

<https://github.com/topspinj/recommender-tutorial/blob/master/part-1-item-item-recommender.ipynb>

```
import numpy as np
import pandas as pd
import sklearn
import matplotlib.pyplot as plt
import seaborn as sns
```

```
import warnings
warnings.simplefilter(action='ignore', category=FutureWarning)
```

```
ratings = pd.read_csv("https://s3-us-west-2.amazonaws.com/recommender-tutorial/ratings.csv")
ratings.head()
```

```
movies =
pd.read_csv("https://s3-us-west-2.amazonaws.com/recommender-tutorial/movies.csv")
movies.head()
```

```
n_ratings = len(ratings)
n_movies = ratings['movieId'].nunique()
n_users = ratings['userId'].nunique()
print(f"Number of ratings: {n_ratings}")
print(f"Number of unique movieId's: {n_movies}")
print(f"Number of unique users: {n_users}")
print(f"Average number of ratings per user: {round(n_ratings/n_users, 2)}")
print(f"Average number of ratings per movie: {round(n_ratings/n_movies, 2)}")
```

```
user_freq = ratings[['userId', 'movieId']].groupby('userId').count().reset_index()
user_freq.columns = ['userId', 'n_ratings']
user_freq.head()
```

```
mean_rating = ratings.groupby('movieId')['rating'].mean()
lowest Rated = mean_rating['rating'].idxmin()
movies.loc[movies['movieId'] == lowest Rated]
```

```
highest Rated = mean_rating['rating'].idxmax()
movies.loc[movies['movieId'] == highest Rated]
```

```
ratings[ratings['movieId']==highest_rated]
```

```
ratings[ratings['movieId']==lowest_rated]
```

```
movie_stats = ratings.groupby('movieId')[['rating']].agg(['count', 'mean'])
```

```
movie_stats.columns = movie_stats.columns.droplevel()
```

```
from scipy.sparse import csr_matrix
```

```
ratings[ratings['movieId']==highest_rated]
```

```
ratings[ratings['movieId']==lowest_rated]
```

```
movie_stats = ratings.groupby('movieId')[['rating']].agg(['count', 'mean'])
```

```
movie_stats.columns = movie_stats.columns.droplevel()
```

```
from scipy.sparse import csr_matrix
```

```
def create_X(df):
```

```
    N = df['userId'].nunique()
```

```
    M = df['movieId'].nunique()
```

```
    user_mapper = dict(zip(np.unique(df["userId"]), list(range(N))))
```

```
    movie_mapper = dict(zip(np.unique(df["movieId"]), list(range(M))))
```

```
    user_inv_mapper = dict(zip(list(range(N)), np.unique(df["userId"])))
```

```
    movie_inv_mapper = dict(zip(list(range(M)), np.unique(df["movieId"])))
```

```
    user_index = [user_mapper[i] for i in df['userId']]
```

```
    movie_index = [movie_mapper[i] for i in df['movieId']]
```

```
    X = csr_matrix((df["rating"], (movie_index, user_index)), shape=(M, N))
```

```
    return X, user_mapper, movie_mapper, user_inv_mapper, movie_inv_mapper
```

```
X, user_mapper, movie_mapper, user_inv_mapper, movie_inv_mapper = create_X(ratings)
```

```
from sklearn.neighbors import NearestNeighbors
```

```
def find_similar_movies(movie_id, X, k, metric='cosine', show_distance=False):
```

```
    """
```

```
    Finds k-nearest neighbours for a given movie id.
```

```
    """
```

```
    neighbour_ids = []
```

```
    movie_ind = movie_mapper[movie_id]
```

```
    movie_vec = X[movie_ind]
```

```
    k+=1
```

```
    kNN = NearestNeighbors(n_neighbors=k, algorithm="brute", metric=metric)
```

```
    kNN.fit(X)
```

```
    if isinstance(movie_vec, (np.ndarray)):
```

```
        movie_vec = movie_vec.reshape(1,-1)
```

```
    neighbour = kNN.kneighbors(movie_vec, return_distance=show_distance)
```

```
    for i in range(0,k):
```

```
        n = neighbour.item(i)
```

```
        neighbour_ids.append(movie_inv_mapper[n])
```

```
    neighbour_ids.pop(0)
```

```
    return neighbour_ids
```

```
movie_titles = dict(zip(movies['movieId'], movies['title']))
```

```
movie_id = 1
```

```
similar_ids = find_similar_movies(movie_id, X, k=10)
```

```
movie_title = movie_titles[movie_id]
```

```
print(f"Because you watched {movie_title}")
```

```
for i in similar_ids:
```

```
    print(movie_titles[i])
```