

Movie Recommendation System



You pick the movie, I'll choose the restaurant...

1. Project Overview

Choosing movies can be stressful, high-stakes endeavor for anyone. Fortunately, this program is here to help. It provides movie recommendations for any new user, based on their movie preferences. The program asks the user to rate five random movies from a database of 9,742 movies. Based on how the new user rates the five random movies, the program utilizes a user-based collaborative filtering model to provide five recommendations.

2. Business Case

With the vast entertainment options available, low engagement and user churn in any social media or streaming service can hurt profit. Considering that 20% of adults are indecisive and 67% of relationship agreements never get resolved, picking a movie can be a daunting task, fueling decision paralysis and disengagement. Luckily, machine learning can relieve indecision by providing recommendations for any user, based on their preferences.

3. Data Understanding

Dataset Background

So, we have our data spanning over 4 separate csv files. We also have a README file which may tell us how this data interacts. Let's open that file to gain some insight.

In [26]:

```
1 file_path = 'data/README.txt'  
2  
3 with open(file_path) as file:  
4     print(file.read())  
service. It contains 100836 ratings and 3683 tag applications across 9742 movies. These data were created by 610 users between March 29, 1996 and September 24, 2018. This dataset was generated on September 26, 2018.
```

Users were selected at random for inclusion. All selected users had rated at least 20 movies. No demographic information is included. Each user is represented by an id, and no other information is provided.

The data are contained in the files `links.csv` , `movies.csv` , `ratings.csv` and `tags.csv` . More details about the contents and use of all these files follows.

This is a *development* dataset. As such, it may change over time and is not an appropriate dataset for shared research results. See available *benchmark* datasets if that is your intent.

This and other GroupLens data sets are publicly available for download at <<http://grouplens.org/datasets/>>.

Summary

This dataset (ml-latest-small) describes 5-star rating and free-text tagging activity from [MovieLens \(http://movielens.org\)](http://movielens.org), a movie recommendation service. It contains 100836 ratings and 3683 tag applications across 9742 movies. These data were created by 610 users between March 29, 1996 and September 24, 2018. This dataset was generated on September 26, 2018.

Users were selected at random for inclusion. All selected users had rated at least 20 movies. No demographic information is included. Each user is represented by an id, and no other information is provided.

The data are contained in the files `links.csv` , `movies.csv` , `ratings.csv` and `tags.csv` . More details about the contents and use of all these files follows.

Data Files

Ratings Data File Structure (ratings.csv)

All ratings are contained in the file `ratings.csv`. Each line of this file after the header row represents one rating of one movie by one user, and has the following format:

```
userId,movieId,rating,timestamp
```

Tags Data File Structure (`tags.csv`)

All tags are contained in the file `tags.csv`. Each line of this file after the header row represents one tag applied to one movie by one user, and has the following format:

```
userId,movieId,tag,timestamp
```

Movies Data File Structure (`movies.csv`)

Movie information is contained in the file `movies.csv`. Each line of this file after the header row represents one movie, and has the following format:

```
movieId,title,genres
```

Links Data File Structure (`links.csv`)

Identifiers that can be used to link to other sources of movie data are contained in the file `links.csv`. Each line of this file after the header row represents one movie, and has the following format:

```
movieId,imdbId,tmdbId
```

Data Inspection

Let's go ahead and see if we can verify some of this data. I'm going to go ahead and import

```
In [27]: 1 import pandas as pd  
          2 import numpy as np  
          3 import random
```

Ratings File Summary

```
In [28]: 1 ratings_df = pd.read_csv('data/ratings.csv')  
          2 ratings_df.head(5)
```

Out[28]:

	userId	movieId	rating	timestamp
0	1	1	4.0	964982703
1	1	3	4.0	964981247
2	1	6	4.0	964982224
3	1	47	5.0	964983815
4	1	50	5.0	964982931

```
In [29]: 1 ratings_df.describe()
```

Out[29]:

	userId	movieId	rating	timestamp
count	100836.000000	100836.000000	100836.000000	1.008360e+05
mean	326.127564	19435.295718	3.501557	1.205946e+09
std	182.618491	35530.987199	1.042529	2.162610e+08
min	1.000000	1.000000	0.500000	8.281246e+08
25%	177.000000	1199.000000	3.000000	1.019124e+09
50%	325.000000	2991.000000	3.500000	1.186087e+09
75%	477.000000	8122.000000	4.000000	1.435994e+09
max	610.000000	193609.000000	5.000000	1.537799e+09

It's interesting that the rating in 25-75 percentile range are from 3.0-4.0, meaning user generally rate movies favorably.

```
In [30]: 1 ratings_df.isna().sum()
```

Out[30]:

```
userId      0
movieId     0
rating      0
timestamp   0
dtype: int64
```

```
In [31]: 1 unique_movies = list(ratings_df['movieId'].unique())
2 print('Number of movies: ', len(unique_movies), '\n')
3
4 unique_users = list(ratings_df['userId'].unique())
5 print('Number of ratings: ', len(unique_users))
6
```

Number of movies: 9724

Number of ratings: 610

So, we have confirmed no null values, as well as 10,0836 movie ratings and a maximum userID of 610. All of our ratings our .5 - 5.0 and... we have 9,724 movies. This looks promising so far and matches our README.

Tags File Summary

```
In [32]: 1 tags_df = pd.read_csv('data/tags.csv')
2 tags_df.head(5)
```

Out[32]:

	userId	movieId	tag	timestamp
0	2	60756	funny	1445714994
1	2	60756	Highly quotable	1445714996
2	2	60756	will ferrell	1445714992
3	2	89774	Boxing story	1445715207
4	2	89774	MMA	1445715200

```
In [33]: 1 tags_df.describe()
```

Out[33]:

	userId	movieId	timestamp
count	3683.000000	3683.000000	3.683000e+03
mean	431.149335	27252.013576	1.320032e+09
std	158.472553	43490.558803	1.721025e+08
min	2.000000	1.000000	1.137179e+09
25%	424.000000	1262.500000	1.137521e+09
50%	474.000000	4454.000000	1.269833e+09
75%	477.000000	39263.000000	1.498457e+09
max	610.000000	193565.000000	1.537099e+09

```
In [34]: 1 tags_df.isna().sum()
```

Out[34]:

```
userId      0
movieId     0
tag         0
timestamp   0
dtype: int64
```

```
In [35]: 1 tags_df['tag'].value_counts()
```

Out[35]:

```
In Netflix queue      131
atmospheric            36
thought-provoking     24
superhero              24
funny                  23
...
action choreography    1
oil                     1
bizzare                1
bears                  1
high fantasy            1
Name: tag, Length: 1589, dtype: int64
```

okay, so... this looks good. Our tags folder contains info for up to 610 user ids, a max movie id of 19365, and no null values. We can already see a few trends with the tags - namely in Netflix queue, atmospheric, and super-hero as the most popular trend. Let's go to the Movie Ids DataFrame.

```
In [36]: 1 unique_items = list(tags_df['tag'].unique())
2 len(unique_items)
3 unique_items[0:20]
```

```
Out[36]: ['funny',
'Highly quotable',
'will ferrell',
'Boxing story',
'MMA',
'Tom Hardy',
'drugs',
'Leonardo DiCaprio',
'Martin Scorsese',
'way too long',
'Al Pacino',
'gangster',
'mafia',
'Mafia',
'holocaust',
'true story',
'twist ending',
'Anthony Hopkins',
'courtroom drama',
'britpop']
```

So, among the first 20 unique tags, we can see a discrepancy between 'mafia' and 'Mafia', so we know that we might need to include the lowercase. We can see, especially with actors, that some appear in lower case "will ferrell" while others appear capitalized, like "Tom Hardy."

Movies File Summary

```
In [37]: 1 movies_df = pd.read_csv('data/movies.csv')
2 movies_df.tail(5)
```

```
Out[37]:
```

	movieId	title	genres
9737	193581	Black Butler: Book of the Atlantic (2017)	Action Animation Comedy Fantasy
9738	193583	No Game No Life: Zero (2017)	Animation Comedy Fantasy
9739	193585	Flint (2017)	Drama
9740	193587	Bungo Stray Dogs: Dead Apple (2018)	Action Animation
9741	193609	Andrew Dice Clay: Dice Rules (1991)	Comedy

```
In [38]: 1 movies_df.isna().sum()
```

```
Out[38]: movieId      0  
          title       0  
          genres      0  
          dtype: int64
```

Again, this look good. It appears that there are about 9742 movies with no null movies, which is exactly what the README said. So we're good. We're not going to concern ourselves with the links ID currently. For now, we'll leave it be.

4. Data Preparation

Data Approach

So we have a lot of data present in these files. For this project, we will only need to use the ratings file, and apply a matrix factorization to it.

For the ratings file, we will drop the timestamp, as we're not as interested in the time series data.

Additional Prep (not needed for this exercise)

Additionally, we will one-hot encode the genre information, as well as the tags.

The genre info is limited to the 19 categories, the tags on the other hand, have over 1589 unique values. We might be able to clean some of these up, but that is still quite a lot.

Dropping timestamp.

I will drop the timestamp from each of the `ratings_df` and `tags_df`.

```
In [39]: 1 ratings = ratings_df.drop('timestamp', axis = 1)  
2 tags = tags_df.drop('timestamp', axis = 1)
```

Lowercase

As we indicated above, we need to convert the tags to lower case.

```
In [40]: 1 #tags['tag']= tags['tag'].str.lower()  
2 #tags['tag']= tags['tag'].apply(str.maketrans('', '', '!@#$'))  
3 tags['tag']= tags['tag'].map(lambda x: x.lower().rstrip('!"!@#$').rstrip('`'))
```

One-hot encoding

Before we merge the files let's go ahead and one-hot encode both the 'genre' category in `movies DataFrame` and the 'tags' in the `tags DataFrame`

```
In [41]: 1 tags_ohe = pd.get_dummies(tags, columns = ['tag'], prefix='', prefix_sep='')
2 tags_ohe
```

Out[41]:

					06 oscar nominated best movie - animation								
					1900s	1920s	1950s	1960s	1970s	1980s	...	world war i	
0	2	60756	0	0	0	0	0	0	0	0	0	...	0
1	2	60756	0	0	0	0	0	0	0	0	0	...	0
2	2	60756	0	0	0	0	0	0	0	0	0	...	0
3	2	89774	0	0	0	0	0	0	0	0	0	...	0
4	2	89774	0	0	0	0	0	0	0	0	0	...	0
...
3678	606	7382	0	0	0	0	0	0	0	0	0	...	0
3679	606	7936	0	0	0	0	0	0	0	0	0	...	0
3680	610	3265	0	0	0	0	0	0	0	0	0	...	0
3681	610	3265	0	0	0	0	0	0	0	0	0	...	0
3682	610	168248	0	0	0	0	0	0	0	0	0	...	0

3683 rows × 1477 columns

5. Modeling & Evaluation

As was noted in the previous section, we can create our model simply from the `ratings.csv` file. To create our baseline model, we're going to use the `surprise` module. We will compare SVD and a variety of KNN based methods within the `surprise` module to determine which is the most accurate for our dataset. For consistency sake, will use RSME (Root Square Mean Error).

We will also establish a baseline of Random prediction ratings to see how the RSME compares.

Random

To see how well our model does, we will compare over a "random" model that predict a user's ratings. THis model will just pick ratings at random between (.5 and 5.0).

```
In [44]: 1 #Let's create a column to test our first random column, which is a random
          2 rand_rate = ratings
          3 rand_rate[ 'predicted' ] = np.random.randint(1,10, rand_rate.shape[0]) / 2
          4
          5 rand_rate
```

Out[44]:

	userId	movieId	rating	predicted
0	1	1	4.0	3.5
1	1	3	4.0	0.5
2	1	6	4.0	4.5
3	1	47	5.0	4.0
4	1	50	5.0	2.5
...
100831	610	166534	4.0	4.0
100832	610	168248	5.0	0.5
100833	610	168250	5.0	3.5
100834	610	168252	5.0	1.0
100835	610	170875	3.0	3.0

100836 rows × 4 columns

We successfully created a new column called 'predicted' for all of the movies. Now, let's see if we can

```
In [43]: 1 rmse = np.sqrt(((rand_rate[ 'predicted' ] - rand_rate[ 'rating' ]) ** 2).mean()
          2 rmse
```

Out[43]: 1.937287850567029

Our RMSE for this random baseline is 1.94. I hope we can beat that in our surprise modules

surprise module models

Now that we have established RSME from random models, let's go ahead and try some of the beefier models in surprise. First, we're going to read in our dataset and establish test and trainsets. Then we will

Reading our Dataset

To begin, we will go through the process of reading in our dataset into the surprise dataset format. This will make the subsequent modeling a little more fluid.

```
In [17]: 1 #import the relevant item from surprise  
2 from surprise.model_selection import cross_validate  
3 from surprise.prediction_algorithms import SVD  
4 from surprise.prediction_algorithms import KNNWithMeans, KNNBasic, KNNBaseline  
5 from surprise.model_selection import GridSearchCV  
6 from surprise.model_selection import train_test_split  
7 from surprise import accuracy
```

As a way to validate the data, we're going to create a test and train set of data.

```
In [18]: 1 #read in dataset to surprise format  
2 from surprise import Reader, Dataset  
3 reader = Reader()  
4 data = Dataset.load_from_df(ratings,reader)  
5  
6 # we will create a test set for validation, this will be used later when we want to evaluate our model.  
7 trainset, testset = train_test_split(data, test_size=0.2)
```

```
In [19]: 1 #check to make sure item's loaded properly and create a new trainset.  
2 dataset = data.build_full_trainset()  
3 print('Number of users: ', dataset.n_users, '\n')  
4 print('Number of ratings: ', dataset.n_items)
```

Number of users: 610

Number of ratings: 9724

This matches our original check so... we've appeared to load the data successfully.

Model-Based Methods (Matrix Factorization) - SVD with surprise module

Below we will use the surprise method to create a SVD model, with tuned hyperparameters. We will utilize GridSearchCV for this.

```
In [20]: 1 ## we will set up a SVD model with appropriate hyperparameters.
2
3 #established some initial hyperparameters
4 params = {'n_factors': [20, 50, 100],
5            'reg_all': [0.02, 0.05, 0.1],
6            'n_epochs': [5, 10, 15],
7            'lr_all': [.002, .005, .010]}
8
9 #instantiate GridSearchCV model
10 g_s_svd = GridSearchCV(SVD,param_grid=params,n_jobs = -1,joblib_verbose=5)
11
12 #fit our ratings dataset "data" onto the model
13 g_s_svd.fit(data)
```

[Parallel(n_jobs=-1)]: Using backend LokyBackend with 4 concurrent workers.
[Parallel(n_jobs=-1)]: Done 10 tasks | elapsed: 14.5s
[Parallel(n_jobs=-1)]: Done 64 tasks | elapsed: 58.3s
[Parallel(n_jobs=-1)]: Done 154 tasks | elapsed: 2.4min
[Parallel(n_jobs=-1)]: Done 280 tasks | elapsed: 5.3min
[Parallel(n_jobs=-1)]: Done 405 out of 405 | elapsed: 9.5min finished

```
In [21]: 1 print(g_s_svd.best_score)
2 print(g_s_svd.best_params)
```

{'rmse': 0.8629065192573494, 'mae': 0.6627385253795459}
{'rmse': {'n_factors': 100, 'reg_all': 0.05, 'n_epochs': 15, 'lr_all': 0.01},
'mae': {'n_factors': 100, 'reg_all': 0.05, 'n_epochs': 15, 'lr_all': 0.01}}

Now let's run the model we have and print the results of our testset.

```
In [22]: 1 svd = SVD(n_factors=100, n_epochs=15, lr_all=0.010, reg_all=0.05)
2 svd.fit(trainset)
3 predictions = svd.test(testset)
4 print(accuracy.rmse(predictions))
```

RMSE: 0.8651
0.8650715807530015

Okay, we see a RMSE of .87. This... isn't bad on a scale of 0.5-5.0. Essentially it's under 1, which feels good, but not under 0.5, which would feel better. Our model had an rmse of 0.87, let's establish that as OUR BASELINE MODEL.

Our optimal parameters are n_factors = 100 and reg_all = .05. This is convenient that these are in the middle of our range. We'll do a few quick spot checks to see if we can improve this.

```
In [23]: 1 svd = SVD(n_factors=100, n_epochs=20, lr_all=0.050, reg_all=0.05)
2 svd.fit(trainset)
3 predictions = svd.test(testset)
4 print(accuracy.rmse(predictions))
```

RMSE: 0.8640
0.8639629266342134

```
In [24]: 1 svd = SVD(n_factors=150, n_epochs=25, lr_all=0.010, reg_all=0.05)
2 svd.fit(trainset)
3 predictions = svd.test(testset)
4 print(accuracy.rmse(predictions))
```

RMSE: 0.8585
0.8585001642972838

```
In [25]: 1 svd = SVD(n_factors=200, n_epochs=30, lr_all=0.010, reg_all=0.05)
2 svd.fit(trainset)
3 predictions = svd.test(testset)
4 print(accuracy.rmse(predictions))
```

RMSE: 0.8597
0.8596641766038512

So... it did go down, but it barely moved. Suffice to say that perhaps we've created a largely optimized model. We can return to this later. Our hyperparameters are {n_factors=150, n_epochs=25, lr_all=0.010, reg_all=0.05}

Memory-Based Methods (Neighborhood-Based) KNN with surprise

To begin with, we can calculate the more simple neighborhood-based approaches. We can start with KNNBasic. With KNNBasic, we'll need a trainset and a testset in order to cross-validate results. We also run a few examples to determine the best hyperparameters

We'll import the relevant first.

```
In [26]: 1 #initiating KNN Basic with pearson similarity matrix and user_based similarity
2 knn_basic = KNNBasic(sim_options={'name':'pearson', 'user_based':True})
3 knn_basic.fit(trainset)
4 predictions = knn_basic.test(testset)
5 print(accuracy.rmse(predictions))
```

Computing the pearson similarity matrix...
Done computing similarity matrix.
RMSE: 0.9753
0.9753330994599322

With the KNN Basic, we have to set some of our hyper parameters. We'll try both "cosine" and "pearson". We'll also establish user based similarity, as there are fewer users than movies so this will save us considerable time. If we had thousands of users and only a handful of movies, we would consider an item based similarity.

Let's try cosine below.

```
In [27]: 1 #initiating KNN Basic with pearson correlation
2 knn_basic = KNNBasic(sim_options={'name':'cosine', 'user_based':True})
3 knn_basic.fit(trainset)
4 predictions = knn_basic.test(testset)
5 print(accuracy.rmse(predictions))
```

Computing the cosine similarity matrix...
Done computing similarity matrix.
RMSE: 0.9752
0.9751602971078402

```
In [28]: 1 #initiating KNN Basic with pearson correlation
2 knn_basic = KNNBasic(sim_options={'name':'pearson', 'user_based':False})
3 knn_basic.fit(trainset)
4 predictions = knn_basic.test(testset)
5 print(accuracy.rmse(predictions))
```

Computing the pearson similarity matrix...
Done computing similarity matrix.
RMSE: 0.9706
0.9706489484223293

```
In [29]: 1 #initiating KNN Basic with pearson correlation
2 knn_basic = KNNBasic(sim_options={'name':'cosine', 'user_based':False})
3 knn_basic.fit(trainset)
4 predictions = knn_basic.test(testset)
5 print(accuracy.rmse(predictions))
```

Computing the cosine similarity matrix...
Done computing similarity matrix.
RMSE: 0.9787
0.9787151489086232

Okay, so we tried to utilize both hyperparameters here, and we got a larger error. Nearly an entire point. We'll sidestep the cross-validation here and see if we can run a different neighborhood based model. This model utilizes ALS (Alternative Linear Squares) method. We'll try both options and see which is better.

```
In [30]: 1 # cross validating with KNNBaseline
2 knn_baseline = KNNBaseline(sim_options={'name':'pearson', 'user_based':True})
3 knn_baseline.fit(trainset)
4 predictions = knn_baseline.test(testset)
5 print(accuracy.rmse(predictions))
```

Estimating biases using als...
Computing the pearson similarity matrix...
Done computing similarity matrix.
RMSE: 0.8801
0.8801035004209004

```
In [31]: 1 # cross validating with KNNBaseline  
2 knn_baseline = KNNBaseline(sim_options={'name': 'cosine', 'user_based':True}  
3 knn_baseline.fit(trainset)  
4 predictions = knn_baseline.test(testset)  
5 print(accuracy.rmse(predictions))
```

```
Estimating biases using als...  
Computing the cosine similarity matrix...  
Done computing similarity matrix.  
RMSE: 0.8814  
0.8813916557808571
```

So, the KNN Baseline module performs better than the KNN Basic, but not better than the SVD.

Summary

The method with the lowest RMSE (0.859) was a user-based, SVD with tuned hyperparameters {n_factors=150, n_epochs=25, lr_all=0.010, reg_all=0.05}.

Let's go ahead and build our recommender using the SVD!!!

6. Implementation

Overview

Now that we have this model (Step 1), we will proceed to develop our recommender system, including how to work through the "cold start" problem. We will do that utilizing the following steps.

Step 1 (previously created): Prior to input from the new user, we created user-based collaborative filtering prediction model to predict how an existing user would rate a movie from the database.

Step 2: Prompt user to rate five movies.

Step 3: Add user's rating to the existing database

Step 4: Use the model from step 1 to predict how new users movie would rate (1-5) for all movies in the database and sort highest to lowest

Step 5: Output the top 5 recommendations

So... let's go to Step 2.

Step 2. Prompt User

Below is the function used to prompt a user to rate movies, (1 - 5) for random movies in the database.

In [32]:

```
1 #create function to be called based on the number of movies created. This will
2 def movie_rater(movie_df, num, last_user, genre=None):
3     userID = last_user + random.randint(0,1000)
4     rating_list = []
5
6     #Loop through for each recommendation
7     while num > 0:
8
9         if genre:
10             movie = movie_df[movie_df['genres'].str.contains(genre)].sample(1)
11         else:
12             movie = movie_df.sample(1)
13
14         print(f"\n {movie.title} {movie.genres}\n")
15         rating = input('How do you rate this movie on a scale of 1-5, press n to skip: ')
16
17         if rating == 'n':
18             continue
19         elif (0 < float(rating) and float(rating) < 5.1):
20             rating_one_movie = {'userId':userID,'movieId':movie['movieId']}
21             rating_list.append(rating_one_movie)
22             num -= 1
23         else:
24             rating_again = input("Please choose either n, if you haven't seen it or a rating between 1-5: ")
25             if rating_again == 'n':
26                 continue
27             elif (0 < float(rating) and float(rating) < 5.1):
28                 rating_one_movie = {'userId':userID,'movieId':movie['movieId']}
29                 rating_list.append(rating_one_movie)
30                 num -= 1
31             else:
32                 print("You're struggling with directions. Let's try a different approach")
33                 continue
34
35 return rating_list
```

Input

```
In [33]:
```

```
1 # Let's call our new function
2 last_user = ratings['userId'].max()
3 user_rating = movie_rater(movies_df, 5, last_user)
```

```
6542    Sydney White (2007)
Name: title, dtype: object 6542    Comedy
Name: genres, dtype: object
```

```
How do you rate this movie on a scale of 1-5, press n if you have not seen:
n
```

```
5153    Man Who Came to Dinner, The (1942)
Name: title, dtype: object 5153    Comedy
Name: genres, dtype: object
```

```
How do you rate this movie on a scale of 1-5, press n if you have not seen:
n
```

```
8655    John Mulaney: New In Town (2012)
Name: title, dtype: object 8655    Comedy
Name: genres, dtype: object
```

```
How do you rate this movie on a scale of 1-5, press n if you have not seen:
5
```

```
6895    Saw V (2008)
Name: title, dtype: object 6895    Crime|Horror|Thriller
Name: genres, dtype: object
```

```
How do you rate this movie on a scale of 1-5, press n if you have not seen:
4
```

```
9491    CHiPS (2017)
Name: title, dtype: object 9491    Action|Comedy|Drama
Name: genres, dtype: object
```

```
How do you rate this movie on a scale of 1-5, press n if you have not seen:
1
```

```
1984    Mummy, The (1959)
Name: title, dtype: object 1984    Horror
Name: genres, dtype: object
```

```
How do you rate this movie on a scale of 1-5, press n if you have not seen:
1
```

```
1261    Starship Troopers (1997)
Name: title, dtype: object 1261    Action|Sci-Fi
Name: genres, dtype: object
```

```
How do you rate this movie on a scale of 1-5, press n if you have not seen:
2
```

Okay, so we prompted the user and got their viewing history. Let's move on to Step 3.

Step 3. Add user ratings to database and rerun model

```
In [34]: 1 ## add the new ratings to the original ratings DataFrame  
2 user_ratings = pd.DataFrame(user_rating)  
3 new_ratings_df = pd.concat([ratings, user_ratings], axis=0)  
4 new_data = Dataset.load_from_df(new_ratings_df,reader)
```

Whelp... that was easy, we now have the user's information here in the database... Let's go to Step 4

Step 4: Predict new user movie preferences

Now that we have a "new" database, similar to the old one, let's rerun our prediction model.

First, we will rerun the model, and then we will create new predictions for all of the movie's in the database, sort from greatest to least.

```
In [36]: 1 # train a model using the new combined DataFrame, recall our parameters from  
2 svd = SVD(n_factors=150, n_epochs=25, lr_all=0.010, reg_all=0.05)  
3 svd.fit(new_data.build_full_trainset())
```

```
Out[36]: <surprise.prediction_algorithms.matrix_factorization.SVD at 0x2cfb3bc0520>
```

```
In [37]: 1 # make predictions for the user, to do this, predict ratings for every movie  
2 list_of_movies = []  
3  
4 for m_id in ratings['movieId'].unique():  
5     list_of_movies.append((m_id,svd.predict(last_user,m_id)[3]))  
6  
7 ranked_movies = sorted(list_of_movies, key=lambda x:x[1], reverse=True)
```

Step 5: Provide recommendations for 5 movies.

Now that we have a "new" database, similar to the old one, let's rerun our prediction model.

First, we will rerun the model, and then we will create new predictions for all of the movie's in the database, sort from greatest to least.

In [38]:

```
1 # return the top n recommendations using the
2 def recommended_movies(user_ratings,movie_title_df,n):
3     for idx, rec in enumerate(user_ratings):
4         title = movie_title_df.loc[movie_title_df['movieId'] == int(rec)]
5         print('Recommendation # ', idx+1, ': ', title, '\n')
6         n-= 1
7         if n == 0:
8             break
9
10 recommended_movies(ranked_movies,movies_df,5)
```

Recommendation # 1 : 659 Godfather, The (1972)
Name: title, dtype: object

Recommendation # 2 : 224 Star Wars: Episode IV - A New Hope (1977)
Name: title, dtype: object

Recommendation # 3 : 898 Star Wars: Episode V - The Empire Strikes Bac
k...
Name: title, dtype: object

Recommendation # 4 : 5621 Neon Genesis Evangelion: The End of Evangelio
n...
Name: title, dtype: object

Recommendation # 5 : 900 Raiders of the Lost Ark (Indiana Jones and th
e...
Name: title, dtype: object

And there we have it! I like all of these movies except for no. 4. I haven't seen it. Maybe I'll watch it later.

In []:

1