeric_ss,columns=df_numeric.columns)
Out[48]:
```

	danceability	energy	key	mode	speechiness	acousticness	instrumentalness	liveness	valence	tempo	duration_ms	loudness
0	0.027493	0.024721	0.128226	0.038577	0.023955	0.021275	0.024854	0.027722	0.031735	0.509062	-0.376258	-0.028872
1	0.020399	0.028894	0.128226	0.038577	0.024042	0.017355	0.022890	0.044857	0.018429	1.310263	1.294579	0.009862
2	0.020936	0.022146	0.014781	0.008243	0.026900	0.033318	0.022847	0.022840	0.022187	-0.719670	0.735136	-0.034398
3	0.020662	0.032575	0.020387	0.038577	0.023896	0.019203	0.022867	0.030166	0.040400	1.008523	-0.144910	0.201838
4	0.016795	0.012063	0.166041	0.000243	0.023835	0.053301	0.057104	0.023147	0.018276	0.694093	-0.020203	-0.163922

Plotting histogram of binned numeric data



Outlier Treatment using z-score

```
In [51]: df_merge.shape
Out[51]: (14680, 15)
```

```
In [52]: df_merge_num = df_merge.select_dtypes(include=np.number)
In [53]: #find absolute value of z-score for each observation
z = np.abs(stats.zscore(df_merge_num))
#only keep rows in dataframe with all z-scores less than absolute value of 3
df_cleaned = df_merge_num[(z<3).all(axis=1)]
#find how many rows are left in the dataframe
df_cleaned.shape
Out[53]: (14530, 13)
```

Number of rows dropped due to outlier treatment = 150

```
In [54]: df_merge.shape[0] = df_cleaned.shape[0]
Out[54]: 150
```

Construction of data free of outliers and missing values

```
In [55]: df_merge = df_merge.index.isin(df_cleaned.index)
df_merge.shape
Out[55]: (14530, 15)
```

```
In [56]: df_merge.isna().sum()
Out[56]:
Track                0
Artist               0
danceability          0
energy                0
key                  0
mode                 0
speechiness           0
acousticness          0
instrumentalness      0
liveness              0
valence               0
tempo                 0
duration_ms           0
loudness              0
is_hit                0
dtype: int64
```

Export data for modeling

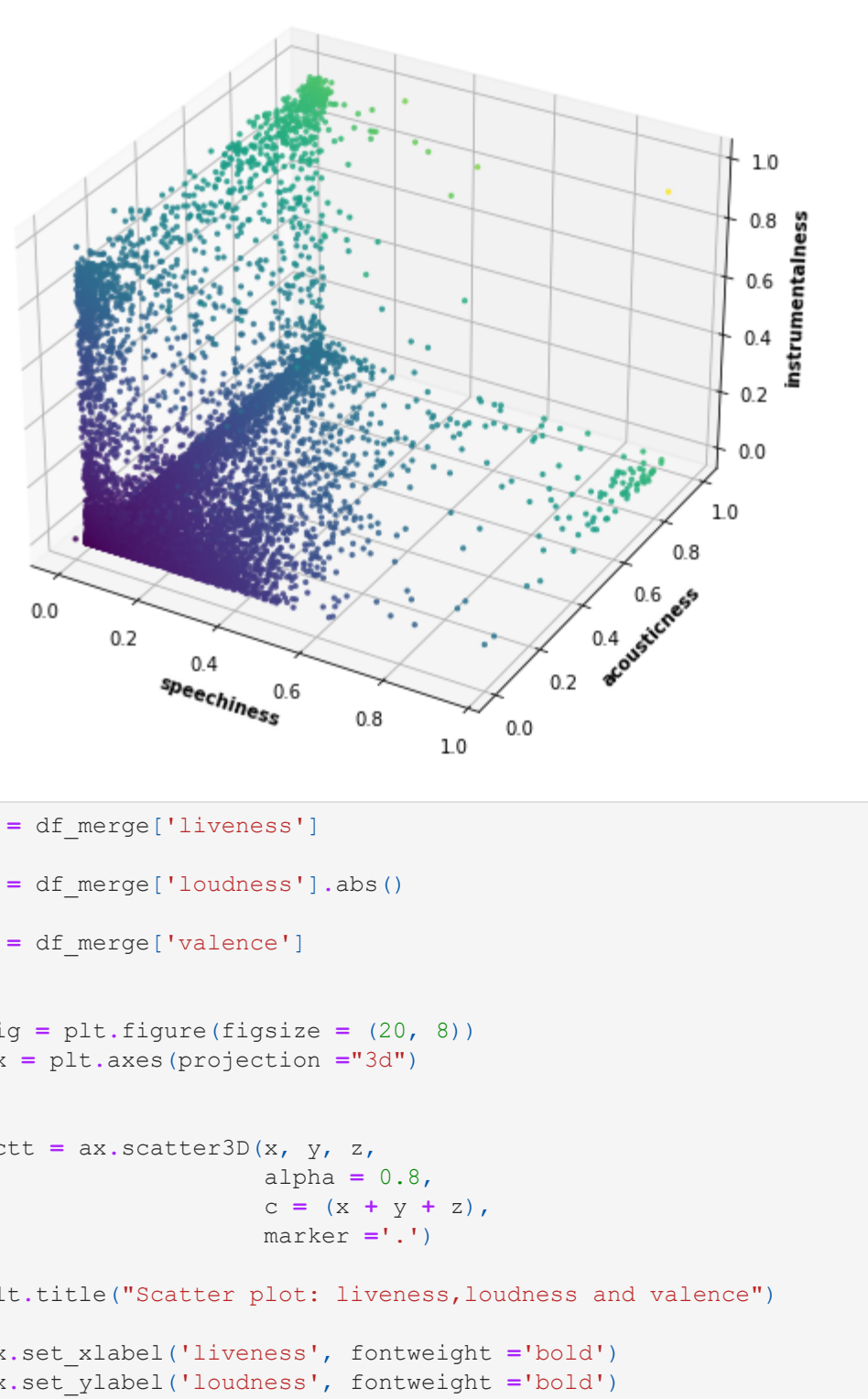
```
In [57]: df_merge.to_csv('bb_merged_exported.csv')
print('Merged dataframe exported.....')
Merged dataframe exported.....
```

Checking if the target data is Imbalanced or not. The ratio is around 63:37

```
In [58]: df_merge['is hit'].value_counts(normalize=True).to_frame()
Out[58]:
```

	is hit
1	0.634893

Scatter plot: speechiness,acousticness and instrumentals



```
In [79]: x = df_merge['liveness']
y = df_merge['loudness'].abs()
z = df_merge['valence']

fig = plt.figure(figsize=(20, 8))
ax = plt.axes(projection='3d')

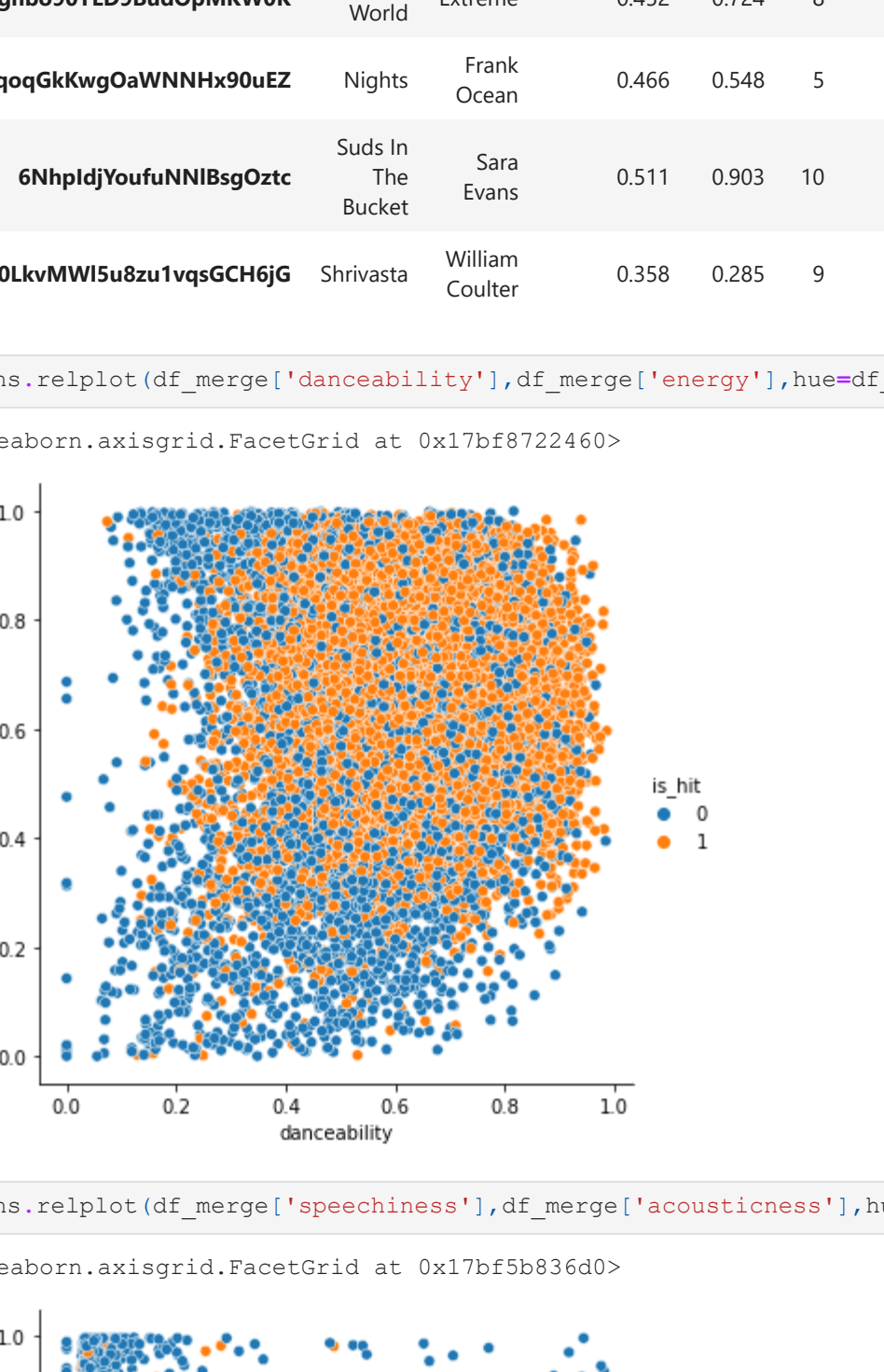
sctt = ax.scatter3D(x, y, z,
                    alpha = 0.8,
                    c = [x + y + z],
                    marker = '+')

plt.title("Scatter plot: liveness,loudness and valence")

ax.set_xlabel('liveness', fontweight='bold')
ax.set_ylabel('loudness', fontweight='bold')
ax.set_xlabel('valence', fontweight='bold')
ax.set_ylabel('valence', fontweight='bold')
fig.colorbar(sctt, ax = ax, shrink = 0.5, aspect = 5)

# show plot
plt.show()
```

Scatter plot: liveness,loudness and valence



Relationship plots (2D scatter plots)

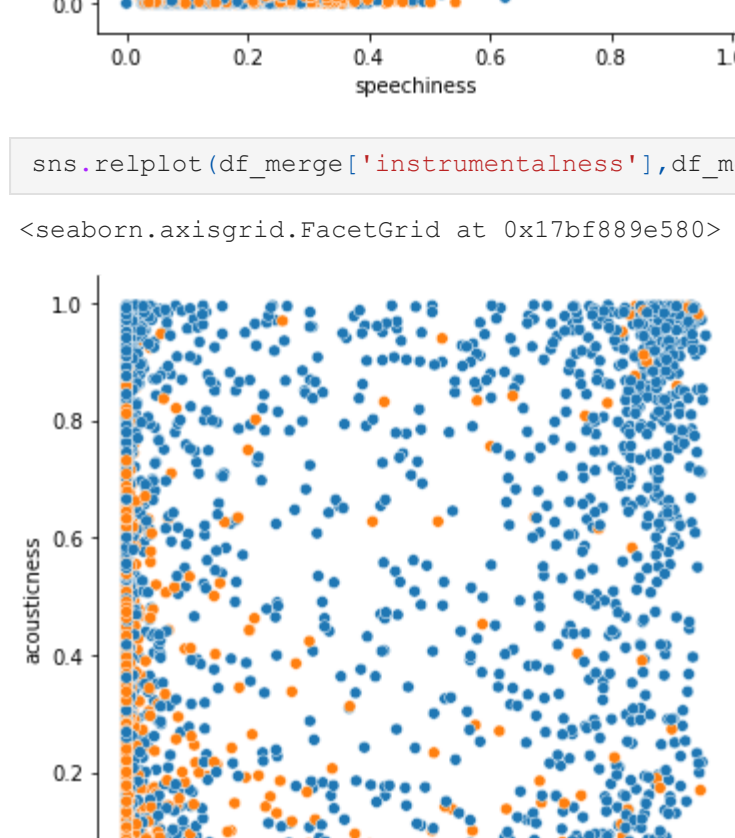
```
In [76]: df_merge.head()

Out[76]:
```

	Track	Artist	danceability	energy	key	mode	speechiness	acousticness	instrumentals	liveness	valence	
	SpotifyID											
	2y5cBrg0zV6WnaJfJLCL	Firedance	Mythic Rhythms Band	0.637	0.876	8	1	0.0415	0.1060	0.313000	0.136	0
	4Vghbo5YLD9BudOpMKW0K	Stop The World	Extreme	0.452	0.724	8	1	0.0435	0.0038	0.001120	0.687	0
	7eqeqGKKwgQaWNNHx9uE2	Nights	Frank Ocean	0.466	0.548	5	0	0.1180	0.4200	0.000001	0.113	0
	6NhpIdjYofuNnNBgOxtc	Suds In The Bucket	Sara Evans	0.511	0.903	10	1	0.0397	0.0520	0.000529	0.304	0
	0LlvMWISu8uTvpqGCH6g	Shrivasta	William Coulter	0.358	0.285	9	0	0.0381	0.9410	0.893000	0.121	0

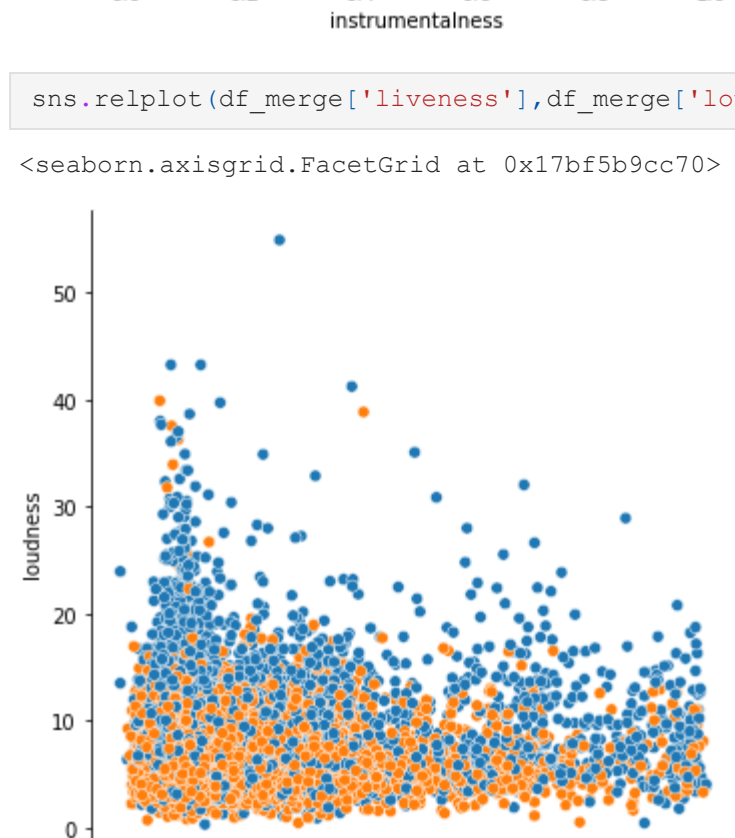
```
In [77]: sns.relplot(df_merge['danceability'],df_merge['energy'],hue=df_merge['is_hit'])

Out[77]: <seaborn.axisgrid.FacetGrid at 0x17bf8722460>
```



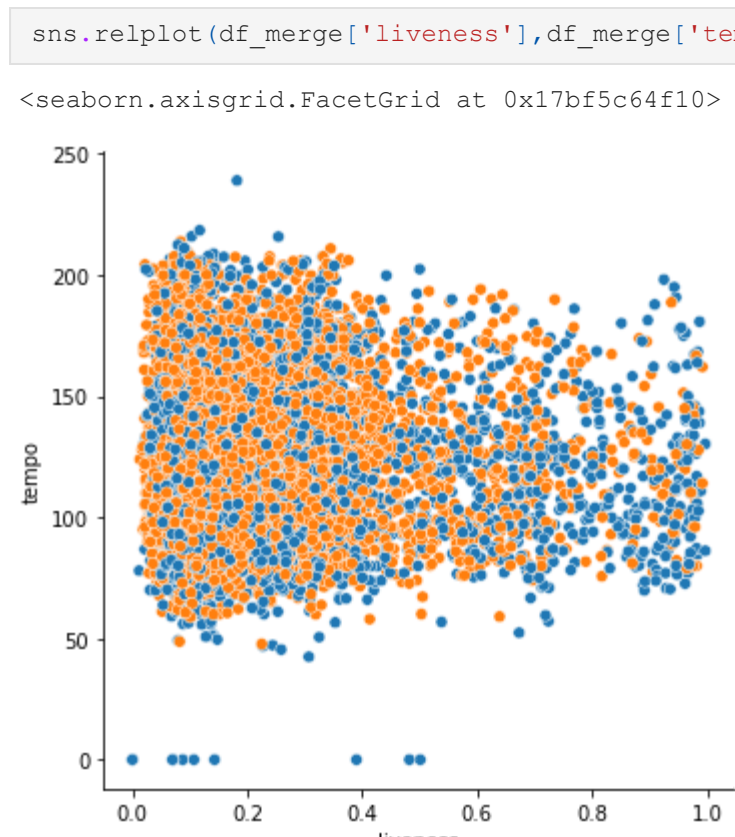
```
In [78]: sns.relplot(df_merge['speechiness'],df_merge['acousticness'],hue=df_merge['is_hit'])

Out[78]: <seaborn.axisgrid.FacetGrid at 0x17bf8b836d0>
```



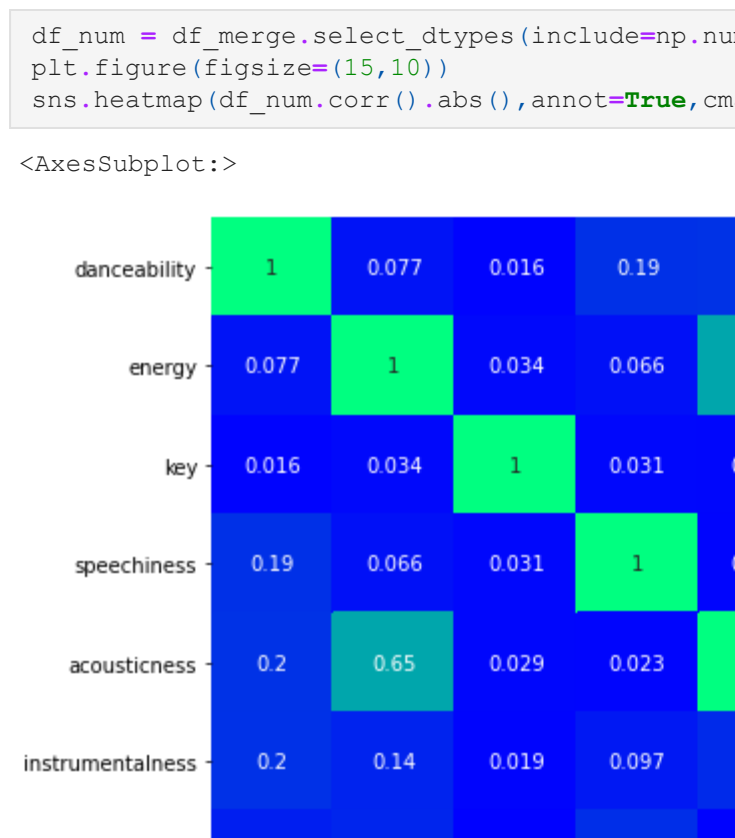
```
In [79]: sns.relplot(df_merge['instrumentals'],df_merge['acousticness'],hue=df_merge['is_hit'])

Out[79]: <seaborn.axisgrid.FacetGrid at 0x17bf889e580>
```



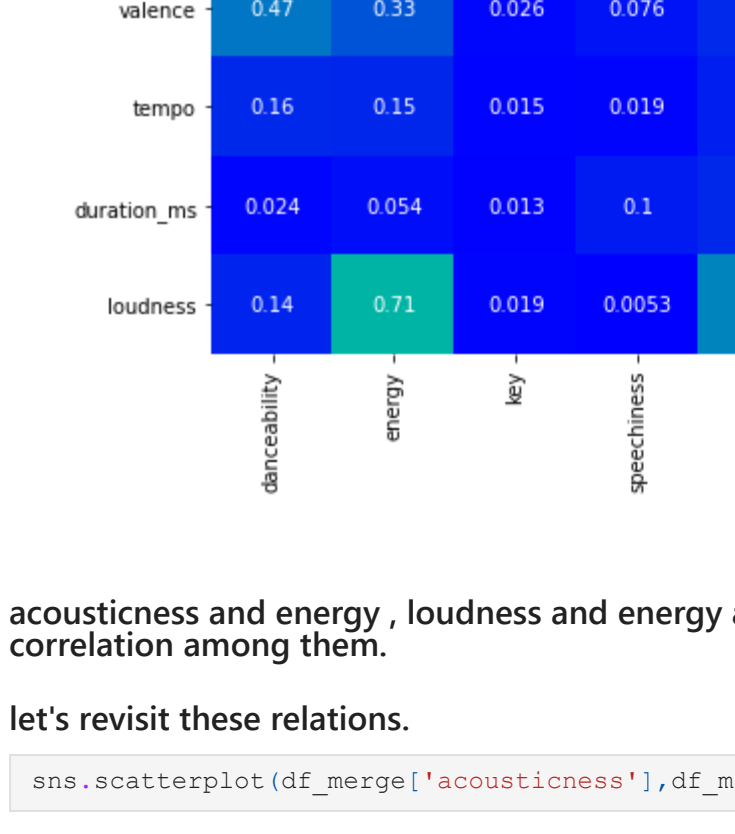
```
In [80]: sns.relplot(df_merge['liveness'],df_merge['loudness'].abs(),hue=df_merge['is_hit'])

Out[80]: <seaborn.axisgrid.FacetGrid at 0x17bf8b9c70>
```



```
In [81]: sns.relplot(df_merge['liveness'],df_merge['tempo'].abs(),hue=df_merge['is_hit'])

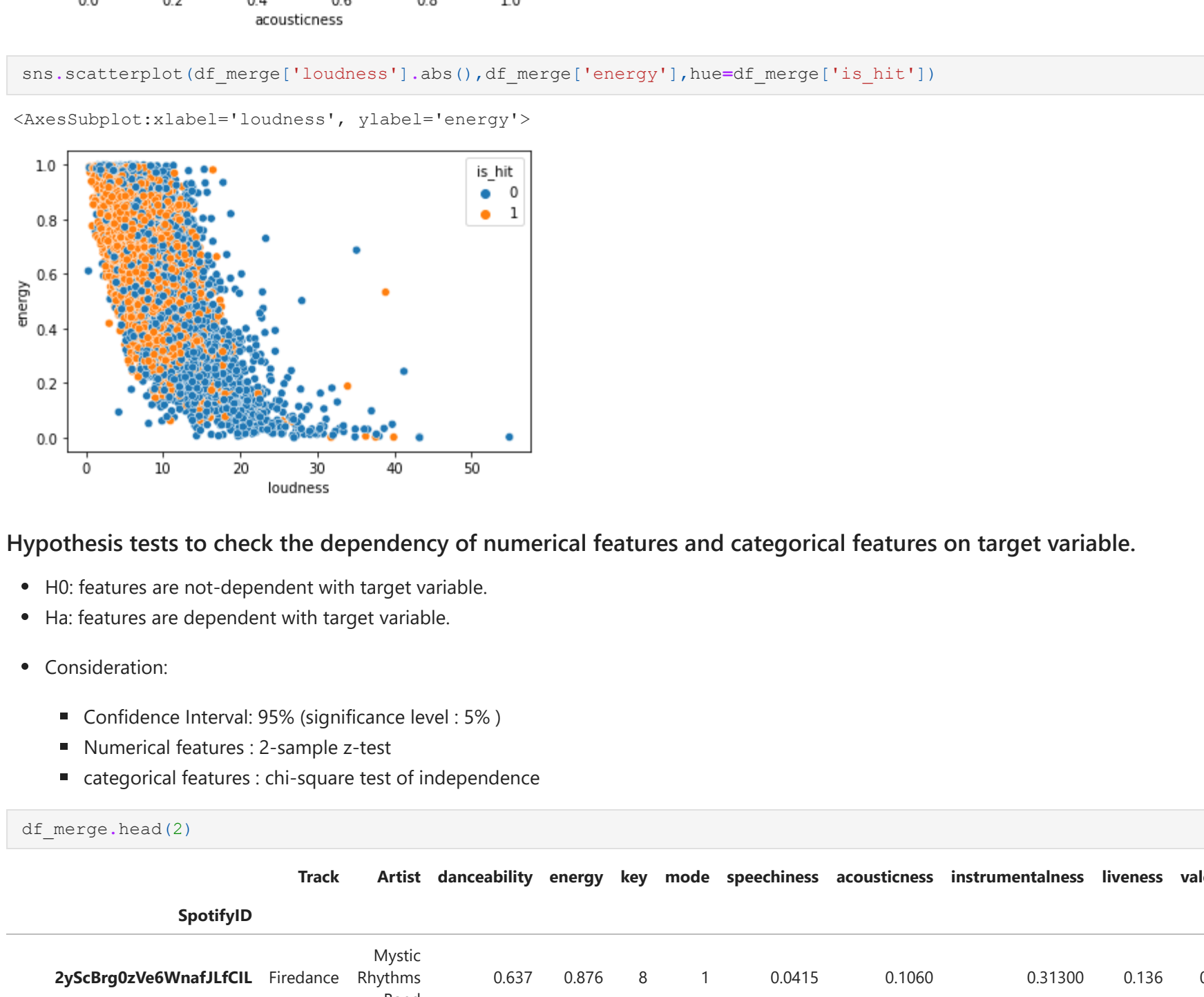
Out[81]: <seaborn.axisgrid.FacetGrid at 0x17bf8f5c64f0>
```



Overall Relationship plot : correlation heatmap.

```
In [82]: df_num = df_merge.select_dtypes(include=[np.number]).drop(['mode','is_hit'],axis=1)
plt.figure(figsize=(15,10))
sns.heatmap(df_num.corr().abs(),annot=True,cmap='winter')

Out[82]: <AxesSubplot>
```

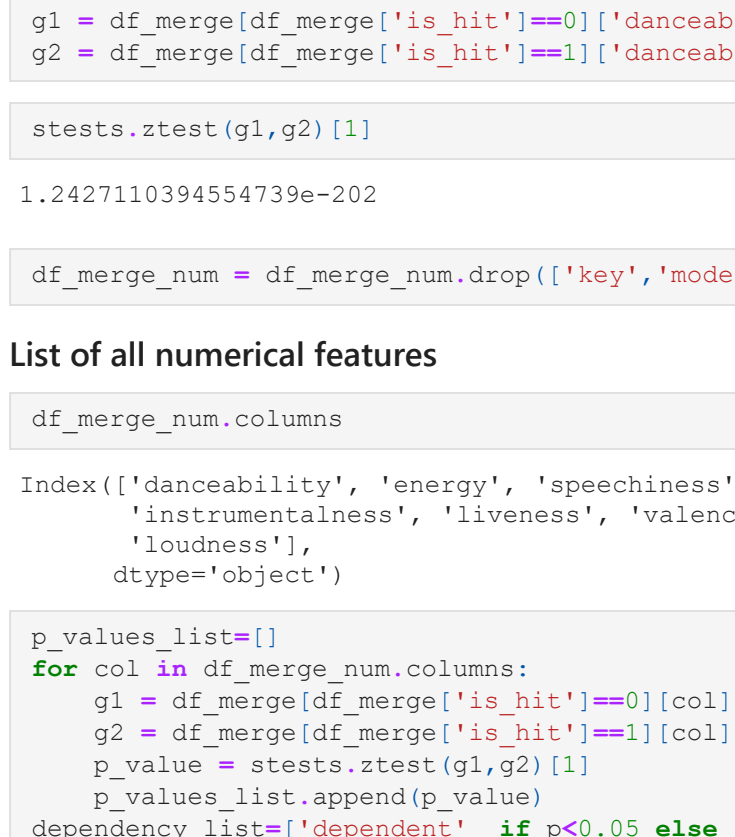


acousticness and energy, loudness and energy are the only pair of features which are having significant correlation among them.

let's revisit these relations.

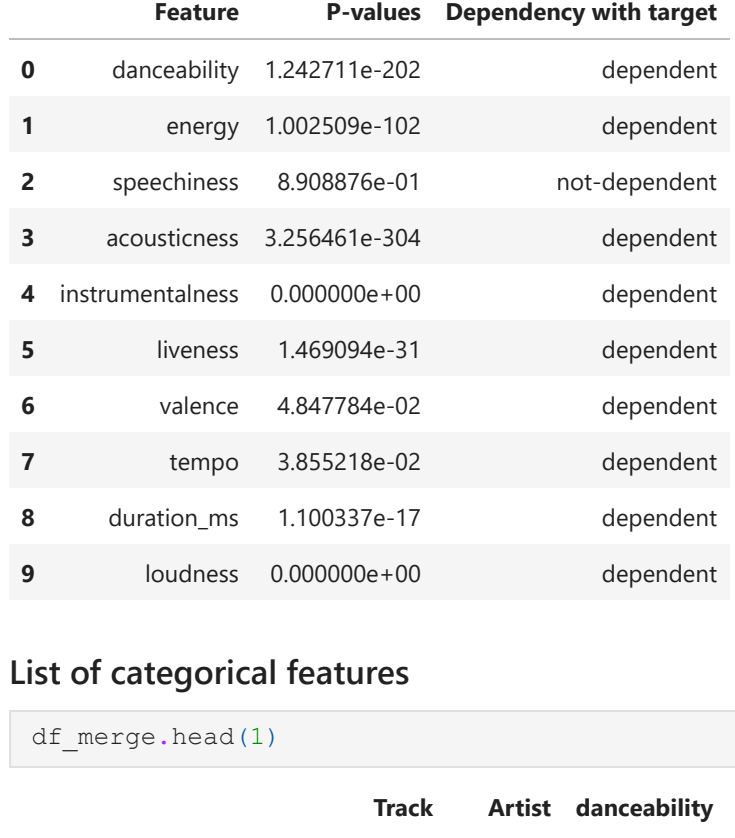
```
In [83]: sns.scatterplot(df_merge['acousticness'],df_merge['energy'],hue=df_merge['is_hit'])

Out[83]: <AxesSubplot:label='acousticness', ylabel='energy'>
```



```
In [84]: sns.scatterplot(df_merge['loudness'].abs(),df_merge['energy'],hue=df_merge['is_hit'])

Out[84]: <AxesSubplot:label='loudness', ylabel='energy'>
```



Hypothesis tests to check the dependency of numerical features and categorical features on target variable.

- H0: Features are not dependent with target variable.
- Ha: Features are dependent with target variable.
- Consideration:
 - Confidence interval: 95% (significance level: 5%)
 - Numerical features: 2-sample z-test
 - categorical features: chi-square test of independence

```
In [85]: df_merge.head(2)

Out[85]:
```

	Feature	P-values	Dependency with target
0	key	3.708598e-33	dependent
1	mode	1.135764e-02	dependent

From the above hypothesis test with 95% Confidence Interval we can clearly see that apart from speechiness all are significant predictors of target variable.

Multicollinearity check : energy column is having highest VIF factor but not significantly large to be dropped.

```
from statsmodels.stats.outliers_influence import variance_inflation_factor

X = df.merge(drop(['Track','Artist','is_hit'],axis=1))

# VIF dataframe
vif_data = pd.DataFrame()
vif_data["Feature"] = X.columns

# calculating VIF for each feature
vif_data["VIF"] = [variance_inflation_factor(X.values, i)
                    for i in range(len(X.columns))]
```

```
In [86]: g1 = df_merge[df_merge['is_hit']==0]['danceability']
g2 = df_merge[df_merge['is_hit']==1]['danceability']

In [87]: stests.ztest(g1,g2)[1]

Out[87]: 1.242711039454739e-202

In [88]: df_merge_num = df_merge_num.drop(['key','mode','is_hit'],axis=1)

Out[88]:
```

List of all numerical features

```
In [89]: df_merge_num.columns

Out[89]: Index(['danceability', 'energy', 'speechiness', 'acousticness', 'instrumentals', 'liveness', 'valence', 'tempo', 'duration_ms', 'loudness'], dtype='object')
```

```
In [90]: p_values_list=[]
for col in df_merge_num.columns:
    g1 = df_merge[df_merge['is_hit']==0][col]
    g2 = df_merge[df_merge['is_hit']==1][col]
    p_value = stats.ztest(g1,g2)[1]
    p_values_list.append(p_value)
dependency_list=['dependent' if p<0.05 else 'not-dependent' for p in p_values_list ]
pd.DataFrame({'Feature':df_merge_num.columns, 'P-values':p_values_list,'Dependency with target':dependency_list})

Out[90]:
```

Feature	P-values	Dependency with target	
0	danceability	1.242711e-202	dependent
1	energy	1.002509e-102	dependent
2	speechiness	8.908876e-01	not-dependent
3	acousticness	3.25641e-304	dependent
4	instrumentals	0.000000e+00	dependent
5	liveness	1.469094e-31	dependent
6	valence	4.847784e-02	dependent
7	tempo	3.855218e-02	dependent
8	duration_ms	1.100337e-17	dependent
9	loudness	0.000000e+00	dependent

List of categorical features

```
In [91]: df_merge.head(1)

Out[91]:
```

	Track	Artist	danceability	energy	key	mode	speechiness	acousticness	instrumentals	liveness	valence	
	SpotifyID											
	2y5cBrg0zV6WnaJfJLCL	Firedance	Mythic Rhythms Band	0.637	0.876	8	1	0.0415	0.1060	0.313000	0.136	0.672

```
In [92]: df_merge_cat = df_merge[['key','mode']]

Out[92]:
```

```
Index(['key', 'mode'], dtype='object')
```

```
In [93]: p_val=[]
for col in df_merge_cat.columns:
    d1 = pd.crosstab(df_merge['is_hit'],df_merge[col])
    p_val = stats.chi2_contingency(d1)[1]
    p_val = stats.ztest(g1,g2)[1]
    p_val = stats.ztest(g1,g2)[1]
    dependency_list_cat=['dependent' if p<0.05 else 'not-dependent' for p in p_val ]
    dependency_list_cat=['dependent' if p<0.05 else 'not-dependent' for p in p_val ]
pd.DataFrame({'Feature':df_merge_cat.columns, 'P-values':p_val,'Dependency with target':dependency_list_cat})

Out[93]:
```

Feature	P-values	Dependency with target	
0	key	3.708598e-33	dependent
1	mode	1.135764e-02	dependent

From the above hypothesis test with 95% Confidence Interval we can clearly see that apart from speechiness all are significant predictors of target variable.

Multicollinearity check : energy column is having highest VIF factor but not significantly large to be dropped.

```
In [124]: from statsmodels.stats.outliers_influence import variance_inflation_factor
X = df_merge.drop(['Track','Artist','is_hit'],axis=1)
# VIF dataframe
vif_data = pd.DataFrame()
vif_data['Feature'] = X.columns

# calculating VIF for each feature
vif_data['VIF'] = [variance_inflation_factor(X.values, i) for i in range(len(X.columns))]

vif_data = vif_data.sort_values(by='VIF',ascending=False)
vif_data

Out[124]:
```

Feature	VIF	
1	energy	19.171371
0	danceability	16.374929
9	tempo	15.011448
10	duration_ms	14.263656
8	valence	8.691957
11	loudness	8.142193
2	key	3.210915
3	mode	3.111841
5	acousticness	2.622257
7	liveness	2.586114
4	speechiness	1.939668
6	instrumentals	1.393568

```
In [130]: plt.figure(figsize=(20,8))
sns.barplot(x='Feature',y='VIF',data=vif_data)

Out[130]:
```



Conclusion :

- Columns Track and Artist needs to be dropped.
- As per statistical significance 'speechiness' column is not a good predictor and can be dropped.
- As per vif factor column 'energy' is involved in multicollinearity and can be dropped.
- Data is moderately balanced, so oversampling not required.
- There is no clear difference of target variable among numerical features visually, no inference can be drawn prior to model building.