

In [1]:

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from surprise import Dataset, Reader
from surprise import accuracy
from surprise.model_selection import train_test_split
from surprise.prediction_algorithms import knns
from surprise.similarities import cosine, pearson
from surprise import accuracy
```

In [2]:

```
#Caricamento dataset
song_df_1 = pd.read_csv('triplets_file.csv')
song_df_1.columns = ['user_id', 'song_id', 'listen_count']

song_df_2 = pd.read_csv('song_data.csv')
song_df_2.drop_duplicates(['song_id'], inplace=True)

#Unione dataset
df_songs = pd.merge(song_df_1, song_df_2, on='song_id', how='left')
```

In [3]:

```
# Numero canzoni ascoltate da ogni utente
user_counts = df_songs.groupby('user_id')['song_id'].count()

# Utenti che hanno ascolato almeno 30 canzoni
user_ten_id = user_counts[user_counts >= 30].index.to_list()
```

In [4]:

```
# Numero di utenti che hanno ascoltato ogni canzone
song_counts = df_songs.groupby('song_id')['user_id'].count()

# Canzoni ascoltate da almeno 200 utenti
song_ten_u_id = song_counts[song_counts >= 200].index.to_list()
```

In [5]:

```
df_songs_ridotto = df_songs[df_songs['user_id'].isin(user_ten_id)] & (df_songs['song_id'].isin(song_ten_id))].reset_index(drop=True)
df_songs_ridotto
```

Out [5]:

	user_id	song_id	listen_count	title	release	artist_name	year
0	b0344d0636cc32127f638f9d413e97c32313a9	SOBXHDL12AB1CD04C0	1	Stronger	Graduation	Kanye West	2007
1	b0344d0636cc32127f638f9d413e97c32313a9	SOBVHAI12AB10BF1D	1	Constellations	In Between Dreams	Jack Johnson	2005
2	b0344d0636cc32127f638f9d413e97c32313a9	SOOACRL12AB13C273	1	Learn To Fly	There Is Nothing Left To Lose	Foo Fighters	1999
3	b0344d0636cc32127f638f9d413e97c32313a9	SOOXRTY12AB0180F3B	1	Paper Gangsta	The Fame Monster	Lady Gaga	2008
4	b0344d0636cc32127f638f9d413e97c32313a9	SOFRTQD12AB1C233C0	1	Sehr kosmesch	Musk von Harmonia	Harmonia	0
...	...	...	...	...	...	...	...
336425	db8b4ec88093773a6e02263c1a01d13b76a92	SOH2PHK12AS8AT7CAE	4	Te Amo	Rated R	Rihanna	2009
336426	db8b4ec88093773a6e02263c1a01d13b76a92	SOIDDNJ12AC3071B90	1	That Should Be Me	My Worlds	Justin Bieber	2010
336427	db8b4ec88093773a6e02263c1a01d13b76a92	SOIOESO12AGD4F21D	4	Unwell (Album Version)	More Than You Think You Are	matchbox twenty	2003
336428	db8b4ec88093773a6e02263c1a01d13b76a92	SOJKQSF12AGD4F5EE9	3	What I've Done (Album Version)	What I've Done	Linkin Park	2007
336429	db8b4ec88093773a6e02263c1a01d13b76a92	SOJUXGA12AC361885C	1	Up	My Worlds	Justin Bieber	2010

336430 rows × 7 columns

In [52]:

```
# Creazione scala valutazione
scala = [0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 2214]

df_songs_ridotto[['listen_count']] = pd.cut(df_songs_ridotto['listen_count'], bins=scala, labels=[1, 2, 3, 4, 5, 6, 7, 8, 9, 10])

listen_counts = pd.DataFrame(df_songs_ridotto.groupby('listen_count').size(), columns=['count']).reset_index(drop=False)

plt.figure(figsize=(10, 8))
sns.barplot(x='listen_count', y='count', palette='Set3', data=listen_counts)
plt.gca().spines['top'].set_visible(False)
plt.gca().spines['right'].set_visible(False)
plt.show();
```

Out [52]:

In [7]:

```
reader = Reader(rating_scale=(1, 10))
data = Dataset.load_from_df(df_songs_ridotto[['user_id', 'title', 'listen_count']], reader)
```

In [8]:

```
train_set, test_set = train_test_split(data, test_size=.25)
```

In [ ]:

```
# Funzione get top n predictions di Surprise
from collections import defaultdict
def get_top_n(predictions, n=10):

    # First map the predictions to each user.
    top_n = defaultdict(list)
    for uid, iid, r_ui, est, _ in predictions:
        top_n[uid].append((iid, est, r_ui))

    # Then sort the predictions for each user and retrieve the k highest ones.
    for uid, user_ratings in top_n.items():
        # print(user_ratings)
        user_ratings = sorted(user_ratings, key=lambda x: x[1], reverse=True)
        top_n[uid] = user_ratings[:n]

    return top_n
```

KNN BASIC

In [9]:

```
# Definizione indice similarità
sim_cos = {'name': 'cosine', 'user_based': False}
```

In [10]:

```
# Addestramento del modello
kbasic = knns.KNNBasic(sim_options=sim_cos)
kbasic.fit(train_set)
```

Out [10]:

Computing the cosine similarity matrix...

Done computing similarity matrix.

<surprise.prediction\_algorithms.knns.KNNBasic at 0x27807b9910>

In [11]:

```
# Test del modello
predictions = kbasic.test(test_set)
```

In [12]:

```
# Metriche di valutazione
rmse = accuracy.rmse(predictions)
mae = accuracy.mae(predictions)

RMSE: 2.3775
MAE: 1.6536
```

In [14]:

```
# Generazione top n previsioni
top_pred_knn = get_top_n(predictions, 10)
```

In [15]:

```
top_pred_knn_df = pd.DataFrame([(iid, pair[0], pair[1]) for iid, row in top_pred_knn.items() for pair in row], columns=['user_id', 'title', 'rat_pred'])
top_pred_knn_df
```

Out [15]:

	user_id	title	rat_pred
0	ed5296869661261650290a01020c0b7c6aa197b	Hochmah (Hochmah)	2.362288
1	ed5296869661261650290a01020c0b7c6aa197b	Kyrgyz Me	2.218035
2	ed5296869661261650290a01020c0b7c6aa197b	You Belong With Me	2.201269
3	K0aa3b73277264c5ba0594826b294858986969	Ma	2.167763
4	K0aa3b73277264c5ba0594826b294858986969	The District Sleeps Alone Tonight (Album)	1.984358
...	...	...	...
73785	5b6a4212405e548e45a35ec315fc11480253e	The Funeral (Album Version)	1.293390
73786	D90c0e478859592030ed132dc7aa1e012c0de	Uha Confusion	2.509336
73787	750cbe758703da09609c097ae2175bede6c76ae	Burden In My Hand	1.450015
73788	7d1c0b9d9e8696c05a42c482c4354006eb	Revelry	1.943061
73789	530cc0c19114061072e33840a679cd4ba3c95a	Jamaica Roots II (Agrona E Sempre)	3.009972

73790 rows × 4 columns

In [16]:

```
top_pred_knn_df.loc[top_pred_knn_df['user_id'] == 'b048f21af05e7467f187bf9f9d413e97c32313a9', ('title', 'rat_pred')]
```

Out [16]:

	title	rat_pred
11402	Life In Technicolor 4	5.237305
11403	Clocks	5.139316
11404	Don't Panic	5.063960
11405	Fix You	4.984168
11406	Lost	4.881390
11407	Every Little Thing She Does Is Magic	4.890238
11408	The Ballad of Michael Valentine	4.890616
11409	I Kissed A Girl	4.785901
11410	Trouble	4.754914
11411	The District Sleeps Alone Tonight (Album)	4.717154

KNN MEANS MODEL

In [17]:

```
# Addestramento del modello
KNNMeans = knns.KNNWithMeans(sim_options=sim_cos)
KNNMeans.fit(train_set)
```

Out [17]:

Computing the cosine similarity matrix...

Done computing similarity matrix.

<surprise.prediction\_algorithms.knns.KNNWithMeans at 0x278d55dfde>

In [18]:

```
# Test del modello
predictions_mean = KNNMeans.test(test_set)
```

In [19]:

```
# Metriche di valutazione
rmse = accuracy.rmse(predictions_mean)
mae = accuracy.mae(predictions_mean)

RMSE: 2.2138
MAE: 1.6075
```

In [20]:

```
# Generazione top n previsioni
top_pred_mean = get_top_n(predictions_mean, 10)
```

In [21]:

```
top_pred_mean['b048f21af05e7467f187bf9f9d413e97c32313a9']
```

Out [21]:

```
[('Clocks', 5.284273965107248, 10.0),
('Life In Technicolor II', 5.172153169602957, 4.0),
('Lost', 5.078141238192131, 8.0),
('Fix You', 5.08931313580844, 7.0),
('Elephant Gun', 4.961543938828061, 1.0),
('I Kissed A Girl', 4.944127378332749, 3.0),
('Don't Panic', 4.91495116323143, 10.0),
('Trouble', 4.761759794603866, 1.0),
('Every Little Thing She Does Is Magic', 4.682350235045847, 5.0),
('The District Sleeps Alone Tonight (Album)', 4.363397786885078, 1.0)]
```

In [22]:

```
top_pred_mean_df = pd.DataFrame([(iid, pair[0], pair[1]) for iid, row in top_pred_mean.items() for pair in row], columns=['user_id', 'title', 'rat_pred'])
top_pred_mean_df.loc[top_pred_mean_df['user_id'] == 'b048f21af05e7467f187bf9f9d413e97c32313a9', ('title', 'rat_pred')]
```

Out [22]:

	title	rat_pred
11402	Clocks	5.204273
11403	Life In Technicolor 4	5.172153
11404	Lost	5.078141
11405	Fix You	5.009313
11406	Elephant Gun	4.961544
11407	I Kissed A Girl	4.944117
11408	Don't Panic	4.914951
11409	Trouble	4.761715
11410	Every Little Thing She Does Is Magic	4.662350
11411	The District Sleeps Alone Tonight (Album)	4.363398

KNN Z SCORE MODEL

In [23]:

```
# Addestramento del modello
KNN_Z = knns.KNNWithZScore(sim_options=sim_cos)
KNN_Z.fit(train_set)
```

Out [23]:

Computing the cosine similarity matrix...

Done computing similarity matrix.

<surprise.prediction\_algorithms.knns.KNNWithZScore at 0x278da8ee978>

In [24]:

```
# Test del modello
predictions_Z = KNN_Z.test(test_set)
```

In [25]:

```
# Metriche di valutazione
rmse = accuracy.rmse(predictions_Z)
mae = accuracy.mae(predictions_Z)

RMSE: 2.3251
MAE: 1.5997
```

In [26]:

```
# Generazione top n previsioni
top_pred_Z = get_top_n(predictions_Z, 10)
```

In [27]:

```
top_pred_Z_df = pd.DataFrame([(iid, pair[0], pair[1]) for iid, row in top_pred_Z.items() for pair in row], columns=['user_id', 'title', 'rat_pred'])
top_pred_Z_df.loc[top_pred_Z_df['user_id'] == 'b048f21af05e7467f187bf9f9d413e97c32313a9', ('title', 'rat_pred')]
```

Out [27]:

	title	rat_pred
11402	Elephant Gun	5.571085
11403	Lost	5.440097
11404	Clocks	5.354612
11405	I Kissed A Girl	5.281291
11406	Life In Technicolor II	5.156917
11407	Fix You	5.015369
11408	Don't Panic	4.916655
11409	Trouble	4.850568
11410	Every Little Thing She Does Is Magic	4.462021
11411	We Will Become Silhouettes (Album)	4.061780

USER BASED

BASIC KNN

In [28]:

```
# Definizione indice similarità
sim_pearson = {'name': 'pearson', 'user_based': True}
```

In [29]:

```
# Addestramento del modello
KNN_basic_u = knns.KNNBasic(sim_options=sim_pearson)
KNN_basic_u.fit(train_set)
```

Out [29]:

Computing the pearson similarity matrix...

Done computing similarity matrix.

<surprise.prediction\_algorithms.knns.KNNBasic at 0x2782c5d3a0>

In [30]:

```
# Test del modello
predictions_basic_u = KNN_basic_u.test(test_set)
```

In [31]:

```
# Metriche di valutazione
rmse = accuracy.rmse(predictions_basic_u)
mae = accuracy.mae(predictions_basic_u)

RMSE: 2.5704
MAE: 1.9333
```

In [32]:

```
# Generazione top n previsioni
top_pred_basic_u = get_top_n(predictions_basic_u, 10)
```

In [33]:

```
top_pred_basic_u_df = pd.DataFrame([(iid, pair[0], pair[1]) for iid, row in top_pred_basic_u.items() for pair in row], columns=['user_id', 'title', 'rat_pred'])
top_pred_basic_u_df.loc[top_pred_basic_u_df['user_id'] == 'b048f21af05e7467f187bf9f9d413e97c32313a9', ('title', 'rat_pred')]
```

Out [33]:

	title	rat_pred
11402	Life In Technicolor II	2.874089
11403	Lost	2.869134
11404	Fix You	2.625000
11405	Clocks	2.500000
11406	I Kissed A Girl	2.485766
11407	Elephant Gun	2.447237
11408	Don't Panic	2.213656
11409	Trouble	2.206178
11410	Every Little Thing She Does Is Magic	1.801824
11411	The Ballad of Michael Valentine	1.696871

WITH MEANS

In [34]:

```
# Addestramento del modello
KNNMeans_u = knns.KNNWithMeans(sim_options=sim_pearson)
KNNMeans_u.fit(train_set)
```

Out [34]:

Computing the pearson similarity matrix...

Done computing similarity matrix.

<surprise.prediction\_algorithms.knns.KNNWithMeans at 0x278df67990d>

In [35]:

```
# Test del modello
predictions_mean_u = KNNMeans_u.test(test_set)
```

In [36]:

```
# Metriche di valutazione
rmse = accuracy.rmse(predictions_mean_u)
mae = accuracy.mae(predictions_mean_u)

RMSE: 2.3704
MAE: 1.6308
```

In [37]:

```
# Generazione top n previsioni
top_pred_mean_u = get_top_n(predictions_mean_u, 10)
```

In [38]:

```
top_pred_mean_u_df = pd.DataFrame([(iid, pair[0], pair[1]) for iid, row in top_pred_mean_u.items() for pair in row], columns=['user_id', 'title', 'rat_pred'])
top_pred_mean_u_df.loc[top_pred_mean_u_df['user_id'] == 'b048f21af05e7467f187bf9f9d413e97c32313a9', ('title', 'rat_pred')]
```

Out [38]:

	title	rat_pred
11402	Lost	5.377371
11403	Fix You	4.693215
11404	Life In Technicolor II	4.663215
11405	Clocks	4.415978
11406	I Kissed A Girl	4.310783
11407	Don't Panic	4.260575
11408	Elephant Gun	4.042171
11409	The Ballad of Michael Valentine	3.865661
11410	We Will Become Silhouettes (Album)	3.865249
11411	Every Little Thing She Does Is Magic	3.700643

Z SCORE KNN

In [39]:

```
# Addestramento del modello
KNN_Z_u = knns.KNNWithZScore(sim_options=sim_pearson)
KNN_Z_u.fit(train_set)
```

Out [39]:

Computing the pearson similarity matrix...

Done computing similarity matrix.

<surprise.prediction\_algorithms.knns.KNNWithZScore at 0x278da1ff2b8>

In [40]:

```
# Test del modello
predictions_Z_u = KNN_Z_u.test(test_set)
```

In [41]:

```
# Metriche di valutazione
rmse = accuracy.rmse(predictions_Z_u)
mae = accuracy.mae(predictions_Z_u)

RMSE: 2.3645
MAE: 1.6052
```

In [42]:

```
# Generazione top n previsioni
top_pred_Z_u = get_top_n(predictions_Z_u, 10)
```

In [43]:

```
top_pred_Z_u_df = pd.DataFrame([(iid, pair[0], pair[1]) for iid, row in top_pred_Z_u.items() for pair in row], columns=['user_id', 'title', 'rat_pred'])
top_pred_Z_u_df.loc[top_pred_Z_u_df['user_id'] == 'b048f21af05e7467f187bf9f9d413e97c32313a9', ('title', 'rat_pred')]
```

Out [43]:

	title	rat_pred
11402	Lost	5.377371
11403	Fix You	4.693215
11404	Life In Technicolor II	4.663215
11405	Clocks	4.415978
11406	I Kissed A Girl	4.310783
11407	Don't Panic	4.260575
11408	Elephant Gun	4.042171
11409	The Ballad of Michael Valentine	3.865661
11410	We Will Become Silhouettes (Album)	3.865249
11411	Every Little Thing She Does Is Magic	3.700643

TABELLA CONFRONTO

In [50]:

```
accuracy_tab = {'RMSE': [2.3775, 2.3138, 2.3251, 2.5704, 2.3784, 2.3645],
                'MAE': [1.6536, 1.6075, 1.5997, 1.9333, 1.6308, 1.6052]
                }

accuracy_tab_df = pd.DataFrame(accuracy_tab, index=['Basic_U', 'WithMeans_U', 'Zscore_U', 'Basic_U', 'WithMeans_U', 'Zscore_U'])
display(accuracy_tab_df)
```

Out [50]:

	RMSE	MAE
Basic_U	2.3775	1.6536
WithMeans_U	2.3138	1.6075
Zscore_U	2.3251	1.5997
Basic_U	2.5704	1.9333
WithMeans_U	2.3784	1.6308
Zscore_U	2.3645	1.6052

In [ ]: