

Virtual Internship Program

Author :- Prateek Kumar Singh

Beginner Level Tasks

Task-3 Music Recommendation

Music recommender systems can suggest songs to users based on their listening patterns.

Datasetlink :- <https://www.kaggle.com/ckkbox-music-recommendation-challenge/data>

1. Importing The Libraries

```
pip install recommender-system
```

```
In [38]: import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
import warnings
warnings.filterwarnings('ignore')
```

2. Load Music Data

```
In [38]: train = pd.read_csv('train.csv')
songs = pd.read_csv('songs.csv')
members = pd.read_csv('members.csv')
```

```
In [39]: df_train.head()
```

	msno	song_id	source_system_tab	source_screen_name	source_type	target
0	FGfVqz1BRPwJyedr2gV78ziAIY/9SmYvIa+HCg=	BBzumQNVLuHKEBOB7mAJuzok+JAlc2RygyzTF6Ik=	explore	Explore	online-playlist	1
1	Xumu+NiS6QYVvDS4t3SawJ7vT9hPKXmfbORLlNx8=	bhpMpSNoqxOIB+BWwPqU6jBt4DpCn3ayXnJqM=	my library	Local playlist more	local-playlist	1
2	Xumu+NiS6QYVvDS4t3SawJ7vT9hPKXmfbORLlNx8=	JNWhrC7zNt7BdMpslSKa4Mw+XVJYnXn3Epw7QgY=	my library	Local playlist more	local-playlist	1
3	Xumu+NiS6QYVvDS4t3SawJ7vT9hPKXmfbORLlNx8=	2A87zhtJTSWgD7gIzHsoIhe4DMzkb6uLzOJKHjN=	my library	Local playlist more	local-playlist	1
4	FGfVqz1BRPwJyedr2gV78ziAIY/9SmYvIa+HCg=	3qm6XTZ8MOCU11k8FVAGH5G5UMKT3ZaWAG1oo2G=	explore	Explore	online-playlist	1

```
In [40]: df_songs.head()
```

	song_id	song_length	genre_ids	artist_name	composer	lyricist	language
0	CkoT1e6t7H+Dn6tUJvccwGRV45CDuZu+YD8JPB+HE=	247640	465	张信哲 (Jeff Chang)	重宝	何啟弘	3.0
1	cdkF7a6zQeivgRvPvqJw6d5zHRUJf7O1Dw0ZDU=	197328	444	BLACKPINK	TEDDY FUTURE BOUNCE	Bekuh BOOM	31.0
2	DwVVwVuz+XpUuccvQV6yPqgUkHrone1RQzP9=	231781	465	SUPER JUNIOR		NaN	31.0
3	dKMBWuZjScdsShkHG+Vt47nc1Bn9q4m58+b4e7dSSE=	273554	465	S.H.E	潘小康	徐世珍	3.0
4	W3wZw5t+veHfHAUfARgW9AvR4f4N5Yzm4M6Ew=	140329	726	黄家驹	Traditional	Traditional	52.0

```
In [41]: df_songs_extra.head()
```

	song_id	name	isrc
0	LPTjL3oJFByu6jww+eL3T+4h+UnBPURchcXJjg=	我們	TMUAT71200043
1	ClacTFa6zOBnuie4bocdNM3rding6CGAGUWcHE=	Let Me Love You	QMZSY1600015
2	u2qay0p6+HrHe2E6LKHh0RHR1kq1TBJoCIW9v+Ts=	原以为	TWAS30867303
3	Q2Fay0p6+HrHe2E6LKHh0RHR1kq1TBJoCIW9v+Ts=	Classic	USMS11301446
4	QQFmz+JyJQ56C1DuYt9HKkG5TUpQnNlwo+Hw=	爱投爱猜	TWAT47130001

```
In [42]: df_test.head()
```

	id	msno	song_id	source_system_tab	source_screen_name	source_type
0	0	V8yV75Gz7Dn3z45D1PpRqgqivmVMKk21q54uM+	WmHk9MplQMc6hNvDMkVvzYVHvFw07725lss=	my library	Local playlist more	local-library
1	1	V8yV75Gz7Dn3z45D1PpRqgqivmVMKk21q54uM+	ynZ8OC7FwKSF3PKZD5nyv4OBu3z614NgE0vM+	my library	Local playlist more	local-library
2	2	AzQ4hAbacv+HWK0zDpF2kxvP9vY7d+gblUwY=	8ZLFOGCVvS85v4vnuLqgH2BvKc7QYlM53Dk4=	discover	NaN	song-based-playlist
3	3	Ia6oolXKAbQ4e59zTVd+KSvIAFpTtUwvWLC1Y0k=	zCCKm1rYs4YVNG3GdLbvXzLmT5mTBVXOO4C9NfVQ=	radio	Radio	radio
4	4	Ia6oolXKAbQ4e59zTVd+KSvIAFpTtUwvWLC1Y0k=	MKVMapRKCQMaFEGvEjCEH5+RZdMYU3uF0dySMBY=	radio	Radio	radio

3. Create New data

```
In [43]: res = df_train.merge(df_songs[['song_id','song_length','genre_ids','artist_name','language']], on=['song_id'], how='left')
res.head()
```

```
Out[43]:
```

	msno	song_id	source_system_tab	source_screen_name	source_type	target	song_length	genre_ids	artist
0	FGfVqz1BRPwJyedr2gV78ziAIY/9SmYvIa+HCg=	BBzumQNVLuHKEBOB7mAJuzok+JAlc2RygyzTF6Ik=	explore	Explore	online-playlist	1	NaN	NaN	
1	Xumu+NiS6QYVvDS4t3SawJ7vT9hPKXmfbORLlNx8=	bhpMpSNoqxOIB+BWwPqU6jBt4DpCn3ayXnJqM=	my library	Local playlist more	local-playlist	1	NaN	NaN	
2	Xumu+NiS6QYVvDS4t3SawJ7vT9hPKXmfbORLlNx8=	JNWhrC7zNt7BdMpslSKa4Mw+XVJYnXn3Epw7QgY=	my library	Local playlist more	local-playlist	1	225396.0	1259	
3	Xumu+NiS6QYVvDS4t3SawJ7vT9hPKXmfbORLlNx8=	2A87zhtJTSWgD7gIzHsoIhe4DMzkb6uLzOJKHjN=	my library	Local playlist more	local-playlist	1	NaN	NaN	
4	FGfVqz1BRPwJyedr2gV78ziAIY/9SmYvIa+HCg=	3qm6XTZ8MOCU11k8FVAGH5G5UMKT3ZaWAG1oo2G=	explore	Explore	online-playlist	1	187802.0	1011	Brett

```
In [44]: train = res.merge(df_songs_extra, on=['song_id'], how='left')
train.head()
```

```
Out[44]:
```

	msno	song_id	source_system_tab	source_screen_name	source_type	target	song_length	genre_ids	artist
0	FGfVqz1BRPwJyedr2gV78ziAIY/9SmYvIa+HCg=	BBzumQNVLuHKEBOB7mAJuzok+JAlc2RygyzTF6Ik=	explore	Explore	online-playlist	1	NaN	NaN	
1	Xumu+NiS6QYVvDS4t3SawJ7vT9hPKXmfbORLlNx8=	bhpMpSNoqxOIB+BWwPqU6jBt4DpCn3ayXnJqM=	my library	Local playlist more	local-playlist	1	NaN	NaN	
2	Xumu+NiS6QYVvDS4t3SawJ7vT9hPKXmfbORLlNx8=	JNWhrC7zNt7BdMpslSKa4Mw+XVJYnXn3Epw7QgY=	my library	Local playlist more	local-playlist	1	225396.0	1259	
3	Xumu+NiS6QYVvDS4t3SawJ7vT9hPKXmfbORLlNx8=	2A87zhtJTSWgD7gIzHsoIhe4DMzkb6uLzOJKHjN=	my library	Local playlist more	local-playlist	1	NaN	NaN	
4	FGfVqz1BRPwJyedr2gV78ziAIY/9SmYvIa+HCg=	3qm6XTZ8MOCU11k8FVAGH5G5UMKT3ZaWAG1oo2G=	explore	Explore	online-playlist	1	187802.0	1011	Brett

```
In [45]: song_id = train.loc[:,['name','target']]
song_id = song_id.groupby(['name','as_of']).count().rename(columns={'target':'listen_count'})
df_train.head()
```

```
Out[45]:
```

	name	listen_count	
0	24karats -type S (與放浪兄弟EXILE、DOBELMAN合作超前衛单曲)	1	
1		A-CHA	2
2		Draw Me Close To You	1
3		I Can Only Imagine	1
4		I Concentrate On You - STACEY KENT and JIM TO...	11

```
In [46]: dataset=train.merge(song_id,on=['name'],how='left')
df=pd.DataFrame(dataset)
```

```
In [47]: df.drop(columns=['source_system_tab','source_screen_name','source_type','target','isrc'],axis=1,inplace=True)
df=df.rename(columns={'msno':'user_id'})
```

4. Loading New data

```
Out[47]: df.head()
```

	user_id	song_id	song_length	genre_ids	artist_name	language	name	listen_count
0	FGfVqz1BRPwJyedr2gV78ziAIY/9SmYvIa+HCg=	BBzumQNVLuHKEBOB7mAJuzok+JAlc2RygyzTF6Ik=	NaN	NaN	NaN	NaN	Good Grief	51.0
1	Xumu+NiS6QYVvDS4t3SawJ7vT9hPKXmfbORLlNx8=	bhpMpSNoqxOIB+BWwPqU6jBt4DpCn3ayXnJqM=	NaN	NaN	NaN	NaN	Lords of Cardboard	1.0
2	Xumu+NiS6QYVvDS4t3SawJ7vT9hPKXmfbORLlNx8=	JNWhrC7zNt7BdMpslSKa4Mw+XVJYnXn3Epw7QgY=	225396.0	1259	Nas	52.0	Hip Hop Is Dead(Album Version (Edited))	1.0
3	Xumu+NiS6QYVvDS4t3SawJ7vT9hPKXmfbORLlNx8=	2A87zhtJTSWgD7gIzHsoIhe4DMzkb6uLzOJKHjN=	NaN	NaN	NaN	NaN	Disco Africa	1.0
4	FGfVqz1BRPwJyedr2gV78ziAIY/9SmYvIa+HCg=	3qm6XTZ8MOCU11k8FVAGH5G5UMKT3ZaWAG1oo2G=	187802.0	1011	Brett Young	52.0	Sleep Without You	56.0

5. Preprocessing of data

```
In [50]: df.shape
Out[50]: (1848575, 8)
```

```
In [51]: #checking null values
df.isnull().sum()
```

```
Out[51]:
```

	user_id	song_id	song_length	genre_ids	artist_name	language	name	listen_count
0	FGfVqz1BRPwJyedr2gV78ziAIY/9SmYvIa+HCg=	BBzumQNVLuHKEBOB7mAJuzok+JAlc2RygyzTF6Ik=	NaN	NaN	NaN	NaN	Good Grief	51.0
1	Xumu+NiS6QYVvDS4t3SawJ7vT9hPKXmfbORLlNx8=	bhpMpSNoqxOIB+BWwPqU6jBt4DpCn3ayXnJqM=	NaN	NaN	NaN	NaN	Lords of Cardboard	1.0
2	Xumu+NiS6QYVvDS4t3SawJ7vT9hPKXmfbORLlNx8=	JNWhrC7zNt7BdMpslSKa4Mw+XVJYnXn3Epw7QgY=	225396.0	1259	Nas	52.0	Hip Hop Is Dead(Album Version (Edited))	1.0
3	Xumu+NiS6QYVvDS4t3SawJ7vT9hPKXmfbORLlNx8=	2A87zhtJTSWgD7gIzHsoIhe4DMzkb6uLzOJKHjN=	NaN	NaN	NaN	NaN	Disco Africa	1.0
4	FGfVqz1BRPwJyedr2gV78ziAIY/9SmYvIa+HCg=	3qm6XTZ8MOCU11k8FVAGH5G5UMKT3ZaWAG1oo2G=	187802.0	1011	Brett Young	52.0	Sleep Without You	56.0

```
In [52]: #fill the null values
df['song_length'].fillna(0,inplace=True)
df['genre_ids'].fillna(0,inplace=True)
df['artist_name'].fillna('none',inplace=True)
df['language'].fillna(0,inplace=True)
df['name'].fillna('none',inplace=True)
df['listen_count'].fillna(0,inplace=True)
```

```
In [53]: #recheck null values
df.isnull().sum()
```

```
Out[53]:
```

	user_id	song_id	song_length	genre_ids	artist_name	language	name	listen_count
0	FGfVqz1BRPwJyedr2gV78ziAIY/9SmYvIa+HCg=	BBzumQNVLuHKEBOB7mAJuzok+JAlc2RygyzTF6Ik=	NaN	NaN	NaN	NaN	Good Grief	51.0
1	Xumu+NiS6QYVvDS4t3SawJ7vT9hPKXmfbORLlNx8=	bhpMpSNoqxOIB+BWwPqU6jBt4DpCn3ayXnJqM=	NaN	NaN	NaN	NaN	Lords of Cardboard	1.0
2	Xumu+NiS6QYVvDS4t3SawJ7vT9hPKXmfbORLlNx8=	JNWhrC7zNt7BdMpslSKa4Mw+XVJYnXn3Epw7QgY=	225396.0	1259	Nas	52.0	Hip Hop Is Dead(Album Version (Edited))	1.0
3	Xumu+NiS6QYVvDS4t3SawJ7vT9hPKXmfbORLlNx8=	2A87zhtJTSWgD7gIzHsoIhe4DMzkb6uLzOJKHjN=	NaN	NaN	NaN	NaN	Disco Africa	1.0
4	FGfVqz1BRPwJyedr2gV78ziAIY/9SmYvIa+HCg=	3qm6XTZ8MOCU11k8FVAGH5G5UMKT3ZaWAG1oo2G=	187802.0	1011	Brett Young	52.0	Sleep Without You	56.0

```
In [32]: print("Total no of songs:",len(df))
Total no of songs: 1848575
```

6. Subset of The Dataset

```
In [54]: df = df.head(10000)
df['song'] = df['name'].map(str) + " - " + df['artist_name']
```

7. Most Popular Songs in The Dataset

```
In [55]: song_gr = df.groupby(['song']).agg({'listen_count': 'count'}).reset_index()
grouped_sum = song_gr['listen_count'].sum()
song_gr['percentage'] = song_gr['listen_count'].div(grouped_sum)*100
song_gr.sort_values(['listen_count', 'song'], ascending = [0, 1])
```

```
Out[55]:
```

	song	listen_count	percentage
3195	告白瓶 - none	62	0.62
3650	醉倒分手 - 周湯豪 (NICKTHEREAL)	54	0.54
5229	謝謝妳愛我 (Thanks For Your Love) - 謝和弦 (R-chord)	53	0.53
427	Closer - The Chainsmokers	38	0.38
2884	你，好不好？ (How Have You Been?) - none	38	0.38
...
5617	حياناً(غOOD TO YOU) - 2NE1	1	0.01
5619	손수레에 왜 이러 - K.Will	1	0.01
5620	칼로리 송 (糖糖糖)라 믿고 싶잖아? Calorie Song - none	1	0.01
5621	헤어질땐 말아줘(I know) we will break up - none	1	0.01
5622	L i a r - ONE OK ROCK	1	0.01

5623 rows x 3 columns

8. Unique Users in The Dataset

```
In [56]: users = df['user_id'].unique()
print("The no. of unique users:", len(users))
The no. of unique users: 1622
```

8. Number of Unique Songs in The Dataset

```
In [57]: songs = df['song'].unique()
len(songs)
```

```
Out[57]: 5623
```

8. Create a Song Recommender

```
In [58]: train_data, test_data = train_test_split(df, test_size = 0.20, random_state=0)
print(train_data.shape)
```

```
Out[58]:
```

	msno	song_id	source_system_tab	source_screen_name	source_type	target
0	FGfVqz1BRPwJyedr2gV78ziAIY/9SmYvIa+HCg=	BBzumQNVLuHKEBOB7mAJuzok+JAlc2RygyzTF6Ik=	explore	Explore	online-playlist	1
1	Xumu+NiS6QYVvDS4t3SawJ7vT9hPKXmfbORLlNx8=	bhpMpSNoqxOIB+BWwPqU6jBt4DpCn3ayXnJqM=	my library	Local playlist more	local-playlist	1
2	Xumu+NiS6QYVvDS4t3SawJ7vT9hPKXmfbORLlNx8=	JNWhrC7zNt7BdMpslSKa4Mw+XVJYnXn3Epw7QgY=	my library	Local playlist more	local-playlist	1
3	Xumu+NiS6QYVvDS4t3SawJ7vT9hPKXmfbORLlNx8=	2A87zhtJTSWgD7gIzHsoIhe4DMzkb6uLzOJKHjN=	my library	Local playlist more	local-playlist	1
4	FGfVqz1BRPwJyedr2gV78ziAIY/9SmYvIa+HCg=	3qm6XTZ8MOCU11k8FVAGH5G5UMKT3ZaWAG1oo2G=	explore	Explore	online-playlist	1

```
Out[58]:
```

	source_screen_name	source_type	target	song_length	genre_ids	language	name
0	Explore	online-playlist	1	NaN	NaN	NaN	Good Grief
1	Local playlist more	local-playlist	1	NaN	NaN	NaN	Lords of Cardboard
2	Local playlist more	local-playlist	1	225396.0	1259	NaN	Hip Hop Is Dead(Album Version (Edited))
3	Local playlist more	local-playlist	1	NaN	NaN	NaN	Disco Africa
4	Explore	online-playlist	1	187802.0	1011	NaN	Sleep Without You

```
Out[58]:
```

	artist_name	language	name
0	NaN	NaN	Good Grief
1	NaN	NaN	Lords of Cardboard
2	Nas	52.0	Hip Hop Is Dead(Album Version (Edited))
3	NaN	NaN	Disco Africa
4	Brett Young	52.0	Sleep Without You

```
Out[58]:
```

	isrc
0	GBUHT1682854
1	USC99910183
2	USUW70518761
3	GBUQH1006063
4	QW3E21686693

9. Data Visualization

```
In [59]: plt.figure(figsize=(10,10))
sns.countplot(x='source_system_tab', hue='source_type', data=train)
```

```
Out[59]: <AxesSubplot:xlabel='source_system_tab', ylabel='count'>
```

```
In [60]: plt.figure(figsize=(10,10))
sns.countplot(x='source_system_tab', hue='target', data=train)
```

```
Out[60]: <AxesSubplot:xlabel='source_system_tab', ylabel='count'>
```

```
In [61]: plt.figure(figsize=(10,10))
plt.xticks(rotation=90)
sns.countplot(x='source_screen_name', hue='target', data=train)
```

```
Out[61]: <AxesSubplot:xlabel='source_screen_name', ylabel='count'>
```

```
In [62]: plt.figure(figsize=(10,10))
plt.xticks(rotation=90)
sns.countplot(x='source_type', hue='source_type', data=train)
```

```
Out[62]: <AxesSubplot:xlabel='source_type', ylabel='count'>
```

```
In [63]: plt.figure(figsize=(10,10))
plt.xticks(rotation=90)
sns.countplot(x='registered_via', hue='registered_via', data=members)
```

```
Out[63]: <AxesSubplot:xlabel='registered_via', ylabel='count'>
```

```
In [64]: ntr = 7000
nts = 3000
names = df['song_id','source_system_tab','source_screen_name','source_type','target']
test1 = pd.read_csv('train.csv',names=names,skiprows=ntr,nrows=nts)
```

```
In [67]: test = test1.drop(['target'],axis=1)
ytr = np.array(test1['target'])
```

```
In [68]: test_name = ['id','msno','song_id','source_system_tab','source_screen_name','source_type']
test['id']=np.arange(nts)
test = test[test_name]
```

```
In [71]: members['registration_year'] = members['registration_init_time'].apply(lambda x: int(str(x)[0:4]))
members['registration_month'] = members['registration_init_time'].apply(lambda x: int(str(x)[4:6]))
members['registration_date'] = members['registration_init_time'].apply(lambda x: int(str(x)[6:8]))
```

```
In [72]: members['expiration_year'] = members['expiration_date'].apply(lambda x: int(str(x)[0:4]))
members['expiration_month'] = members['expiration_date'].apply(lambda x: int(str(x)[4:6]))
members['expiration_date'] = members['expiration_date'].apply(lambda x: int(str(x)[6:8]))
members = members.drop(['registration_init_time'],axis=1)
```

```
In [73]: members_cols = members.columns
train = train.merge(members[members_cols], on='msno', how='left')
test = test.merge(members[members_cols], on='msno', how='left')
```

```
In [74]: train = train.fillna(-1)
test = test.fillna(-1)
```

```
In [75]: import gc
del members, songs; gc.collect();
import warnings
warnings.filterwarnings('ignore')
```

```
In [76]: cols = list(train.columns)
cols.remove('target')
```

```
In [81]: from sklearn.preprocessing import LabelEncoder
```

```
In [83]: unique_songs = range(max(train['song_id'].max()), test['song_id'].max())
song_popularity = pd.DataFrame({'song_id': unique_songs, 'popularity':0})
train_sorted = train.sort_values('song_id')
train_sorted.reset_index(drop=True, inplace=True)
test_sorted = test.sort_values('song_id')
test_sorted.reset_index(drop=True, inplace=True)
```

10. data Cleaning

```
In [66]: ntr = 7000
nts = 3000
names = df['song_id','source_system_tab','source_screen_name','source_type','target']
test1 = pd.read_csv('train.csv',names=names,skiprows=ntr,nrows=nts)
```

```
In [67]: test = test1.drop(['target'],axis=1)
ytr = np.array(test1['target'])
```

```
In [68]: test_name = ['id','msno','song_id','source_system_tab','source_screen_name','source_type']
test['id']=np.arange(nts)
test = test[test_name]
```

```
In [71]: members['registration_year'] = members['registration_init_time'].apply(lambda x: int(str(x)[0:4]))
members['registration_month'] = members['registration_init_time'].apply(lambda x: int(str(x)[4:6]))
members['registration_date'] = members['registration_init_time'].apply(lambda x: int(str(x)[6:8]))
```