

Task 2

```
In [1]: import os

import matplotlib.pyplot as plt
import numpy as np
import pandas as pd
from IPython.display import display, Markdown

from surprise import Dataset
from surprise import Reader
from surprise.model_selection import train_test_split, cross_validate
from surprise import accuracy
from surprise import SVD, KNNBasic
from sklearn.metrics import confusion_matrix, precision_score, recall_score, classification_report
import warnings
warnings.filterwarnings('ignore')
```

Part-2

```
In [2]: print('Downloaded files are..')
os.listdir('datasets/MovieRating')
```

Downloaded files are..

```
Out[2]: ['credits.csv',
'keywords.csv',
'links.csv',
'links_small.csv',
'movies_metadata.csv',
'ratings.csv',
'ratings_small.csv']
```

```
In [3]: data = pd.read_csv('datasets/MovieRating/ratings_small.csv')
data.sample(5)
```

```
Out[3]:
```

	userId	movieId	rating	timestamp
53751	388	423	3.0	946534602
32686	236	1206	4.5	1109968692
79563	547	4881	3.5	1053907106
29739	213	7458	2.5	1462634314
99999	671	6268	2.5	1065579370

Part-3

```
In [4]: data = data.drop('timestamp', axis=1)
```

```
In [5]: data.shape
```

```
Out[5]: (100004, 3)
```

```
In [6]: data.sample(5)
```

```
Out[6]:
```

	userId	movieId	rating
54440	388	38886	4.5
78246	544	89492	5.0
58777	427	2759	2.0
5397	30	1951	4.0
39447	287	68791	5.0

```
In [7]: reader = Reader(rating_scale=(1, 5))
ratings = Dataset.load_from_df(data, reader)
```

```
In [8]: type(ratings)
```

```
Out[8]: surprise.dataset.DatasetAutoFolds
```

3(c) Compute the average MAE and RMSE of the Probabilistic Matrix Factorization (PMF), User based Collaborative Filtering, Item based Collaborative Filtering, under the 5-folds cross-validation (10 points)

Use Probabilistic Matrix Factorization(PMF)

```
In [9]: model_pmf = SVD()
model_pmf_cv = cross_validate(model_pmf, ratings, measures=['RMSE', 'MAE'], cv = 5, verb
```

Evaluating RMSE, MAE of algorithm SVD on 5 split(s).

	Fold 1	Fold 2	Fold 3	Fold 4	Fold 5	Mean	Std
RMSE (testset)	0.8920	0.9077	0.8915	0.8953	0.8979	0.8969	0.0059
MAE (testset)	0.6875	0.6988	0.6856	0.6923	0.6910	0.6910	0.0046
Fit time	0.89	0.94	0.97	0.94	0.92	0.93	0.03
Test time	0.11	0.17	0.12	0.17	0.11	0.14	0.03

```
In [10]: avg_pmf_rmse = np.average(model_pmf_cv['test_rmse'])
avg_pmf_mae = np.average(model_pmf_cv['test_mae'])
print('Average of RMSE for Probabilistic Matrix Factorization(PMF) = ', avg_pmf_rmse)
print('Average of MAE for Probabilistic Matrix Factorization(PMF) = ', avg_pmf_mae)
```

Average of RMSE for Probabilistic Matrix Factorization(PMF) = 0.8968896033085734
Average of MAE for Probabilistic Matrix Factorization(PMF) = 0.6910420967843592

Use User-based collaborative filtering (UCF)

```
In [11]: sim_options = {'name': 'cosine', 'user_based': True}
model_ucf = KNNBasic(sim_options=sim_options)
model_ucf_cv = cross_validate(model_ucf, ratings, measures=['RMSE', 'MAE'], cv = 5, verb
```

Computing the cosine similarity matrix...
Done computing similarity matrix.
Computing the cosine similarity matrix...
Done computing similarity matrix.
Computing the cosine similarity matrix...
Done computing similarity matrix.
Computing the cosine similarity matrix...
Done computing similarity matrix.
Computing the cosine similarity matrix...

Done computing similarity matrix.
Evaluating RMSE, MAE of algorithm KNNBasic on 5 split(s).

	Fold 1	Fold 2	Fold 3	Fold 4	Fold 5	Mean	Std
RMSE (testset)	0.9864	0.9996	0.9804	0.9962	0.9962	0.9918	0.0072
MAE (testset)	0.7634	0.7718	0.7591	0.7692	0.7711	0.7670	0.0049
Fit time	0.52	0.45	0.42	0.44	0.56	0.48	0.05
Test time	1.34	1.30	1.31	1.73	1.39	1.42	0.16

```
In [12]: avg_ucf_rmse = np.average(model_ucf_cv['test_rmse'])
avg_ucf_mae = np.average(model_ucf_cv['test_mae'])
text_1 = f"Average of RMSE for User-based collaborative filtering (UCF) = {avg_ucf_rmse}"
display(Markdown(text_1))
```

Average of RMSE for User-based collaborative filtering (UCF) = 0.9917525439607067

Average of MAE for User-based collaborative filtering (UCF) = 0.7669527198626307

Use Item-based collaborative filtering (ICF)

```
In [13]: sim_options = {'name': 'cosine', 'user_based': False}
model_icf = KNNBasic(sim_options=sim_options)
model_icf_cv = cross_validate(model_icf, ratings, measures=['RMSE', 'MAE'], cv = 5, verb
```

Computing the cosine similarity matrix...
Done computing similarity matrix.
Computing the cosine similarity matrix...
Done computing similarity matrix.
Computing the cosine similarity matrix...
Done computing similarity matrix.
Computing the cosine similarity matrix...
Done computing similarity matrix.
Computing the cosine similarity matrix...
Done computing similarity matrix.
Computing the cosine similarity matrix...
Done computing similarity matrix.
Evaluating RMSE, MAE of algorithm KNNBasic on 5 split(s).

	Fold 1	Fold 2	Fold 3	Fold 4	Fold 5	Mean	Std
RMSE (testset)	1.0017	0.9912	0.9936	0.9961	0.9840	0.9933	0.0058
MAE (testset)	0.7779	0.7711	0.7736	0.7782	0.7685	0.7739	0.0038
Fit time	5.84	5.66	6.29	6.60	6.29	6.14	0.34
Test time	5.92	7.03	9.24	8.04	8.26	7.70	1.13

```
In [14]: avg_icf_rmse = np.average(model_icf_cv['test_rmse'])
avg_icf_mae = np.average(model_icf_cv['test_mae'])
print('Average of RMSE for Item-based collaborative filtering (ICF) =', avg_icf_rmse)
print('Average of MAE for Item-based collaborative filtering (ICF) =', avg_icf_mae)
```

Average of RMSE for Item-based collaborative filtering (ICF) = 0.9933200887032789

Average of MAE for Item-based collaborative filtering (ICF) = 0.7738791298206834

3(d) Compare the average (mean) performances of User-based collaborative filtering, item-based collaborative filtering, PMF with respect to RMSE and MAE. Which ML model is the best in the movie rating data? (10 points)

```
In [15]: text = f"""|Algo|Mean RMSE|Mean MAE|
|----|-----|-----|
|Probability Matrix Factorization|{avg_pmf_rmse}|{avg_pmf_mae}|
|User-based collaborative filtering|{avg_ucf_rmse}|{avg_ucf_mae}|
|Item-based collaborative filtering|{avg_icf_rmse}|{avg_icf_mae}|"""
display(Markdown(text))
```

Algo	Mean RMSE	Mean MAE
Probability Matrix Factorization	0.8968896033085734	0.6910420967843592
User-based collaborative filtering	0.9917525439607067	0.7669527198626307
Item-based collaborative filtering	0.9933200887032789	0.7738791298206834

Looking at the table, we see User-based collaborative filtering and Item-based collaborative filtering are pretty similar in performance.

But on pure numbers perspective, Item-based collaborative filtering is the better model for movie rating data for both RMSE and MAE

3(e) Examine how the cosine, MSD (Mean Squared Difference), and Pearson similarities impact the performances of User based Collaborative Filtering and Item based Collaborative Filtering. Plot your results. Is the impact of the three metrics on User based Collaborative Filtering consistent with the impact of the three metrics on Item based Collaborative Filtering? (10 points)

In [16]: *# Item-based collaborative filtering*

```
sim_options_cosine = {
    "name": 'cosine',
    'user_based': False
}

sim_options_msd = {
    "name": 'msd',
    'user_based': False
}

sim_options_pearson = {
    "name": 'pearson',
    'user_based': False
}
```

In [17]:

```
model_icf_cosine = KNNBasic(sim_options=sim_options_cosine)
model_icf_cosine_cv = cross_validate(algo=model_icf_cosine, data=ratings, measures=['RMS

model_icf_msd = KNNBasic(sim_options=sim_options_msd)
model_icf_msd_cv = cross_validate(algo=model_icf_msd, data=ratings, measures=['RMSE'], c

model_icf_pearson = KNNBasic(sim_options=sim_options_pearson)
model_icf_pearson_cv = cross_validate(algo=model_icf_pearson, data=ratings, measures=['R

avg_model_icf_cosine_cv = np.average(model_icf_cosine_cv['test_rmse'])
avg_model_icf_msd_cv = np.average(model_icf_msd_cv['test_rmse'])
avg_model_icf_pearson_cv = np.average(model_icf_pearson_cv['test_rmse'])
```

```
Computing the cosine similarity matrix...
Done computing similarity matrix.
Computing the cosine similarity matrix...
Done computing similarity matrix.
Computing the cosine similarity matrix...
Done computing similarity matrix.
Computing the cosine similarity matrix...
Done computing similarity matrix.
```

Computing the cosine similarity matrix...
 Done computing similarity matrix.
 Evaluating RMSE of algorithm KNNBasic on 5 split(s).

	Fold 1	Fold 2	Fold 3	Fold 4	Fold 5	Mean	Std
RMSE (testset)	0.9919	1.0016	0.9940	0.9936	0.9971	0.9956	0.0034
Fit time	5.82	5.70	5.64	5.51	5.53	5.64	0.12
Test time	5.78	6.02	6.66	5.95	6.54	6.19	0.35

Computing the msd similarity matrix...

Done computing similarity matrix.

Computing the msd similarity matrix...

Done computing similarity matrix.

Computing the msd similarity matrix...

Done computing similarity matrix.

Computing the msd similarity matrix...

Done computing similarity matrix.

Computing the msd similarity matrix...

Done computing similarity matrix.

Evaluating RMSE of algorithm KNNBasic on 5 split(s).

	Fold 1	Fold 2	Fold 3	Fold 4	Fold 5	Mean	Std
RMSE (testset)	0.9402	0.9290	0.9282	0.9370	0.9388	0.9346	0.0050
Fit time	3.53	3.74	3.70	3.78	3.69	3.69	0.09
Test time	7.47	8.20	8.76	8.26	8.44	8.22	0.43

Computing the pearson similarity matrix...

Done computing similarity matrix.

Computing the pearson similarity matrix...

Done computing similarity matrix.

Computing the pearson similarity matrix...

Done computing similarity matrix.

Computing the pearson similarity matrix...

Done computing similarity matrix.

Computing the pearson similarity matrix...

Done computing similarity matrix.

Evaluating RMSE of algorithm KNNBasic on 5 split(s).

	Fold 1	Fold 2	Fold 3	Fold 4	Fold 5	Mean	Std
RMSE (testset)	0.9861	0.9909	0.9914	0.9829	0.9922	0.9887	0.0036
Fit time	10.55	9.28	13.31	12.10	9.74	11.00	1.50
Test time	8.16	11.10	12.54	10.02	9.43	10.25	1.48

In [18]: *# User-based collaborative filtering*

```
sim_options_cosine = {
    "name": 'cosine',
    'user_based': True
}

sim_options_msd = {
    "name": 'msd',
    'user_based': True
}

sim_options_pearson = {
    "name": 'pearson',
    'user_based': True
}
```

In [19]:

```
model_ucf_cosine = KNNBasic(sim_options=sim_options_cosine)
model_ucf_cosine_cv = cross_validate(algo=model_ucf_cosine, data=ratings, measures=['RMS

model_ucf_msd = KNNBasic(sim_options=sim_options_msd)
model_ucf_msd_cv = cross_validate(algo=model_ucf_msd, data=ratings, measures=['RMSE'], c

model_ucf_pearson = KNNBasic(sim_options=sim_options_pearson)
model_ucf_pearson_cv = cross_validate(algo=model_ucf_pearson, data=ratings, measures=['R
```

```
avg_model_ucf_cosine_cv = np.average(model_ucf_cosine_cv['test_rmse'])
avg_model_ucf_msd_cv = np.average(model_ucf_msd_cv['test_rmse'])
avg_model_ucf_pearson_cv = np.average(model_ucf_pearson_cv['test_rmse'])
```

```
Computing the cosine similarity matrix...
Done computing similarity matrix.
Computing the cosine similarity matrix...
Done computing similarity matrix.
Computing the cosine similarity matrix...
Done computing similarity matrix.
Computing the cosine similarity matrix...
Done computing similarity matrix.
Computing the cosine similarity matrix...
Done computing similarity matrix.
Evaluating RMSE of algorithm KNNBasic on 5 split(s).
```

	Fold 1	Fold 2	Fold 3	Fold 4	Fold 5	Mean	Std
RMSE (testset)	0.9982	0.9892	0.9917	0.9986	0.9934	0.9942	0.0037
Fit time	0.55	0.59	0.56	0.63	0.58	0.58	0.03
Test time	1.90	2.01	2.00	2.09	1.94	1.99	0.06

```
Computing the msd similarity matrix...
Done computing similarity matrix.
Computing the msd similarity matrix...
Done computing similarity matrix.
Computing the msd similarity matrix...
Done computing similarity matrix.
Computing the msd similarity matrix...
Done computing similarity matrix.
Computing the msd similarity matrix...
Done computing similarity matrix.
Evaluating RMSE of algorithm KNNBasic on 5 split(s).
```

	Fold 1	Fold 2	Fold 3	Fold 4	Fold 5	Mean	Std
RMSE (testset)	0.9622	0.9732	0.9675	0.9643	0.9689	0.9672	0.0038
Fit time	0.23	0.26	0.27	0.28	0.27	0.26	0.02
Test time	1.97	2.04	1.89	1.77	1.98	1.93	0.09

```
Computing the pearson similarity matrix...
Done computing similarity matrix.
Computing the pearson similarity matrix...
Done computing similarity matrix.
Computing the pearson similarity matrix...
Done computing similarity matrix.
Computing the pearson similarity matrix...
Done computing similarity matrix.
Computing the pearson similarity matrix...
Done computing similarity matrix.
Evaluating RMSE of algorithm KNNBasic on 5 split(s).
```

	Fold 1	Fold 2	Fold 3	Fold 4	Fold 5	Mean	Std
RMSE (testset)	0.9875	1.0076	0.9950	0.9997	0.9980	0.9976	0.0065
Fit time	0.67	0.67	0.84	0.78	0.67	0.73	0.07
Test time	1.95	2.09	2.14	2.19	1.83	2.04	0.13

```
In [20]: final_res = {
    'RMSE Cosine Similarity - UCF': avg_model_ucf_cosine_cv,
    'RMSE Pearson Similarity - UCF': avg_model_ucf_pearson_cv,
    'RMSE MSD Similarity - UCF': avg_model_ucf_msd_cv,
    'RMSE Cosine Similarity - ICF': avg_model_icf_cosine_cv,
    'RMSE Pearson Similarity - ICF': avg_model_icf_pearson_cv,
    'RMSE MSD Similarity - ICF': avg_model_icf_msd_cv
}

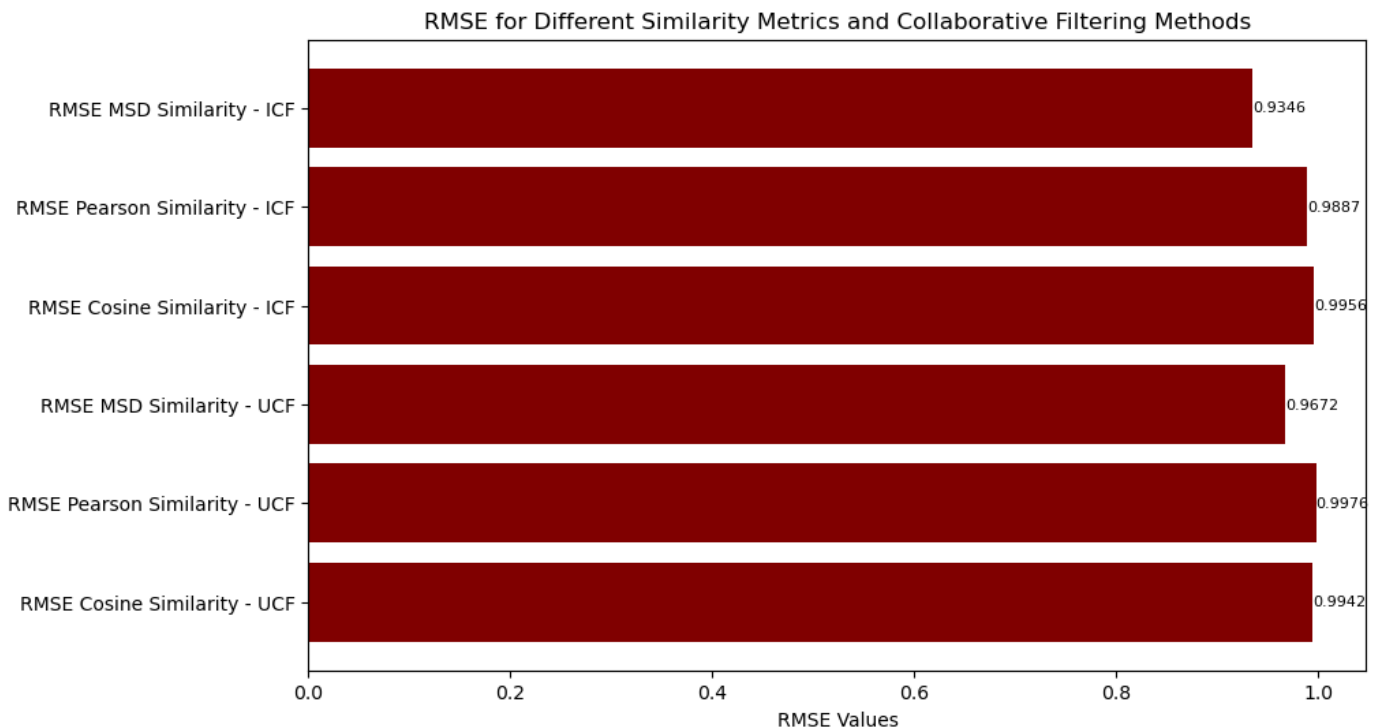
fig, ax = plt.subplots(figsize=(10, 6))
bars = ax.barh(
    list(final_res.keys()), list(final_res.values()),
```

```

color='maroon',
label=list(final_res.values()))
)
for bar in bars:
    plt.text(bar.get_width(), bar.get_y() + bar.get_height() / 2, f'{bar.get_width():.4f}'
            va='center', ha='left', fontsize=8, color='black')

ax.set_xlabel('RMSE Values')
ax.set_title('RMSE for Different Similarity Metrics and Collaborative Filtering Methods')
plt.show()

```



```

In [21]: text = ""
text += f"Looking at the graph above Item CF ({avg_model_icf_msd_cv}) achieves **lower R"
text += f" Accuracy than User CF ({avg_model_ucf_msd_cv}) for **MSD** similarity measure"
text += f" than User CF ({avg_model_ucf_pearson_cv}) for **Pearson similarity measure**"
text += f" than User CF ({avg_model_ucf_cosine_cv}) for **Cosine** similarity measure."

display(Markdown(text))

```

Looking at the graph above Item CF (0.9346341553269231) achieves **lower RMSE** Accuracy than User CF (0.9672341143744596) for **MSD** similarity measure, **lower RMSE** Accuracy(0.9886948589375084) than User CF (0.9975578954803224) for **Pearson similarity measure** and higher RMSE Accuracy(0.9956496541116596) than User CF(0.9942404700457219) for **Cosine** similarity measure.

While the difference is not significant for Pearson and Cosine, but looking at the pure numbers, the impact of the 3 metrics is not significantly consistent between User-based collaborative filtering and Item-based collaborative filtering.

3(f) Examine how the number of neighbors impacts the performances of User based Collaborative Filtering and Item based Collaborative Filtering? Plot your results.

```

In [22]: k_values = np.arange(5, 41, 5)

def evaluate_rmse_for_different_k(model: KNNBasic, data: any):
    cv_results = cross_validate(model, data, measures=['RMSE'], cv=5)

```

```

return np.average(cv_results['test_rmse'])

sim_options_icf = {
    "name": 'cosine',
    "user_based": False
}

sim_options_ucf = {
    "name": 'cosine',
    'user_based': True
}

results = {}
for k_val in k_values:
    model_icf = KNNBasic(k=k_val, sim_options=sim_options_icf)
    model_ucf = KNNBasic(k=k_val, sim_options=sim_options_ucf)

    acc_icf_rmse = evaluate_rmse_for_different_k(model=model_icf, data=ratings)
    acc_ucf_rmse = evaluate_rmse_for_different_k(model=model_ucf, data=ratings)
    results[k_val] = (acc_icf_rmse, acc_ucf_rmse)

```

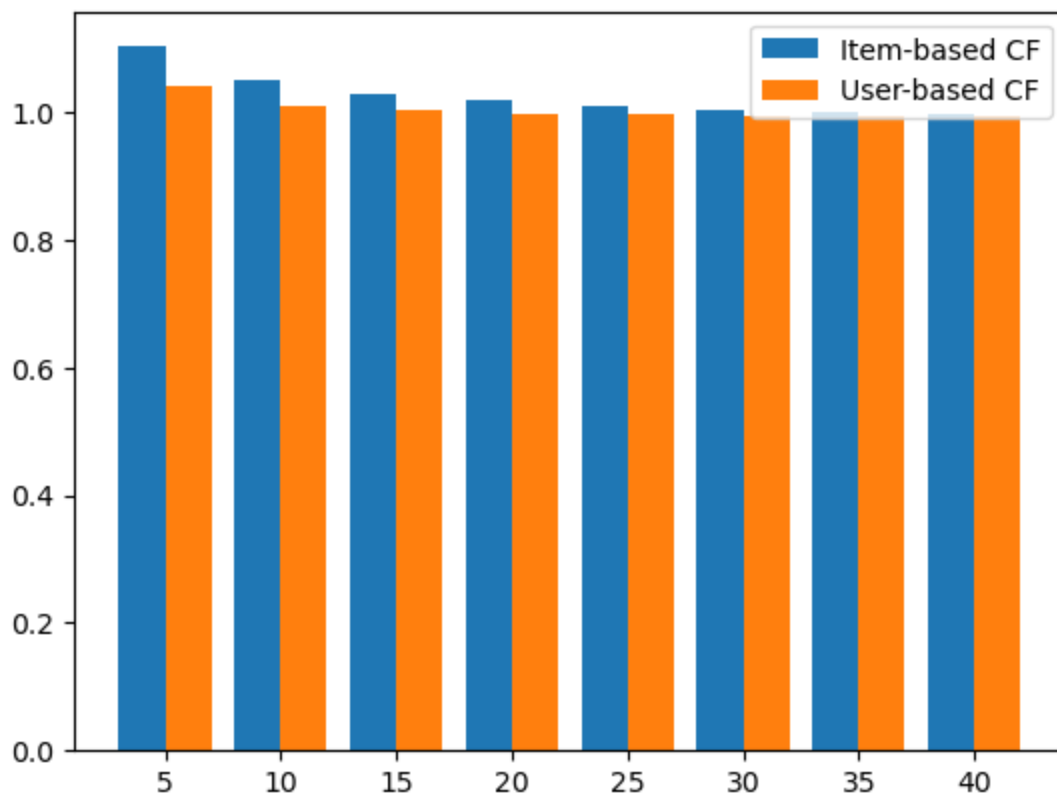
[illegible]

[illegible]

[illegible]

```
In [23]: fig, ax = plt.subplots()
X = list(results.keys())
X_axis = np.arange(len(X))
bar1 = ax.bar(X_axis - 0.2, [x[0] for x in results.values()], width = 0.4, label='Item-b
bar2 = ax.bar(X_axis + 0.2, [x[1] for x in results.values()], width = 0.4, label='User-b

ax.set_xticks(X_axis, X)
ax.legend()
plt.show()
print(results)
```



```
{5: (1.102664241201464, 1.041639865719923), 10: (1.0499401731271873, 1.010180773787399),
15: (1.0289903039569348, 1.0034415215135046), 20: (1.0185826859248295, 0.9978063708348737),
25: (1.0097188188543549, 0.9956883413588061), 30: (1.0043792775861513, 0.992922246293297),
35: (0.9989631624759054, 0.9941996662168922), 40: (0.9959949120649633, 0.9938738720676639)}
```

Examining the graph, we observe a consistent accuracy gap between Item-based Collaborative Filtering (ICF) and User-based Collaborative Filtering (UCF), with ICF consistently outperforming UCF. The number of neighbors does not appear to significantly impact the consistency of performance. However, it is notable that the RMSE accuracy keeps reducing after K value is computed.

3(g) Identify the best number of neighbor (denoted by K) for User/Item based collaborative filtering in terms of RMSE. Is the best K of User based collaborative filtering the same with the best K of Item based collaborative filtering? (10 points)

```
In [24]: k_for_min_ucf_rmse = min(results, key=lambda x : results[x][1])
k_for_min_icf_rmse = min(results, key=lambda x : results[x][0])
text = ""
text += f"K = {k_for_min_icf_rmse} has the lowest RMSE so is the best number of neighbor
text += f" K = {k_for_min_ucf_rmse} is the lowest RMSE and best number of neighbors for
display(Markdown(text))
```

K = 40 has the lowest RMSE so is the best number of neighbors for Item based collaborative filtering. K = 30 is the lowest RMSE and best number of neighbors for User-based collaborative filtering.

No, the best K of User based collaborative filtering is not the same with the best K of Item based collaborative filtering.