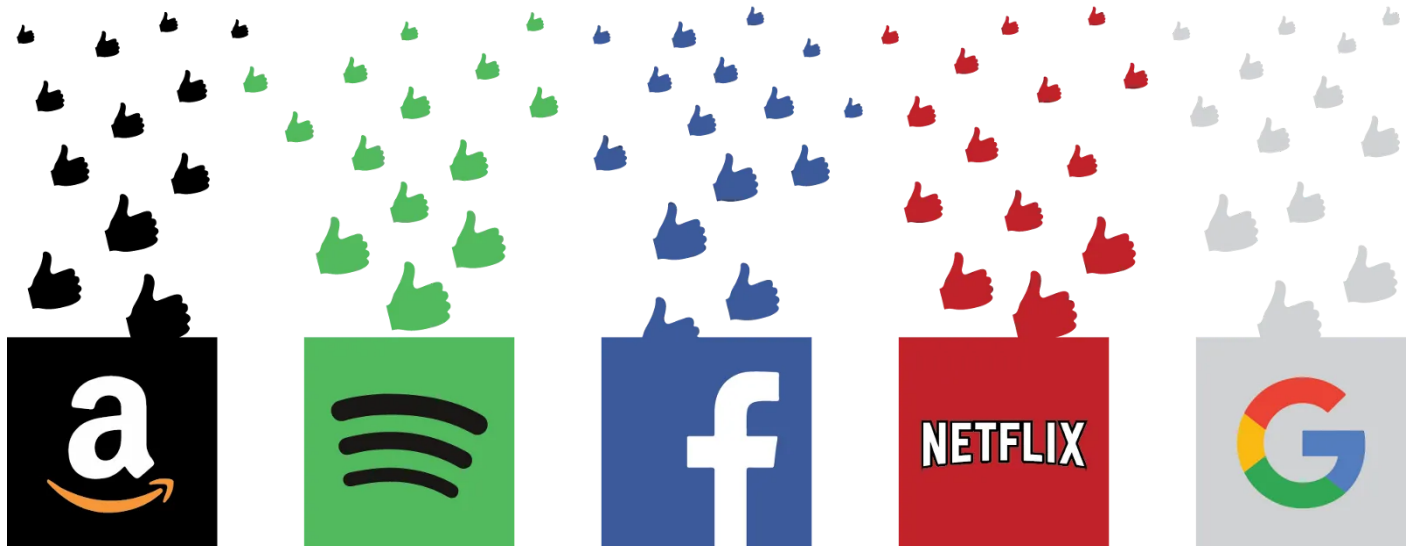


Project T Final: Collaborative Filtering

CS 189/289A: Introduction to Machine Learning, Fall 2020

Team MA: Maxwell Chen and Abinav Routhu



(The Data Scientist: What is the right way to build a recommender system for a startup? May 3, 2018)

Introduction

This notebook will cover various methods to construct a recommender system through the process of collaborative filtering -- algorithms and techniques that are concerned with finding similarities between users and items, and calculating numerical ratings to quantify this similarity. We will ground ourselves in a standard and accessible application of recommender systems -- recommending movies to users on a service such as Netflix.

Historical background

The modern age is undergoing rapid and intense changes due to the vast amounts of data being generated in the early 21st century thanks to the Information Age and inventions such as MOSFETs, digital electronics, and the internet. Understanding and leveraging said data has led to the rise of machine learning and data science -- practitioners of these fields are becoming indispensable to virtually every industry. One such "industry" we will focus on in particular is advertising and marketing, which have radically changed through the inception of recommender systems. We see this everywhere -- Amazon products, Spotify songs, Facebook friends, YouTube videos, Google ads -- all these companies are using the concept of learning from data to predict new products to users and customers.

The 2006 Netflix Prize was a \$1,000,000 challenge run by Netflix to find the best collaborative filtering algorithm that could improve Netflix movie recommendations. This ran for three years until a team comprised of research scientists bested Netflix's own prediction accuracy by 10.06%. The contest ignited interest in recommendation and perhaps led to the growth of machine learning competitions through website such as Kaggle.

Learning Objectives

This notebook serves to introduce and explore the topic of Collaborative Filtering through mathematical methods, along with practical application to the task of recommending movies to users.

Collaborative Filtering is a process or algorithm to filter information or patterns through the collaboration of multiple users, agents, or data sources.

We shall approach this through two paradigms:

1. Memory-Based Approaches (Clustering, KNN)
2. Model-Based Approaches (Matrix Factorization)

Table of Contents

- |—Introduction
- |—Table of Contents ☆
- |—Data Analysis
 - |—Q1: Loading the Dataset
 - |—Q2: Exploratory Data Analysis: Understanding and Visualizing the Dataset
- |—Q3: The User-Interaction Matrix
- |—Q4: Memory-Based Approaches
 - |—Q4a: Cosine Similarity
 - |—Q4b: K-Nearest Neighbors
- |—Q5: Model-Based Approaches
 - |—Gradient Descent
- |—Extensions (Optional)
 - |—Surpriselib Package
 - |—Deep Learning
 - |—Regularization
 - |—Issues
 - |—Cold Start

```
In [ ]: # Load Packages
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
```

To begin, we will load the MovieLens Dataset. MovieLens was a research project launched by GroupLens Research at the University of Minnesota, and was one of the earliest modern projects that investigated personalized recommendations via recommender systems. We will be using their Dataset for a similar purpose: recommending a user which movie to watch based on their own interests or preferences.

Question 1: Loading the Dataset

1.1: Import `zipfile` and from `urllib.request` import `urlretrieve` . Use these libraries to load the [MovieLens Dataset \(http://files.grouplens.org/datasets/movielens/ml-100k.zip\)](http://files.grouplens.org/datasets/movielens/ml-100k.zip) -- this is the "small" Dataset containing 100,000 ratings. If you are up for it, you can also load the [expanded MovieLens Dataset \(http://files.grouplens.org/datasets/movielens/ml-latest.zip\)](http://files.grouplens.org/datasets/movielens/ml-latest.zip) -- this contains 27,000,000 ratings. For the purposes of this assignment, loading either Dataset will work.

Reference Material:

- [DataCamp Tutorial on zipfile module \(https://www.datacamp.com/community/tutorials/zip-file\)](https://www.datacamp.com/community/tutorials/zip-file)
- [GeeksForGeeks Tutorial on zipfile module \(https://www.geeksforgeeks.org/working-zip-files-python/\)](https://www.geeksforgeeks.org/working-zip-files-python/)
- [urllib.request Documentation \(https://docs.python.org/3/library/urllib.request.html\)](https://docs.python.org/3/library/urllib.request.html)

```
In [ ]: from urllib.request import urlretrieve
import zipfile
#### BEGIN CODE ####
...
#### END CODE ####
```

1.2: We now have a raw .csv file containing our Dataset. As with many other problems involving machine learning or data mining, we must manipulate our raw data to a form that we can use.

First, investigate the structure of the zipped dataset we just downloaded. Open up each of the unzipped files on DataHub or your local machine, and describe the contents of each file:

Answer: ...

1.3: Use your knowledge of data cleaning and processing from the first week of 16ML to load the different .csv files into multiple Pandas DataFrames. Use the provided columns stored in `user_features` , `ratings_features` , and `movie_features` . Use appropriate naming conventions for these DataFrames, such as "movies", for example. Then combine the Dataframes into a single DataFrame, using `user_id` as a primary key.

Hint #1: When using `pd.read_csv` , you MUST use the flag `encoding='latin-1'` to properly read from the files.

Hint #2: Use `sep="|"` when reading in the csv file for users and movies, but `sep="\t"` for ratings

```
In [ ]: user_features = ["user_id", "age", "sex", "occupation", "zip_code"]
        ### BEGIN USERS CODE ###
        users = pd.read_csv(...)
        ### END USERS CODE ###

        ratings_features = ["user_id", "movie_id", "rating", "unix_timestamp"]
        ### BEGIN RATINGS CODE ###
        ratings = pd.read_csv(...)
        ### END RATINGS CODE ###

        movie_features = ['movie_id', 'title', 'release_date', "video_release_date",
                           "imdb_url", "genre_unknown", "Action",
                           "Adventure", "Animation", "Children", "Comedy", "Crime", "Do
                           cumentary", "Drama", "Fantasy",
                           "Film-Noir", "Horror", "Musical", "Mystery", "Romance", "Sci
                           -Fi", "Thriller", "War", "Western"]
        ### BEGIN MOVIES CODE ###
        movies = pd.read_csv(...)
        ### END MOVIES CODE ###

        ### BEGIN MERGE CODE ###
        all_data = ratings.merge(...).merge(...)
        ### END MERGE CODE ###
        all_data.head()
```

Question 2: Understanding and Visualizing the Dataset

2.1: Distribution of Movie Genres

Make a plot of the frequency of each distribution in the dataset. Refer back to material from Week 1 and 2 if you need a refresher on using Pandas and Matplotlib.

Hint: Try a bar plot

```
In [ ]: ### BEGIN CODE ###
movie_column_labels = ...
movie_genres = ...
genre_frequency = ...
### END CODE ###

plt.figure(figsize=(10, 10));
plt.title("Frequency of Genres in Movielens Dataset");
plt.xlabel("Genre");
plt.ylabel("Proportion");
genre_frequency.plot.bar();
```

2.2: It is important to identify biases in our dataset that can skew our results or impact how generalizable our recommendation system is to novel users and novel movies. What might be some issues we run into by using this dataset?

Answer: ...

2.3: Distribution of User Ratings

Plot the distribution of user ratings for movies from the Children, Fantasy, and Film-Noir genres -- that is, for each genre, plot a distribution describing the number of ratings from 1 to 5 received by movies belonging to that genre. Note any similarities or differences between your plots -- how does this inform us about biases in the dataset, and how could such bias affect our predictions?

[Hint: Try multiple histograms or bar plots]

```
In [ ]: ### BEGIN CODE ###
children_ratings = ...
fantasy_ratings = ...
film_noir_ratings = ...
### END CODE ###

### BEGIN PLOTTING ###
sns.histplot(data=children_ratings, label="Children")
plt.title("Children Movie Ratings")
plt.xlabel("Ratings")
plt.ylabel("Count")
plt.show()

sns.histplot(data=fantasy_ratings, label="Fantasy")
plt.title("Fantasy Movie Ratings")
plt.xlabel("Ratings")
plt.ylabel("Count")
plt.show()

sns.histplot(data=film_noir_ratings, label="Film Noir")
plt.title("Film Noir Movie Ratings")
plt.xlabel("Ratings")
plt.ylabel("Count")
plt.show()
### END PLOTTING ###
```

Answer: ...

Question 3: The User-Interaction Matrix

3.1: How many unique users and unique movies are there in our dataset? Assign your answers to `num_users` and `num_movies`, respectively.

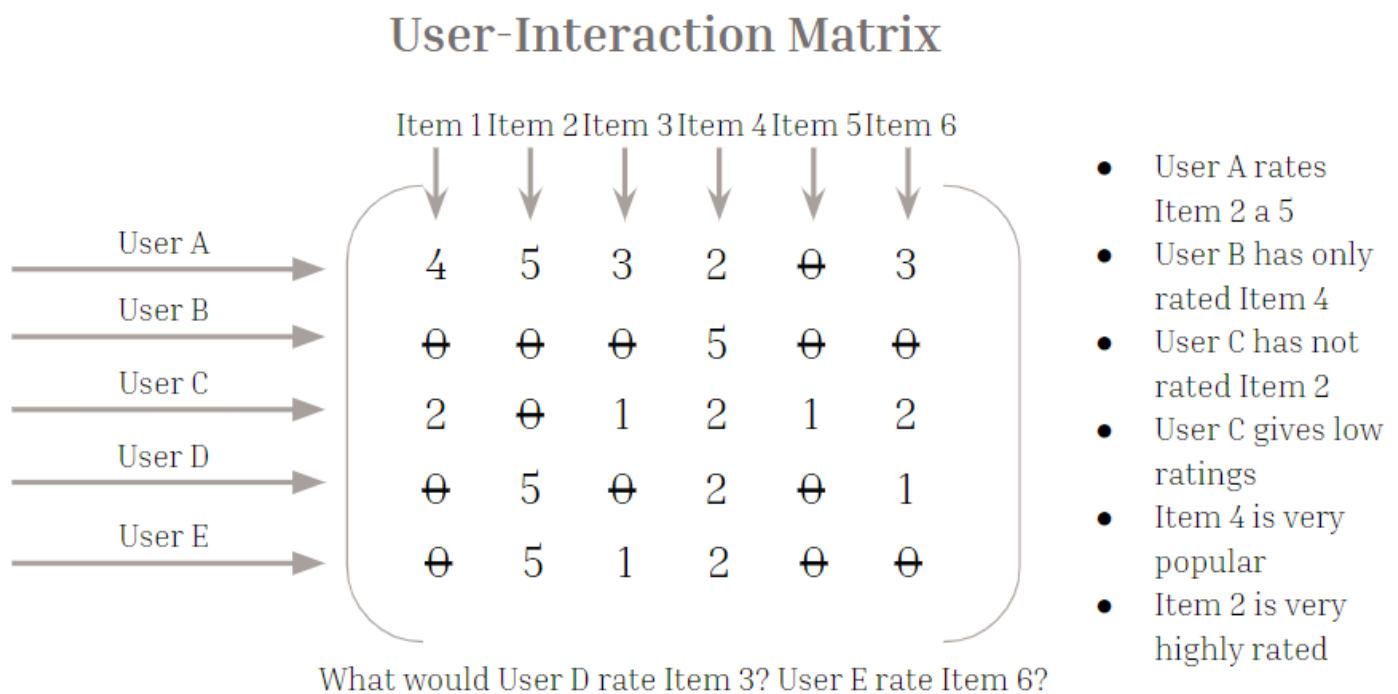
Hint: What data structure did you see this week in CS 61B that could be helpful here?

```
In [ ]: ### BEGIN CODE ###
unique_users = ...
unique_movies = ...
num_users = ...
num_movies = ...
### END CODE ###

print("Number of ratings:", len(all_data))
print("Number of unique users:", num_users)
print("Number of unique movies:", num_movies)
```

3.2: You should see that the values for the number of unique users and unique movies are much smaller than the dimensions of our raw data matrix. What does that tell us about how many movies each user rated? What would you expect to be the most common number in our raw data matrix?

Answer: ...



3.3: Recall the structure of the User-Item Interaction Matrix taken from this week's slides. For `all_data`, construct the corresponding User-Item Interaction Matrix using Pandas. Call it `interaction_matrix`. Print out its dimensions and the first few rows to confirm that the dimensions match with the number of unique users and movies you found in **3.1**.

Hint: Look into the Pandas function `df.pivot` (<https://pandas.pydata.org/pandas-docs/stable/reference/api/pandas.DataFrame.pivot.html>) (<https://pandas.pydata.org/pandas-docs/stable/reference/api/pandas.DataFrame.pivot.html>)

```
In [ ]: ### BEGIN CODE ###
truncated_data = ...
interaction_matrix = ...
### END CODE ###

print("Interaction Matrix Dimensions: ", interaction_matrix.shape);
display(interaction_matrix.head());
```


You should see that most of the values in the matrix are 0. This reflects the fact that the original `all_data` matrix was sparse.

Question 4: Memory-Based Approaches and Clustering

Question 4a: Cosine Similarity

To introduce the idea of comparing users and movies, we will look into a commonly-used similarity measure known as "cosine similarity". It boils down to determining the cosine of the angle between every pair of feature vectors, which can be expressed using the following equation:

For vectors i and j of length n ,

$$\cos(\theta) = \frac{v_i^T v_j}{|v_i| |v_j|} = \sum_{k=1}^n \frac{v_{i,k} v_{j,k}}{\sqrt{\sum_{k=1}^n v_{i,k}^2} \sqrt{\sum_{k=1}^n v_{j,k}^2}}$$

4.1: Fill out the following function `cosine_similarity(U, V)` to compute the cosine similarity of the vectors in a matrix:

Note: You can use `np.finfo(float).eps` to add a "fudge" factor and prevent issues with ratings of 0.

```
In [ ]: def cosine_similarity(U):  
    '''  
    Inputs:  
        - U: Data of interest represented as vectors in a matrix  
    Output:  
        - similarity: Matrix where a_{i,j} represents the cosine similarity be  
tween vectors v_i and v_j  
    '''  
    ### BEGIN CODE ###  
    similarity = ...  
    ### END CODE ###  
  
    return similarity
```

4.2: Use the completed function to find the cosine similarity between users. Assign this to `user_similarity`.

Hint: Treat each row of `interaction_matrix` as a vector.

```
In [ ]: ### BEGIN CODE ###
        user_similarity = ...
        ### END CODE ###

        user_similarity.head()
```

4.3: What does each value in the matrix represent? Why is the diagonal 1_m , i.e., why are there 1's along the diagonal?

Answer: ...

4.4: How could we now use this matrix to recommend movies to a user? What drawbacks are there with using this approach?

Answer: ...

4.5: How can we change the input to `cosine_similarity` to compute the cosine similarity between movies, rather than users?

Answer: ...

Question 4b: K-Nearest Neighbors

When we recommend similar users based on the ratings they give, that is a type of **user-user collaborative filtering**. Alternatively, when we look for similar items based on the ratings given to them by users, that is a type of **item-item collaborative filtering**.

We've just seen that we can use similarity functions to quantify the similarity between users or items. The next step is to make a recommendation by grouping together users who are most similar to one another. In this way, we create groups, or clusters, that represent users who give similar ratings. This boils down to the algorithm known as **K-Nearest Neighbors (KNN)**, which is used to partition data into K clusters of greatest similarity. (This is spiritually similar to the K-Means Clustering you saw in lecture earlier this week -- K-Means Clustering is an unsupervised learning technique that assigns the data into K clusters, whereas KNN looks at the K data points most similar to a certain training point in order to assign it a class or label).

For the purposes of this assignment, we will look at a pre-implemented [KNN Algorithm \(https://scikit-learn.org/stable/modules/generated/sklearn.neighbors.NearestNeighbors.html#sklearn.neighbors.NearestNeighbors\)](https://scikit-learn.org/stable/modules/generated/sklearn.neighbors.NearestNeighbors.html#sklearn.neighbors.NearestNeighbors) from [scikit-learn \(https://scikit-learn.org/\)](https://scikit-learn.org/), the machine learning library we have been working with over the last few weeks of 16ML. The KNN Algorithm we use here is a bit simpler than what KNN really is -- here, instead of assigning labels based on the K nearest neighbors, we will use this to simply identify the K nearest neighbors.

[Note: Later on in EECS 16B, you will revisit K-Means Clustering. This will be used in lab to classify voice commands to control your car.]

4.6: Read the linked documentation to understand how to use the KNN Algorithm from Scikit-Learn. What are the relevant functions and return values we can use?

Answer: ...

4.7: Run KNN on `interaction_matrix` from the earlier part of this question. For now, use $k = 10$.

```
In [ ]: ### BEGIN CODE ###
distances, indices = ...
indices += 1
### END CODE ###

print("Shape of distances matrix:", distances.shape)
display(distances)
print("Shape of indices matrix:", indices.shape)
display(indices)
```

4.8: Let's look at the first user in our `interaction_matrix`. What other users is our user most similar to? What are the IDs of the 5 movies these users liked most?

```
In [ ]: ### BEGIN CODE ###
        first_user = ...
        similar_users = ...

        print("IDs of top movies:", ...)
        ### END CODE ###

        print("First user is similar to:", first_user)
        display(similar_users)
```

As seen above, this approach with similarity metrics and K-Nearest Neighbors allows us to identify the users any given user is most similar to, and then look at those users' ratings to determine what that given user should watch next.

Question 5: Model-Based Approaches -- Matrix Factorization and Embeddings

Now, we will investigate another paradigm of collaborative filtering: Model-Based Approaches. These methods instead leverage matrix factorization to break apart matrices into other matrices or vectors that have special meanings or interpretations.

5.1: What Matrix Decompositions you have seen in EECS 16A and EECS 16B? What are the tradeoffs or different use cases for one over the other?

Answer: ...

Here is a quick refresher on the SVD:

For a matrix $A \in \mathbb{R}^{m \times n}$, the "Full SVD" is the following matrix product:

$$A = U \Sigma V^T$$

where

$$U \in \mathbb{R}^{m \times m}$$

$$\Sigma \in \mathbb{R}^{m \times n}$$

$$V^T \in \mathbb{R}^{n \times n}$$

Alternatively, there is a "Compact SVD" that involves truncating these matrices to remove zero-value singular vectors corresponding to the Nullspace of A :

$$A = U_c \Sigma_c V_c^T$$

where

$$U_c \in \mathbb{R}^{m \times c}$$

$$\Sigma_c \in \mathbb{R}^{c \times c}$$

$$V_c^T \in \mathbb{R}^{n \times c}$$

As popularized during the 2006 Netflix Prize, there is a famed "SVD" Matrix Factorization Algorithm used for creating a recommender system. **NOTE: THIS IS A DIFFERENT SVD!** Even though they are both called "SVD", and are spiritually related in the sense that they are both related to matrix decomposition / factorization, they are in fact, different algorithms entirely. They are related in the sense that, given full information of A , the outputs of this process should converge to the results we would obtain from SVD.

Our goal will be to factorize `interaction_matrix` and **approximate** it as the product $A = UV$ (You may think of the sigma matrix Σ as being "absorbed" into either the U or V matrix). For a matrix $A \in \mathbb{R}^{m \times n}$, it will be equal to the product $U \in \mathbb{R}^{m \times d}$ and $V \in \mathbb{R}^{d \times n}$.

In the context of movie recommendations, U is referred to as the **User Embeddings**, while V is referred to as the **Movie Embeddings**. Each row of U , represented by U_i , denotes the "essence" of user i , while each column of V , represented by V_j , denotes the "essence" of movie j .

5.2: How can we approximate the matrix A using a sum instead of a matrix product?

Answer: ...

5.3: Why can't we use the actual SVD to compute U and V ?

Answer: ...

5.4: Calculate A^* , the approximation for `interaction_matrix`, using gradient descent to learn randomly-initialized embedded vectors for U and V . Assume a latent space of dimension $d = 2$, i.e., that our embedded vectors are of length 2.

```
In [ ]: ### BEGIN CODE ###  
U = ...  
V = ...  
### END CODE ###
```

```
In [ ]: gamma = 0.04  
  
### BEGIN CODE ###  
  
### END CODE ###
```

```
In [ ]: A_star = np.dot(U, V.T)
```

You should see that `A_star - interaction_matrix` returns a matrix where all the values are fairly close to 0. This symbolizes the fact that we have attained appropriate embedding vectors, and thus appropriate user/item embeddings U and V , that best approximate our interaction matrix.

To go further with these embeddings, we can return to the formula in **5.2** and compute the expected rating for a given user for a given movie, by computing $a_{i,j} = \overset{\rightarrow{U}}{u_i} \cdot \overset{\rightarrow{V}}{v_j}$.

Summary and Extensions

Congratulations! We've reached the end of the assignment. We hope that this was an eye-opening experience into a very applied topic of interest in machine learning. To summarize what we've covered so far:

- **History of Recommender Systems and the Netflix Prize**
- **Loading Datasets**
- **Exploratory Data Analysis**
- **User-Interaction Matrices**
- **Memory-Based Approaches** - We looked at how we could use similarity measurements such as cosine similarity to determine the similarity between users/items in our user-interaction matrix. Using such measurements allows us to use algorithms such as K-Nearest Neighbors to identify the users or items most similar to a given user or item, respectively.
 - Cosine Similarity
 - KNN
- **Model-Based Approaches** - We looked at another paradigm in collaborative filtering that leverages matrix decompositions. The "SVD" Matrix Factorization Algorithm (which isn't actually the SVD!) was used to decompose our interaction matrix into user and item embeddings. We then generated random embeddings and learned correct values for them by using gradient descent.
 - User and Movie Embeddings
 - Gradient Descent

If you are interested in further investigating collaborative filtering algorithms and building recommender systems, here are some relevant links to sophisticated packages used in production settings, cutting-edge improvements, and current, unsolved issues:

Surpriselib

Surprise is a Package for SciPy. It is a dedicated SciPy package for building and analyzing Recommender Systems. In particular, it has native access to various datasets such as MovieLens, and also has a wide array of prediction algorithms.

- [SurpriseLib Website \(http://surpriselib.com/\)](http://surpriselib.com/)
- [GitHub Link \(https://github.com/NicolasHug/Surprise\)](https://github.com/NicolasHug/Surprise)

Here is a demo for using Surprise to perform prediction using cosine similarity and KNN on a small example training set:

```
In [ ]: !pip install scikit-surprise
```

```

In [ ]: from surprise import Dataset
        from surprise import Reader
        from surprise import KNNWithMeans

        '''
        Surprise_Ratings contains mappings of (item, user, rating) pairs
        - item: Name/ID for item
        - user: Name of user
        - rating: Rating user assigns to corresponding item
        '''

surprise_ratings = {
    "item": ["A", "B",
             "A", "B",
             "A", "B",
             "A", "B",
             "A"],
    "user": ["Allen", "Allen",
             "Bill", "Bill",
             "Cathy", "Cathy",
             "Devin", "Devin",
             "Evan"],
    "rating": [1, 2,
               2, 4,
               2.5, 4,
               4.5, 5,
               3],
}

reader = Reader(rating_scale=(1, 5))
df = pd.DataFrame(surprise_ratings)

# Load elements from surprise_ratings into form that Surpriselib can use
surprise_data = Dataset.load_from_df(df[["user", "item", "rating"]], reader)

sim_options = {
    "name": "cosine", # Use item-based cosine similarity
    "user_based": False, # Flag to determine whether to do user-user or item-i
    tem similarity -- here we do item-item similarity
}

# Generate KNN Algorithm and train on data to create predictor
surprise_algorithm = KNNWithMeans(sim_options=sim_options)
surprise_training_data = surprise_data.build_full_trainset()
surprise_algorithm.fit(surprise_training_data)

# Predict Evan's rating for item "B"
surprise_prediction = algo.predict("Evan", "B")
print(surprise_prediction.est)

```

Alternative and Advanced Methods

Deep Learning and Deep Neural Networks

Later on in 16ML, we will introduce the idea of Deep Learning, which relies heavily on using Neural Networks with many, many layers to perform complex computations. These are considered state-of-the-art models that are at the forefront of recent advances in many different fields of Machine Learning, and are a valid option for developing a recommender system.

Regularization

When performing Model-Based Collaborative Filtering, we can run into the issue of underregularization -- in other words, the way we perform such similarity measurements and clustering techniques make it so that we are very likely to overfit. This is a concept that was discussed earlier in 16ML during Week 4. In essence, regularization allows us to "smooth" out our data by "lifting" up small, near-zero values in our data that can cause numerical instability and/or unexpected magnification of small values. Regularization can help with Collaborative Filtering by dealing with movies that were not rated by many users, or users who did not rate many movies. In this way, it can help reduce error in predicting or recommending new movies to users, or conversely, identifying new users who might like a given movie.

Issues

Cold Start

A fundamental flaw with recommender systems as they currently exist is a concept called "Cold Start." Consider the following: we've been working with a massive matrix containing ratings given by various users to various movies. What happens when we add a new movie to the database? Any new movies will have no ratings from any user, meaning that our collaborative filtering algorithms have no information whatsoever on how a user will like it. That means that these movies will never be recommended with a naive algorithm. This means that we must carefully consider how we add new movies to our database. Some considerations include assigning it the same user ratings as a movie we believe to be very similar, or perhaps using some metric such as median or average to get a representative value of the new movie's genre. We could also apply content-based filtering when first looking at a new movie by using certain factors, such as genre or length, to establish initial ratings.

As we saw earlier, the matrices we are working with are very sparse, which will also inherently contribute to the Cold Start problem.

References

- The Netflix Prize. https://en.wikipedia.org/wiki/Netflix_Prize (https://en.wikipedia.org/wiki/Netflix_Prize)
- Baptiste Rocca. Introduction to recommender systems. <https://towardsdatascience.com/introduction-to-recommender-systems-6c66cf15ada,2019> (<https://towardsdatascience.com/introduction-to-recommender-systems-6c66cf15ada,2019>).
- Simon Funk. Netflix update: Try this at home. <https://sifter.org/simon/journal/20061211.html> (<https://sifter.org/simon/journal/20061211.html>), 2006.
- Yehuda Koren. The bellkor solution to the netflix grand prize. Published on Netflix PrizeForums, 2009.
- Build a Recommendation Engine With Collaborative Filtering. <https://realpython.com/build-recommendation-engine-collaborative-filtering/> (<https://realpython.com/build-recommendation-engine-collaborative-filtering/>)
- Prince Grover. Various Implementations of Collaborative Filtering <https://towardsdatascience.com/various-implementations-of-collaborative-filtering-100385c6dfe0> (<https://towardsdatascience.com/various-implementations-of-collaborative-filtering-100385c6dfe0>)
- Intro to Recommender Systems: Collaborative Filtering <https://www.ethanrosenthal.com/2015/11/02/intro-to-collaborative-filtering/> (<https://www.ethanrosenthal.com/2015/11/02/intro-to-collaborative-filtering/>)
- Alesha Tony. All You Need to Know About Collaborative Filtering <https://www.digitalvidya.com/blog/collaborative-filtering/> (<https://www.digitalvidya.com/blog/collaborative-filtering/>)
- MovieLens Dataset <https://grouplens.org/datasets/movielens/> (<https://grouplens.org/datasets/movielens/>)
- Surprise: A Python scikit for recommender systems. <http://surpriselib.com/> (<http://surpriselib.com/>)
- Movies Recommender System with Surpriselib. <https://medium.com/analytics-vidhya/movies-recommender-system-with-surpriselib-33580ae9b47c> (<https://medium.com/analytics-vidhya/movies-recommender-system-with-surpriselib-33580ae9b47c>)