

Homework 2

```
In [1]: """
Project members:
Rezaei Yoo (by2350)
Kushal Wijesundera (kcu2144)
James Ding (jd3703)
Ryan McNally (rcm2109)
"""
```

```
In [2]: import os
import time
import sklearn
import numpy as np
import pandas as pd
import surprise
import seaborn as sns
import matplotlib.pyplot as plt

from utility import *
from surprise import SVD
from surprise import NMF
from surprise import Reader
from surprise import Dataset
from surprise import KNNWithMeans
from collections import defaultdict, Counter
from surprise.model_selection import GridSearchCV
from surprise.model_selection import cross_validate
from surprise.model_selection import train_test_split
matplotlib inline
```

Part 1: Business Objectives

- 1. Our objective**
- The objective of this recommendation system is to increase users' experience/stickiness and ultimately their time spent on our site by making their search for the movies easier as well as by making them constantly interested in our movie inventory. To achieve that, the recommendation system will provide a list of n (customizable, default 5 and we use 1 for this case study) movies to the users every time they log in, and the recommendation list will be a combination of movies that are most relevant to the users as well as the movies that are novel to the users.
- 2. Other metrics we wish to optimize, in addition to accuracy**
- To echo back to the objective of our recommendation system, we will use item-coverage and novelty as the two metrics we optimize in addition to accuracy
 - Therefore, we will simultaneously optimize for coverage and novelty, the hypothesis being that a sufficiently diverse and new set of recommendations with high accuracy will lead to deeper engagement, more time on the site, more clicks through the recommender, and higher high quality data for improving the model in later stages.
- 3. The intended user**
- The intended user for the business is any movie watcher & content consumer. Our objective is to maximize the engagement and personalization of this coming user. As a core piece of the user experience uplift strategy, this recommendation system intends to serve all users funneling to our site to search for movies to watch, which means this recommendation system should be able to make recommendations to both new users and existing users. New users with little to no engagement data will be recommended very diverse and popular movies while their behavior and tastes become evident to the model over time.
- 4. Business rules we think will be important**
- Never recommend items with less than 20% popularity. Our overall strategy will be to 'discard' the movies that are the least popular, 'maintain' the movies that are moderately popular and 'boost' the movies that are in the middle of the popularity spectrum.
 - Always include the most popular movies in the recommendation to users who are at the bottom 25% in terms of activeness(measured by the number of movies rated). This will serve as a baseline cold-start strategy.
 - Always recommend movies that the users have started but haven't finished if the relevant user data is accessible.
 - Always recommend movies that are trending now. For example, recommend romance movies during Valentine's Day.
 - Always include some level of novelty in the recommendations, meaning movies from different genres.
- 5. Performance requirements**
- The time complexity of rating making should be $O(N)$, where N denotes the number of movie items.
- 6. Notable successives**
- As it pertains to the final recommendation, there is room for future improvement in testing the recommendation of unpopular movies on subsets of users, and then boosting them pending positive feedback. As it stands, we require a critical mass of ratings for a film to be considered for recommendation by the model (20th percentile of ratings received). A new movie that is never recommended has no hope of getting to that critical mass.
 - As it pertains to **optimization and accuracy**: We are willing to sacrifice marginal 'accuracy' for the deeper and broader goal of attaining a holistic understanding of our users and their tastes. A maximally accurate but non-novel or high-coverage model might converge and recommend the user's favorite movie back to them n times repeatedly; the chances of a user remaining on the site or clicking through over and over in that situation are low.
 - As it pertains to **data sampling for model building**: Our sampling method biases for existing popularity, this will ignore currently unpopular users and movies (<50th percentile in the dataset). This is mainly for ensuring richness of signal in the final model which will be recommending to sparse datapoint users. This is a business decision, we would rather have a functional beta model to begin with than to erroneously recommend movies with a low network effect, at least for an initial v1. We cannot let the perfect be the enemy of the good. (We delve into tuning this parameter in Part 6 with impacts on training performance and model performance.)

Part 2: Exploratory Data Analysis and Two models building

Based on the instructions, we downloaded the entire movie data set at <https://grouplens.org/datasets/movielens/latest/> with 27M ratings for exploratory data analysis. Within the movie data, the ratings data set consists of 'usersid', 'movieid', 'ratings' and 'timestamp' along with the movies dataset, which consists of 'movieids', 'titles' and 'genres'.

```
In [3]: # load the data set for full ratings
ratings = pd.read_csv('./ml-latest/ratings.csv',
                    usecols=['userId', 'movieId', 'rating'],
                    dtype={'userId': 'int32', 'movieId': 'int32', 'rating': 'float128'},
                    engine='c')
movies = pd.read_csv('./ml-latest/movies.csv',
                    dtype={'movieId': 'int32'})
print("Number of distinct participated users: ", ratings['userId'].nunique())
print("Number of distinct rated movies: ", ratings['movieId'].nunique())
```

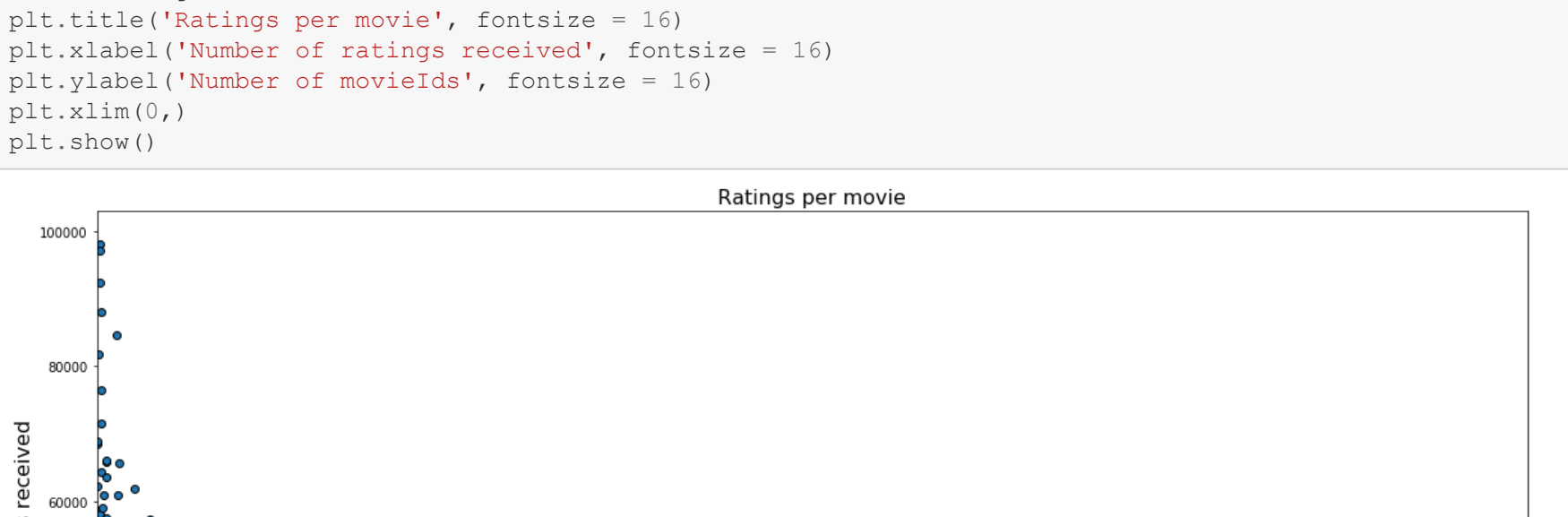
number of distinct participated users: 283228
number of distinct rated movies: 53889

```
In [4]: # merge data
movie_ratings = ratings.merge(movies, on = 'movieId', how = 'inner')
# numerical data statistics
movie_ratings.describe(include='number').transpose()
```

	count	mean	std	min	25%	50%	75%	max
userid	27753444	141942.015571	81707.400091	1.0	71176.0	142022.0	212459.0	283228.0
movieid	27753444	18487.998834	35102.623248	1.0	1097.0	2716.0	7150.0	193886.0
rating	27753444	3.530445	1.060353	0.5	3.0	3.5	4.0	5.0

```
In [5]: # categorical data statistics
movie_ratings.describe(include='object').transpose()
```

	count	unique	top	freq
title	27753444	53817	Shawshank Redemption, The (1994	97999
genres	27753444	1610	Drama	1959308



The distribution of the ratings for the entire data set is shown above. It presents a negatively skewed distribution where the mean rating and median is above 3. (note: Number of ratings in millions)

For experiment and development purposes, we will extract a small subset of users and movie items from the "full" dataset. We will first select ~1000 items and ~25000 users for model development and we will then select ~200 items and ~1500 users from the development set for testing purposes.

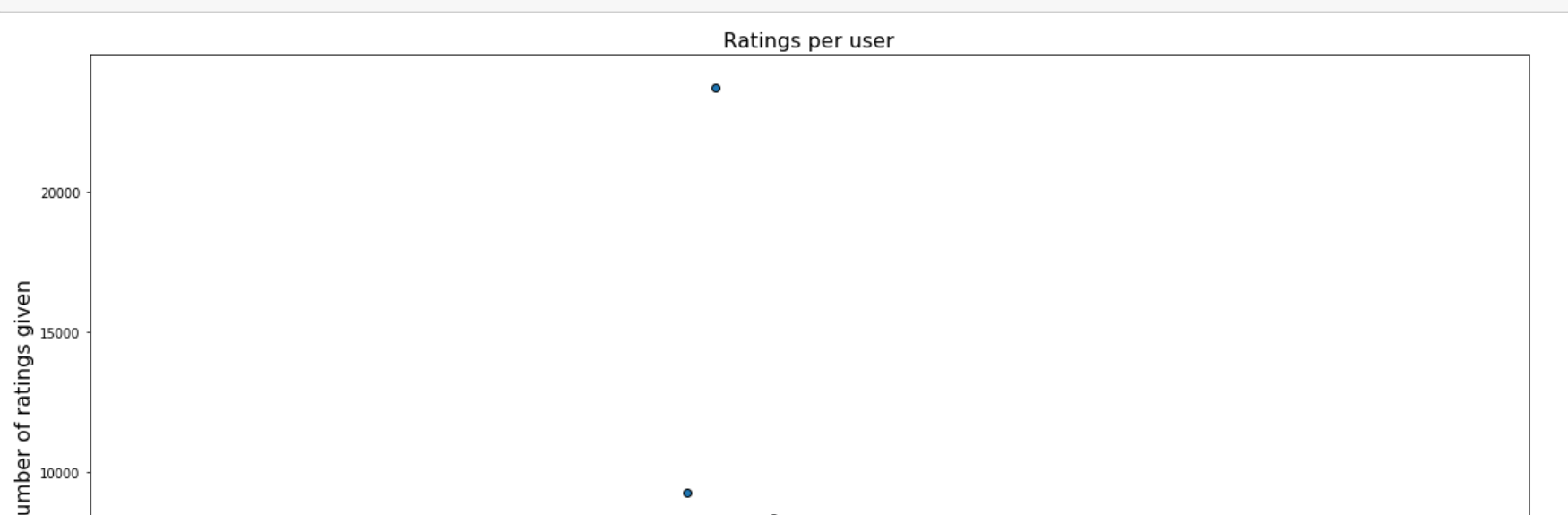
The methodology of selecting these users and items will be discussed in the data exploration section.

Data exploration and dataset sampling

Understanding how the data is distributed is crucial for designing the subsampling strategy and we believe that, in a real world setting, data collected from explicit feedbacks like ratings can be very sparse and data points are mostly collected from very popular movies and highly active users.

The following data exploration will justify our hypothesis.

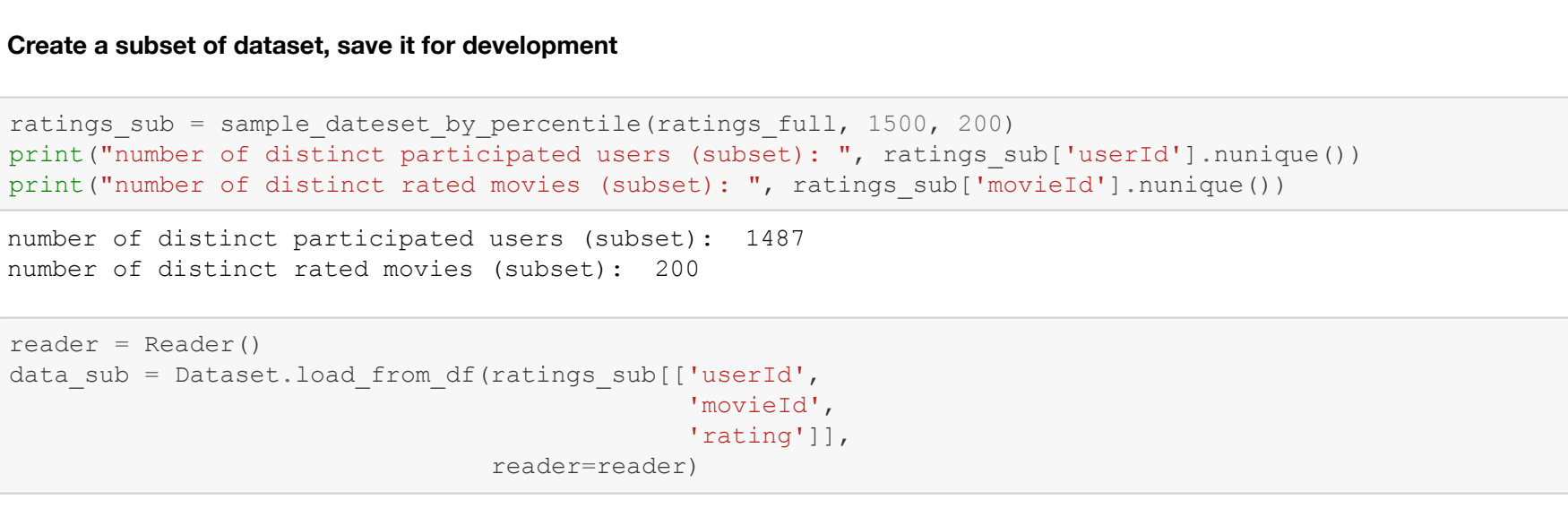
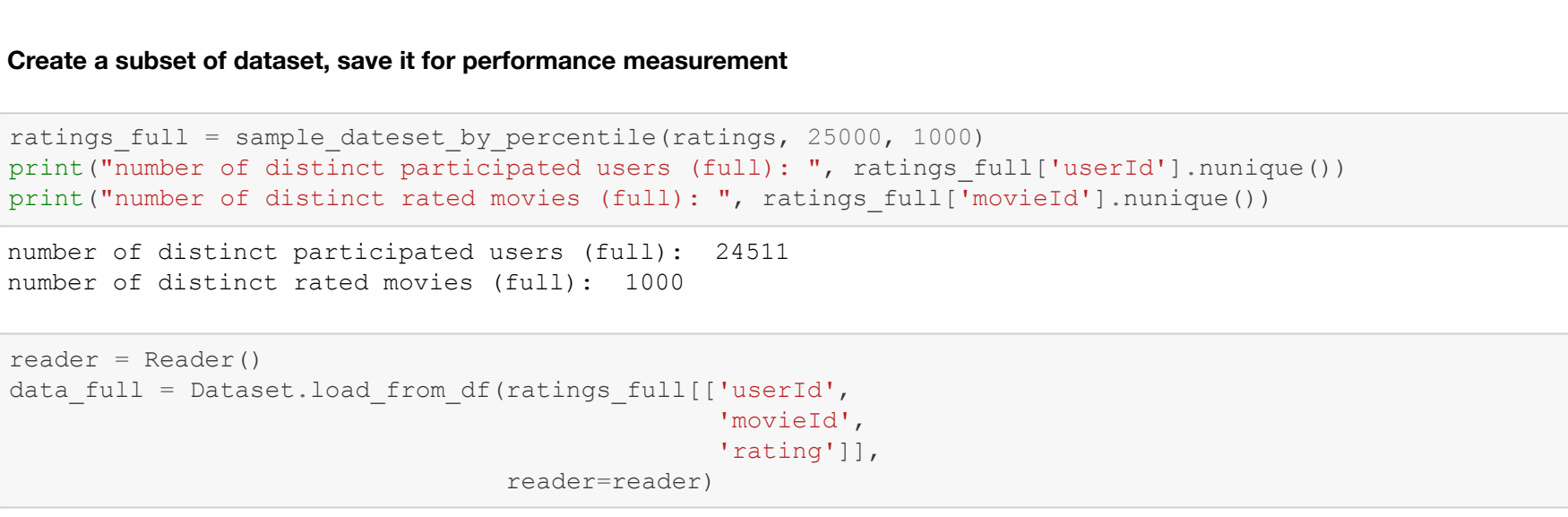
```
In [7]: # visualize the popularity by each movie
user_activeness = ratings[ratings['userId'].isin('movieId')].groupby('movieId').count()
movie_popularity.columns = ['num_ratings']
plt.figure(figsize=(20, 10))
plt.scatter(user_activeness.index,
            movie_popularity.num_ratings,
            edgecolor='black')
plt.xlim(0, movie_popularity.index.max())
plt.ylim(0, movie_popularity.index.max())
plt.title('Ratings per movie', fontsize = 16)
plt.xlabel('Number of ratings received', fontsize = 16)
plt.ylabel('Number of movies', fontsize = 16)
plt.show()
```



```
# visualize the distribution of the popularity across all movies
plt.figure(figsize=(20, 10))
plt.hist(movie_popularity.num_ratings,
        bins=200,
        edgecolor='black',
        log=True)
plt.title('Ratings per movie', fontsize = 16)
plt.xlabel('Number of ratings received', fontsize = 16)
plt.ylabel('Number of movies', fontsize = 16)
plt.xlim(0, movie_popularity.index.max())
plt.show()
```



```
In [8]: # visualize the popularity by each movie
user_activeness = ratings[ratings['userId'].isin('movieId')].groupby('userId').count()
movie_popularity.columns = ['num_ratings']
plt.figure(figsize=(20, 10))
plt.scatter(user_activeness.index,
            movie_popularity.num_ratings,
            edgecolor='black')
plt.xlim(0, movie_popularity.index.max())
plt.ylim(0, movie_popularity.index.max())
plt.title('Ratings per user', fontsize = 16)
plt.xlabel('Number of ratings given', fontsize = 16)
plt.ylabel('Number of users', fontsize = 16)
plt.show()
```



- The graphs above preliminarily confirm our assumptions regarding the distribution of the number of ratings received per movie and the distribution of the number of ratings given per user with a few particularly enthusiastic viewers, which are two extremely long-tailed distributions. Based on this observation, we ideally want to draw a subset of movie items and users that can represent the original dataset without giving unrealistically good or bad experiment results.
- We specifically want to get rid of the movies whose popularity is at the bottom 50 percentile, because we believe that there is a high possibility that those movies are only favored by an extremely small group of users, and more importantly, that the average ratings received by those movies are more likely to be biased. The same strategy also applies to users.
- For example, if certain movies are rated only once and received an average rating of 5, in an item-based collaborative filtering model, those movies will accidentally become part of the recommendations for many if not all users, which is not ideal.

Create a subset of dataset, save it for performance measurement

```
In [9]: ratings_full = sample_dataset_by_percentile(ratings, 25000, 1000)
print("Number of distinct participated users (full): ", ratings_full['userId'].nunique())
print("Number of distinct rated movies (full): ", ratings_full['movieId'].nunique())
number of distinct participated users (full): 24511
number of distinct rated movies (full): 1000
```

```
In [10]: reader = Reader()
data_full = Dataset.load_from_df(ratings_full[['userId',
                                             'movieId',
                                             'rating']],
                                reader=reader)
```

Create a subset of dataset, save it for development

```
In [11]: ratings_sub = sample_dataset_by_percentile(ratings_full, 1500, 200)
print("Number of distinct participated users (subset): ", ratings_sub['userId'].nunique())
print("Number of distinct rated movies (subset): ", ratings_sub['movieId'].nunique())
number of distinct participated users (subset): 1487
number of distinct rated movies (subset): 200
```

```
In [12]: reader = Reader()
data_sub = Dataset.load_from_df(ratings_sub[['userId',
                                             'movieId',
                                             'rating']],
                                reader=reader)
```

After careful manipulation, we have successfully keep 1487 users from top active levels and 200 movie items from top popularity levels, for development purposes.

Part 2.1: Item-based neighborhood method

```
In [13]: trainvalset, testset = train_test_split(data_sub, test_size=0.2)
```

```
In [14]: # baseline model constructing
item_model = KNNWithMeans(k=4, sim_options={'name': 'pearson',
                                           'user_based': False,
                                           'verbose': False})
```

```
In [15]: item_model.fit(trainvalset);
Computing the pearson similarity matrix...
Done computing similarity matrix.
```

```
In [16]: uid = str(699)
lid = str(296)
item_model_pred = item_model.predict(uid,
                                     lid,
                                     verbose=False)
print("Prediction for rating: ", item_model_pred)
```

Prediction for rating: user: 699 item: 296 r_ui = None est = 3.49 ['was_impossible': True, 'reason': 'User and/or item is unknown.']

- In this project, we'll use the surprise package, a popular package for building recommendation systems in Python.
- At this part, we have successfully built and tested a baseline item-based neighborhood model by making a demo rating prediction for user 699 and item 296.
- The item-based neighborhood model built above uses pearson as the similarity metric and considers the ratings of four closest neighbors while making predictions.

Part 2.2: Matrix Factorization

```
In [17]: # baseline model constructing
NMF_model = NMF(n_factors=20,
               n_epochs=10,
               biased=True)
```

```
In [18]: NMF_model.fit(trainvalset);
```

```
In [19]: uid = str(699)
lid = str(296)
NMF_model_pred = NMF_model.predict(uid,
                                    lid,
                                    verbose=False)
print("Prediction for rating: ", NMF_model_pred)
```

Prediction for rating: user: 699 item: 296 r_ui = None est = 3.49 ['was_impossible': False]

- At this part, we have successfully built and tested a baseline matrix factorization model(NMF) by making a demo rating prediction for user 699 and item 296.
- The NMF model built above decomposes the characteristic matrix in latent space size of 20.

Part 3: Model evaluation

Part 3.1.1: CV setup for item-based neighborhood method

For this task, we will use RMSE(Root Mean Squared Error) as the primary accuracy metric and MAE(Mean Absolute Error) as the secondary accuracy metric.

```
In [20]: # Run 5-fold cross-validation and print results
result = cross_validate(item_model,
                        data_full,
                        measures=['RMSE', 'MAE'],
                        cv=5,
                        return_train_measures=True,
                        verbose=True)
```

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, MAE of algorithm KNNWithMeans on 5 split(s).

	Fold 1	Fold 2	Fold 3	Fold 4	Fold 5	Mean	Std
RMSE (testset)	0.9448	0.9478	0.9566	0.9518	0.9558	0.9514	0.0045
MAE (testset)	0.7128	0.7154	0.7207	0.7203	0.7225	0.7184	0.0037
RMSE (trainset)	0.4180	0.4189	0.4174	0.4193	0.4185	0.4184	0.0007
MAE (trainset)	0.3128	0.3130	0.3121	0.3135	0.3127	0.3128	0.0004
F1 time	0.31	0.38	0.99	1.00	1.01	0.98	0.03
Test time	0.13	0.07	2.18	2.13	2.11	2.12	0.04

The results above give an preliminary review of how well the matrix fits the training set as well as how well the model generalizes on the test set. According to the results above, the average testset RMSE is 0.95 and the average testset MAE is 0.72.

Part 3.1.2: Coverage of item-based neighborhood method

```
In [21]: trainvalset, testset = train_test_split(data_sub, test_size=0.2)
trainvalset_testify = trainvalset.build_anti_testset()
```

```
In [22]: user_list_train = set(trainvalset.all_users())
user_list_test = set((item['movie_id'] for item in trainvalset_testify))
```

```
In [23]: item_model.fit(trainvalset);
Computing the pearson similarity matrix...
Done computing similarity matrix.
```

```
In [24]: predictions_train = item_model.test(trainvalset_testify)
predictions_test = item_model.test(testset)
top_n_train = get_top_n(predictions_train, n=1)
top_n_test = get_top_n(predictions_test, n=1)
```

```
In [25]: recommendation_train = []
# Print the recommended items for each user
for uid, user_ratings in top_n_train.items():
    recommendation_train.append((uid for (lid, _) in user_ratings))
recommendation_train_list = []
for item in recommendation_train:
    for rec_list in recommendation_train:
        recommendation_train_list.append(item)
recommendation_train_list = set(recommendation_train_list)
```

```
In [26]: recommendation_test = []
# Print the recommended items for each user
for uid, user_ratings in top_n_test.items():
    recommendation_test.append((lid for (lid, _) in user_ratings))
recommendation_test_list = []
for item in recommendation_test:
    for rec_list in recommendation_test:
        recommendation_test_list.append(item)
recommendation_test_list = set(recommendation_test_list)
```

```
In [27]: print('Coverage on train set', '{:.2f}'.format(len(recommendation_train_list)/len(item_list_train)))
print('Coverage on test set', '{:.2f}'.format(len(recommendation_test_list)/len(item_list_test)))
coverage on train set 0.31
coverage on test set 0.66
```

- As we can see with the above discrepancy in coverage between the train and test sets, our recommender is converging on a few top performers in the train set but is more dispersed in the test set. We will test this hypothesis below by checking how often, on a percentage basis, the top movies were recommended in each case.
- With the lower coverage, we would expect that the top recommendations are recommended more often in the train set.

```
In [28]: train_counter = pd.DataFrame(Counter([item for sublist in recommendation_train for item in sublist]).most_common(n), columns=['movie_id', 'count'])
train_counter['pct_recommended'] = train_counter['count']/len(recommendation_train)
test_counter = pd.DataFrame(Counter([item for sublist in recommendation_test for item in sublist]).most_common(n), columns=['movie_id', 'count'])
test_counter['pct_recommended'] = test_counter['count']/len(recommendation_test)
```

```
In [29]: train_counter.join(movies).head()
Out[29]:
```

	title	count	pct recommended
0	Toy Story (1995)	702	0.482474
1	Jurassic Park (1993)	184	0.126460
2	Grumpier Old Men (1995)	80	0.054983
3	Waiting to Exhale (1995)	55	0.037001
4	Father of the Bride Part II (1995)	45	0.030928

```
In [30]: test_counter.join(movies).head()
Out[30]:
```

	title	count	pct recommended
0	Toy Story (1995)	205	0.215789
1	Jurassic Park (1993)	50	0.026332
2	Grumpier Old Men (1995)	42	0.044211
3	Waiting to Exhale (1995)	35	0.036842
4	Father of the Bride Part II (1995)	34	0.035789

- Our hypothesis is confirmed, we see that the test case recommended the most popular title (Toy Story) 19.89% of the time, while the train set was recommended it over 28.6% of the time.
- The discrepancy is driven by the 'richer rich' effect of the top performers nabbing the n=1 top spot.
- The next hypothesis will be to increase n to 5 and we will ideally see higher coverage in both sets, and relative closer parity between coverages.

```
In [31]: predictions_train = item_model.test(trainvalset_testify)
predictions_test = item_model.test(testset)
#new hypothesis n=5
top_n_train = get_top_n(predictions_train, n=5)
top_n_test = get_top_n(predictions_test, n=5)
```

```
In [32]: recommendation_train = []
# Print the recommended items for each user
for uid, user_ratings in top_n_train.items():
    recommendation_train.append((lid for (lid, _) in user_ratings))
recommendation_train_list = []
for item in recommendation_train:
    for rec_list in recommendation_train:
        recommendation_train_list.append(item)
recommendation_train_list = set(recommendation_train_list)
```

```
In [33]: recommendation_test = []
# Print the recommended items for each user
for uid, user_ratings in top_n_test.items():
    recommendation_test.append((lid for (lid, _) in user_ratings))
recommendation_test_list = []
for item in recommendation_test:
    for rec_list in recommendation_test:
        recommendation_test_list.append(item)
recommendation_test_list = set(recommendation_test_list)
```

```
In [34]: print('Coverage on train set', '{:.2f}'.format(len(recommendation_train_list)/len(item_list_train)))
print('Coverage on test set', '{:.2f}'.format(len(recommendation_test_list)/len(item_list_test)))
coverage on train set 0.12
coverage on test set 0.73
```

```
In [41]: train_counter = pd.DataFrame(Counter([item for sublist in recommendation_train for item in sublist]).most_common(n), columns=['movie_id', 'count'])
train_counter['pct_recommended'] = train_counter['count']/len(recommendation_train)
test_counter = pd.DataFrame(Counter([item for sublist in recommendation_test for item in sublist]).most_common(n), columns=['movie_id', 'count'])
test_counter['pct_recommended'] = test_counter['count']/len(recommendation_test)
```

```
In [42]: train_counter.join(movies).head()
Out[42]:
```

	title	count	pct recommended
0	Toy Story (1995)	1121	0.770447
1	Jurassic Park (1993)	683	0.469416
2	Grumpier Old Men (1995)	641	0.440550
3	Waiting to Exhale (1995)	504	0.343092
4	Father of the Bride Part II (1995)	483	0.331659

```
In [43]: test_counter.join(movies).head()
Out[43]:
```

	title	count	pct recommended
0	Toy Story (1995)	629	0.423202
1	Jurassic Park (1993)	47	0.049474
2	Grumpier Old Men (1995)	44	0.046316
3	Waiting to Exhale (1995)	35	0.036842
4	Father of the Bride Part II (1995)	31	0.032332

- Here we see we have much the same effect as with matrix factorization as we do with KNN. Higher coverage on test than train, and concentration among top performers in train.
- Below we run the same experiment, increasing n to 5:

```
In [44]: top_n_train = get_top_n(predictions_train, n=5)
top_n_test = get_top_n(predictions_test, n=5)
```

```
In [45]: recommendation_train = []
# Print the recommended items for each user
for uid, user_ratings in top_n_train.items():
    recommendation_train.append((lid for (lid, _) in user_ratings))
recommendation_train_list = []
for item in recommendation_train:
    for rec_list in recommendation_train:
        recommendation_train_list.append(item)
recommendation_train_list = set(recommendation_train_list)
```

```
In [46]: recommendation_test = []
# Print the recommended items for each user
for uid, user_ratings in top_n_test.items():
    recommendation_test.append((lid for (lid, _) in user_ratings))
recommendation_test_list = []
for item in recommendation_test:
    for rec_list in recommendation_test:
        recommendation_test_list.append(item)
recommendation_test_list = set(recommendation_test_list)
```

```
In [47]: print('Coverage on train set', '{:.2f}'.format(len(recommendation_train_list)/len(item_list_train)))
print('Coverage on test set', '{:.2f}'.format(len(recommendation_test_list)/len(item_list_test)))
coverage on train set 0.24
coverage on test set 0.97
```


[illegible]