'users.dat', sep='::', names=['UserID','Gender','Age', 'Occupation', 'zip-code'], engine='python', header=**None** ratings_df = pd.read_csv('ratings.dat', sep='::', names=['UserID', 'MovieID', 'Rating', 'Timestamp'], parse_dates=['Timestamp'], engine='python', header=**None** In [9]: movies_df.head() Out[9]: MovieID Title Genres 0 Toy Story (1995) Animation|Children's|Comedy 1 2 Jumanji (1995) Adventure|Children's|Fantasy 2 3 Grumpier Old Men (1995) Comedy|Romance 3 Waiting to Exhale (1995) Comedy|Drama 5 Father of the Bride Part II (1995) Comedy In [10]: users_df.head() Out[10]: UserID Gender Age Occupation zip-code 48067 M 56 70072 16 M 25 55117 M 45 02460 55455 In [11]: ratings_df.head() Out[11]: UserID MovieID Rating Timestamp 1193 5 978300760 661 3 978302109 3 978301968 914 4 978300275 5 978824291 2355 In [12]: ###Create a new dataset [Master_Data] with the following columns: ###MovieID Title, UserID, Age, Gender, Occupation, Rating. ###(Hint: (i) Merge two tables at a time. (ii) Merge the tables using two primary keys Movie ID & UserId) movie_ratings_df = pd.merge(movies_df, ratings_df, on='MovieID') movie_ratings_df.info() <class 'pandas.core.frame.DataFrame'> Int64Index: 1000209 entries, 0 to 1000208 Data columns (total 6 columns): Column Non-Null Count Dtype MovieID 1000209 non-null int64 1000209 non-null object Title 1 1000209 non-null object Genres UserID 1000209 non-null int64 3 4 Rating 1000209 non-null int64 Timestamp 1000209 non-null int64 5 dtypes: int64(4), object(2) memory usage: 53.4+ MB In [13]: movie_ratings_df.head() Out[13]: MovieID Title Genres UserID Rating Timestamp 1 Toy Story (1995) Animation|Children's|Comedy 5 978824268 1 1 Toy Story (1995) Animation|Children's|Comedy 4 978237008 2 1 Toy Story (1995) Animation|Children's|Comedy 4 978233496 3 1 Toy Story (1995) Animation|Children's|Comedy 9 5 978225952 5 978226474 1 Toy Story (1995) Animation|Children's|Comedy 10 In [14]: movie_ratings_users_df = pd.merge(movie_ratings_df, users_df, on='UserID' movie_ratings_users_df.info() <class 'pandas.core.frame.DataFrame'> Int64Index: 1000209 entries, 0 to 1000208 Data columns (total 10 columns): Column Non-Null Count # Dtype -----MovieID 1000209 non-null int64 Title 1000209 non-null object 1 1000209 non-null object Genres 2 3 UserID 1000209 non-null int64 4 Rating 1000209 non-null int64 1000209 non-null int64 5 Timestamp 1000209 non-null object 6 Gender 7 Age 1000209 non-null int64 Occupation 1000209 non-null int64 1000209 non-null object zip-code dtypes: int64(6), object(4) memory usage: 83.9+ MB In [15]: movie_ratings_users_df.head() Out[15]: zip-Genres UserID Rating Timestamp Gender Age Occupation MovielD Title Toy Story Animation|Children's|Comedy 5 978824268 10 48067 (1995)**Pocahontas** 48 Animation|Children's|Musical|Romance 5 978824351 1 10 48067 (1995)Apollo 13 150 Drama 5 978301777 10 48067 (1995)Star Wars: Episode IV - A 260 Action|Adventure|Fantasy|Sci-Fi 4 978300760 10 48067 New Hope (1977)Schindler's List Drama|War 5 978824195 10 48067 (1993)In [16]: Master_Data = movie_ratings_users_df.drop(['zip-code', 'Timestamp'], axis=1 Master_Data.head() Out[16]: MovieID Title Genres UserID Rating Gender Age Occupation Animation|Children's|Comedv 0 Toy Story (1995) 10 1 1 48 Pocahontas (1995) Animation|Children's|Musical|Romance 5 10 2 150 Apollo 13 (1995) 10 Drama 5 1 Star Wars: Episode IV - A New Hope 260 Action|Adventure|Fantasy|Sci-Fi 3 4 1 10 (1977)Schindler's List (1993) 527 Drama|War 10 Master_Data.describe(include='all') Out[17]: Occupation MovieID Title Genres UserID Rating Gender Age count 1.000209e+06 1000209 1000209 1.000209e+06 1.000209e+06 1000209 1.000209e+06 1.000209e+06 unique NaN 3706 301 NaN NaN 2 NaN NaN Comedy NaN American Beauty (1999) NaN Μ NaN NaN NaN top freq 3428 116883 NaN NaN 753769 NaN NaN mean 1.865540e+03 NaN 3.024512e+03 3.581564e+00 NaN 2.973831e+01 8.036138e+00 NaN std 1.096041e+03 NaN NaN 1.728413e+03 1.117102e+00 1.175198e+01 6.531336e+00 min 1.000000e+00 NaN 1.000000e+00 1.000000e+00 NaN 1.000000e+00 0.000000e+00 NaN 25% 1.030000e+03 NaN 1.506000e+03 3.000000e+00 2.500000e+01 2.000000e+00 **50%** 1.835000e+03 NaN 3.070000e+03 4.000000e+00 NaN 2.500000e+01 7.000000e+00 NaN **75%** 2.770000e+03 NaN 4.476000e+03 4.000000e+00 3.500000e+01 1.400000e+01 max 3.952000e+03 NaN NaN 6.040000e+03 5.000000e+00 NaN 5.600000e+01 2.000000e+01 In [18]: Master_Data.info() <class 'pandas.core.frame.DataFrame'> Int64Index: 1000209 entries, 0 to 1000208 Data columns (total 8 columns): Column Non-Null Count Dtype 0 1000209 non-null int64 MovieID 1 Title 1000209 non-null object Genres 1000209 non-null object UserID 1000209 non-null int64 3 Rating 1000209 non-null int64 Gender 1000209 non-null object Age 1000209 non-null int64 Occupation 1000209 non-null int64 dtypes: int64(5), object(3) memory usage: 68.7+ MB In [19]: ### Explore the datasets using visual representations (graphs or tables), ###also include your comments on the following: ###1. User Age Distribution In [20]: plt.figure(figsize=(8,6)) users_df.Age.hist() plt.title('User Age Distribution') plt.xlabel('Age') plt.ylabel('Number of Users') plt.show() User Age Distribution 2000 1750 1500 of Users 혈 1000 750 500 250 20 30 In [24]: ###User rating of the movie "Toy Story" plt.figure(figsize=(8,6)) movies_grouped = movie_ratings_df.groupby('Title') toy_story = movies_grouped.get_group('Toy Story (1995)') toy_story['Rating'].hist() plt.title('User rating of the movie "Toy Story"') plt.xlabel('Rating') plt.ylabel('Number of Users') plt.show() User rating of the movie "Toy Story" 800 700 600 of Users 500 Number 400 300 200 100 In [25]: ###Top 25 movies by viewership rating rating_avg = movie_ratings_df.groupby('Title')['Rating'].mean() rating_avg.head() Out[25]: Title \$1,000,000 Duck (1971) 3.027027 'Night Mother (1986) 3.371429 'Til There Was You (1997) 2.692308 'burbs, The (1989) 2.910891 ...And Justice for All (1979) 3.713568 Name: Rating, dtype: float64 In [26]: rating avg = rating avg.sort values(ascending=False) rating_avg.head() Out[26]: Title Gate of Heavenly Peace, The (1995) 5.0 Lured (1947) 5.0 Ulysses (Ulisse) (1954) 5.0 Smashing Time (1967) 5.0 Follow the Bitch (1998) 5.0 Name: Rating, dtype: float64 In [28]: rating_count = movie_ratings_df.groupby('Title')['Rating'] rating_count = rating_count.count().sort_values(ascending=False) rating_count[:25] Out[28]: Title 3428 American Beauty (1999) Star Wars: Episode IV - A New Hope (1977) 2991 Star Wars: Episode V - The Empire Strikes Back (1980) 2990 Star Wars: Episode VI - Return of the Jedi (1983) 2883 Jurassic Park (1993) 2672 Saving Private Ryan (1998) 2653 Terminator 2: Judgment Day (1991) 2649 Matrix, The (1999) 2590 2583 Back to the Future (1985) Silence of the Lambs, The (1991) 2578 Men in Black (1997) 2538 Raiders of the Lost Ark (1981) 2514 Fargo (1996) 2513 Sixth Sense, The (1999) 2459 Braveheart (1995) 2443 Shakespeare in Love (1998) 2369 Princess Bride, The (1987) 2318 Schindler's List (1993) 2304 L.A. Confidential (1997) 2288 Groundhog Day (1993) 2278 E.T. the Extra-Terrestrial (1982) 2269 Star Wars: Episode I - The Phantom Menace (1999) 2250 Being John Malkovich (1999) 2241 Shawshank Redemption, The (1994) 2227 Godfather, The (1972) 2223 Name: Rating, dtype: int64 In [29]: rating_avg_count = pd.DataFrame(data=rating_avg) rating_avg_count['number_of_ratings'] = pd.DataFrame(rating_count) rating_avg_count.head() Out[29]: Rating number_of_ratings **Gate of Heavenly Peace, The (1995)** 5.0 Lured (1947) 5.0 1 Ulysses (Ulisse) (1954) 5.0 **Smashing Time (1967)** Follow the Bitch (1998) 5.0 In [30]: rating_avg_count.describe() Out[30]: Rating number_of_ratings count 3706.000000 3706.000000 3.238892 269.889099 mean std 0.672925 384.047838 1.000000 1.000000 min **25**% 2.822705 33.000000 **50%** 3.331546 123.500000 **75**% 3.740741 350.000000 5.000000 3428.000000 max In [31]: | filter_data = rating_avg_count[rating_avg_count['number_of_ratings'] > 10] filter_data[:25] Out[31]: Rating number_of_ratings Title 69 Sanjuro (1962) 4.608696 Seven Samurai (The Magnificent Seven) (Shichinin no samurai) (1954) 4.560510 628 Shawshank Redemption, The (1994) 4.554558 2227 Godfather, The (1972) 4.524966 2223 Close Shave, A (1995) 4.520548 657 Usual Suspects, The (1995) 4.517106 1783 Schindler's List (1993) 4.510417 2304 Wrong Trousers, The (1993) 4.507937 882 Sunset Blvd. (a.k.a. Sunset Boulevard) (1950) 4.491489 470 **Raiders of the Lost Ark (1981)** 4.477725 2514 **Rear Window (1954)** 4.476190 1050 Paths of Glory (1957) 4.473913 230 **Star Wars: Episode IV - A New Hope (1977)** 4.453694 2991 Third Man, The (1949) 4.452083 480 Dr. Strangelove or: How I Learned to Stop Worrying and Love the Bomb (1963) 4.449890 1367 For All Mankind (1989) 4.444444 27 Wallace & Gromit: The Best of Aardman Animation (1996) 4.426941 438 **To Kill a Mockingbird (1962)** 4.425647 928 **Double Indemnity (1944)** 4.415608 551 Casablanca (1942) 4.412822 1669 World of Apu, The (Apur Sansar) (1959) 4.410714 56 2459 Sixth Sense, The (1999) 4.406263 **Yojimbo (1961)** 4.404651 215 Pather Panchali (1955) 4.404255 47 **Lawrence of Arabia (1962)** 4.401925 831 In [32]: ###Find the ratings for all the movies reviewed by for a particular user of user id = 2696 user_2696 = movie_ratings_users_df[movie_ratings_users_df['UserID'] == 2696] user_2696 Out[32]: zip-MovieID Title Genres UserID Rating Timestamp Gender Age Occupation code Client, The M 25 991035 350 Drama|Mystery|Thriller 2696 3 973308886 7 24210 (1994)Lone Star 991036 800 Drama|Mystery 2696 5 973308842 M 25 7 24210 (1996)**Basic Instinct** 991037 1092 Mystery|Thriller 2696 4 973308886 M 25 7 24210 (1992)E.T. the Extra-Children's|Drama|Fantasy|Sci-2696 991038 1097 Terrestrial 3 973308690 M 25 7 24210 (1982)Shining, The 1258 991039 Horror 2696 4 973308710 M 25 7 24210 (1980)Back to the 1270 2 973308676 991040 Comedy|Sci-Fi 2696 M 25 7 24210 Future (1985) Cop Land 991041 1589 Crime|Drama|Mystery 2696 3 973308865 M 25 7 24210 (1997)L.A. Confidential Crime|Film-1617 991042 2696 4 973308842 M 25 7 24210 Noir|Mystery|Thriller (1997)Game, The 1625 991043 Mystery|Thriller 2696 4 973308842 M 25 7 24210 (1997)I Know What 1644 2 973308920 7 24210 991044 You Did Last Horror|Mystery|Thriller 2696 M 25 Summer (1997) Devil's 991045 1645 Advocate, The Crime|Horror|Mystery|Thriller 2696 4 973308904 M 25 7 24210 (1997)Midnight in the 1711 Comedy|Crime|Drama|Mystery 2696 4 973308904 M 25 7 24210 991046 Garden of Good and Evil (1997) 991047 1783 2696 4 973308865 7 24210 Palmetto (1998) Film-Noir|Mystery|Thriller M 25 Wild Things 1805 991048 2696 4 973308886 7 24210 Crime|Drama|Mystery|Thriller M 25 (1998)Perfect Murder, 1892 991049 2696 4 973308904 7 24210 Mystery|Thriller 25 A (1998) I Still Know What You Did 991050 2338 2696 2 973308920 7 24210 Horror|Mystery|Thriller M 25 Last Summer (1998)2389 991051 Psycho (1998) Crime|Horror|Thriller 2696 4 973308710 7 24210 Lake Placid 991052 2713 2696 1 973308710 7 24210 Horror|Thriller M 25 (1999)Talented Mr. 991053 3176 Ripley, The Drama|Mystery|Thriller 2696 4 973308865 M 25 7 24210 (1999)3386 991054 JFK (1991) Drama|Mystery 2696 1 973308842 M 25 7 24210 ###Feature Engineering: Use column genres: ###Find out all the unique genres ###(Hint: split the data in column genre making a list and then process the data to find out ###only the unique categories of genres) movie_ratings_df['Genres'].value_counts().head() Out[34]: Comedy 116883 Drama 111423 Comedy|Romance 42712 42245 Comedy|Drama Drama|Romance 29170 Name: Genres, dtype: int64 In [35]: movie_ratings_df['Genres'].unique() Out[35]: array(["Animation|Children's|Comedy", "Adventure|Children's|Fantasy", 'Comedy|Romance', 'Comedy|Drama', 'Comedy', 'Action|Crime|Thriller', "Adventure|Children's", 'Action', 'Action|Adventure|Thriller', 'Comedy|Drama|Romance', 'Comedy|Horror', "Animation|Children's", 'Drama', 'Action|Adventure|Romance', 'Drama|Thriller', 'Drama|Romance', 'Thriller', 'Action|Comedy|Drama', 'Crime|Drama|Thriller', 'Drama|Sci-Fi', 'Romance', 'Adventure|Sci-Fi', 'Adventure|Romance', "Children's | Comedy | Drama", 'Documentary', 'Drama | War', 'Action|Crime|Drama', 'Action|Adventure', 'Crime|Thriller', "Animation|Children's|Musical|Romance", "Children's|Comedy", 'Drama|Mystery', 'Sci-Fi|Thriller', 'Action|Comedy|Crime|Horror|Thriller', 'Drama|Musical', 'Crime|Drama|Romance', 'Adventure|Drama', 'Action|Thriller', "Adventure|Children's|Comedy|Musical", 'Action|Drama|War', 'Action|Adventure|Crime', 'Crime', 'Drama|Mystery|Romance', 'Action|Drama', 'Drama|Romance|War', 'Horror', 'Action|Adventure|Comedy|Crime', 'Comedy|War' 'Action|Adventure|Mystery|Sci-Fi', 'Drama|Thriller|War', 'Action|Romance|Thriller', 'Crime|Film-Noir|Mystery|Thriller', 'Action|Adventure|Drama|Romance', "Adventure|Children's|Drama", 'Action|Sci-Fi|Thriller', 'Action|Adventure|Sci-Fi', "Action|Children's", 'Horror|Sci-Fi', 'Action|Crime|Sci-Fi', 'Western', "Animation|Children's|Comedy|Romance", "Children's|Drama", 'Crime|Drama', 'Drama|Fantasy|Romance|Thriller', 'Drama|Horror', 'Comedy|Sci-Fi', 'Mystery|Thriller', "Adventure|Children's|Comedy|Fantasy|Romance", 'Action|Adventure|Fantasy|Sci-Fi', 'Drama|Romance|War|Western', 'Action|Drama|Thriller', 'Crime|Drama|Romance|Thriller', 'Action|Adventure|Western', 'Horror|Thriller', "Children's | Comedy | Fantasy", 'Film-Noir | Thriller', 'Action|Comedy|Musical|Sci-Fi', "Children's", 'Drama|Mystery|Thriller', 'Comedy|Romance|War', 'Action|Comedy', "Adventure|Children's|Romance", "Animation|Children's|Musical" 'Comedy|Crime|Fantasy', 'Action|Comedy|Western', 'Action|Sci-Fi', 'Action|Adventure|Comedy|Romance', 'Comedy|Thriller', 'Horror|Sci-Fi|Thriller', 'Mystery|Romance|Thriller', 'Comedy|Western', 'Drama|Western', 'Action|Adventure|Crime|Thriller', 'Action|Comedy|War', 'Comedy|Mystery', 'Comedy|Mystery|Romance', 'Comedy|Drama|War', 'Action|Drama|Mystery', 'Comedy|Crime|Horror', 'Film-Noir|Sci-Fi', 'Comedy|Romance|Thriller', "Action|Adventure|Children's|Sci-Fi", "Children's|Comedy|Musical", 'Action|Adventure|Comedy', 'Action|Crime|Romance', "Action|Adventure|Animation|Children's|Fantasy", "Animation|Children's|Comedy|Musical", 'Adventure|Drama|Western', 'Action|Adventure|Crime|Drama', 'Action|Adventure|Animation|Horror|Sci-Fi', 'Action|Horror|Sci-Fi', 'War', 'Action|Adventure|Mystery', 'Mystery', 'Action|Adventure|Fantasy', "Adventure|Animation|Children's|Comedy|Fantasy", 'Sci-Fi', 'Documentary|Drama', 'Action|Adventure|Comedy|War', 'Crime|Film-Noir|Thriller', 'Animation', 'Action|Adventure|Romance|Thriller', 'Animation|Sci-Fi', 'Animation|Comedy|Thriller', 'Film-Noir', 'Sci-Fi|War', 'Adventure', 'Comedy|Crime', 'Action|Sci-Fi|War', 'Comedy|Fantasy|Romance|Sci-Fi', 'Fantasy', 'Action|Mystery|Thriller', 'Comedy|Musical', 'Action|Adventure|Sci-Fi|Thriller', "Children's|Drama|Fantasy",
'Adventure|War', 'Musical|Romance', 'Comedy|Musical|Romance',
'Comedy|Mystery|Romance|Thriller', 'Film-Noir|Mystery', 'Musical', "Adventure|Children's|Drama|Musical", 'Drama|Mystery|Sci-Fi|Thriller', 'Romance|Thriller', 'Film-Noir|Romance|Thriller', 'Crime|Film-Noir|Mystery', 'Adventure|Comedy', 'Action|Adventure|Romance|War', 'Romance|War', 'Action|Drama|Western', 'Action|Crime', "Children's|Comedy|Western", "Adventure|Children's|Comedy",
"Children's|Comedy|Mystery", "Adventure|Children's|Fantasy|Sci-Fi", 'Adventure|Animation|Children's|Musical", "Adventure|Children's|Musical", 'Crime|Film-Noir', "Adventure|Children's|Comedy|Fantasy", "Children's|Drama|Fantasy|Sci-Fi", 'Action|Romance', 'Adventure|Western', 'Comedy|Fantasy', 'Animation|Comedy' 'Crime|Drama|Film-Noir', 'Action|Adventure|Drama|Sci-Fi|War', 'Action|Sci-Fi|Thriller|War', 'Action|Western', "Action|Animation|Children's|Sci-Fi|Thriller|War", 'Action|Adventure|Romance|Sci-Fi|War', 'Action|Horror|Sci-Fi|Thriller', 'Action|Adventure|Comedy|Horror|Sci-Fi', 'Action|Comedy|Musical', 'Mystery|Sci-Fi', 'Film-Noir|Mystery|Thriller', 'Adventure|Comedy|Drama', 'Action|Adventure|Comedy|Horror', 'Action|Drama|Mystery|Romance|Thriller', 'Comedy|Mystery|Thriller', 'Adventure|Animation|Sci-Fi|Thriller', 'Action|Drama|Romance', 'Action|Adventure|Drama', 'Comedy|Drama|Musical', 'Documentary|War', 'Drama|Musical|War', 'Action|Horror', 'Horror|Romance', 'Action|Comedy|Sci-Fi|War', 'Crime|Drama|Sci-Fi', 'Action|Romance|War', 'Action|Comedy|Crime|Drama', 'Action|Drama|Thriller|War', "Action|Adventure|Children's", "Action|Adventure|Children's|Fantasy", "Adventure|Animation|Children's|Comedy|Musical", 'Action|Adventure|Comedy|Sci-Fi', "Children's|Fantasy", 'Crime|Drama|Mystery', 'Action|Mystery|Sci-Fi|Thriller', 'Action|Mystery|Romance|Thriller', 'Adventure|Thriller', 'Action|Thriller|War', 'Action|Crime|Mystery', 'Horror|Mystery|Thriller', 'Crime|Horror|Mystery|Thriller', 'Comedy|Drama|Thriller', 'Drama|Sci-Fi|Thriller', 'Drama|Romance|Thriller', 'Action|Adventure|Sci-Fi|War', 'Comedy|Crime|Drama|Mystery', 'Comedy|Crime|Mystery|Thriller', 'Film-Noir|Sci-Fi|Thriller', 'Adventure|Sci-Fi|Thriller', 'Crime|Drama|Mystery|Thriller', 'Comedy|Crime|Drama', 'Comedy|Documentary', 'Documentary|Musical', 'Action|Drama|Sci-Fi|Thriller', "Adventure | Animation | Children's | Fantasy", 'Adventure|Comedy|Romance', 'Mystery|Sci-Fi|Thriller', 'Action|Comedy|Crime', "Animation|Children's|Fantasy|War", 'Action|Crime|Drama|Thriller', 'Comedy|Sci-Fi|Western',
"Children's|Fantasy|Musical", 'Fantasy|Sci-Fi',
"Children's|Comedy|Sci-Fi", "Action|Adventure|Children's|Comedy", "Adventure|Children's|Drama|Romance", "Adventure|Children's|Sci-Fi", "Adventure|Children's|Comedy|Fantasy|Sci-Fi", "Animation|Children's|Comedy|Musical|Romance", "Children's | Musical", 'Drama | Fantasy', "Animation|Children's|Fantasy|Musical", 'Adventure|Comedy|Musical', "Children's|Sci-Fi", "Children's|Horror", 'Comedy|Fantasy|Romance', 'Comedy|Crime|Thriller', "Adventure|Animation|Children's|Sci-Fi", 'Action|Crime|Mystery|Thriller', 'Adventure|Musical', "Animation|Children's|Drama|Fantasy", "Children's|Fantasy|Sci-Fi", 'Adventure|Fantasy|Romance', 'Crime|Horror', 'Action|Adventure|Horror', 'Adventure|Fantasy|Sci-Fi', 'Drama|Film-Noir|Thriller', 'Action|Comedy|Fantasy', 'Sci-Fi|Thriller|War', 'Action|Adventure|Sci-Fi|Thriller|War', 'Action|Adventure|Drama|Thriller', 'Crime|Horror|Thriller', 'Animation|Musical', 'Action|War',
'Action|Comedy|Romance|Thriller', 'Comedy|Horror|Thriller', 'Drama|Horror|Thriller', 'Action|Sci-Fi|Thriller|Western', 'Drama|Romance|Sci-Fi', 'Action|Adventure|Horror|Thriller', 'Comedy|Film-Noir|Thriller', 'Comedy|Horror|Musical|Sci-Fi', 'Comedy|Romance|Sci-Fi', 'Action|Comedy|Sci-Fi|Thriller', 'Action|Sci-Fi|Western', 'Comedy|Horror|Musical', 'Crime|Mystery', 'Animation|Mystery', 'Action|Horror|Thriller', 'Action|Drama|Fantasy|Romance', 'Horror|Mystery', "Adventure|Animation|Children's", 'Musical|Romance|War', 'Adventure|Drama|Romance', 'Adventure|Animation|Film-Noir', 'Action|Adventure|Animation', 'Comedy|Drama|Western', 'Adventure | Comedy | Sci-Fi', 'Drama | Romance | Western', 'Comedy|Drama|Sci-Fi', 'Action|Drama|Romance|Thriller', 'Adventure|Romance|Sci-Fi', 'Film-Noir|Horror', 'Crime|Drama|Film-Noir|Thriller', 'Action|Adventure|War', 'Romance|Western', "Action|Children's|Fantasy" 'Adventure|Drama|Thriller', 'Adventure|Fantasy', 'Musical|War', 'Adventure|Musical|Romance', 'Action|Romance|Sci-Fi', 'Drama|Film-Noir', 'Comedy|Horror|Sci-Fi', 'Adventure|Drama|Romance|Sci-Fi', 'Adventure|Animation|Sci-Fi', 'Adventure|Crime|Sci-Fi|Thriller'], dtype=object) In [36]: ###Create a separate column for each genre category with a one-hot encoding (1 and 0) ###whether or not the movie belongs to that genre. movie_ratings_selected_df = movie_ratings_users_df[['Gender', 'Age', 'Occupation', 'Rating', 'Genres']] In [48]: Genre = movie_ratings_selected_df['Genres'] Genre = Genre.str.get_dummies().add_prefix('Genres_') movie_ratings_genres_df = pd.concat([movie_ratings_selected_df.drop(['Genres'], axis=1), Genre], axis=1) In [76]: movie_ratings_genres_df =movie_ratings_genres_df.reindex(movie_ratings_selected_df.index) movie_ratings_genres_df.head() Out[76]: Age Occupation Rating Genres_Action Genres_Adventure Genres_Animation Genres_Children's Genres_Comedy Genre 0 1 10 5 0 1 1 1 1 10 5 0 0 1 1 0 0 10 5 1 1 10 4 1 0 0 0 1 10 5 0 5 rows × 23 columns In [60]: ###Determine the features affecting the ratings of any particular movie. movie_ratings_genres_df.dtypes Out[60]: Age int64 Occupation int64 Rating int64 Genres_Action int64 Genres_Adventure int64 Genres_Animation int64 Genres_Children's int64 Genres_Comedy int64 Genres_Crime int64 Genres_Documentary int64 Genres_Drama int64 Genres_Fantasy int64 Genres_Film-Noir int64 Genres_Horror int64 Genres_Musical int64 int64 Genres_Mystery Genres_Romance int64 Genres_Sci-Fi int64 Genres_Thriller int64 int64 Genres_War Genres_Western int64 Gender_F uint8 uint8 Gender_M dtype: object In [61]: from sklearn.linear_model import LinearRegression from sklearn.model_selection import train_test_split from sklearn import metrics lineReg = LinearRegression(copy_X=True, fit_intercept=True, n_jobs=1, normalize=**False** In [62]: movie_ratings_users_sample_df = movie_ratings_genres_df.sample(n=50000, random_state=0 movie_ratings_users_sample_df.head() Out[62]: Age Occupation Rating Genres_Action Genres_Adventure Genres_Animation Genres_Children's Genres_Comedy 324271 18 0 0 0 0 1 0 818637 18 4 3 0 0 1 1 148677 18 0 0 14 7 0 0 0 778790 50 4 1 0 525489 25 0 5 rows × 23 columns In [64]: x = movie_ratings_users_sample_df.drop('Rating', axis=1) y = movie_ratings_users_sample_df['Rating'] x.shape Out[64]: (50000, 22) In [65]: x_train, x_test, y_train, y_test = train_test_split(Х, У, test_size=0.20, random_state=0 In [68]: linear_reg = LinearRegression() In [69]: linear_reg.fit(x_train, y_train) Out[69]: LinearRegression() In [70]: y_pred = linear_reg.predict(x_test) In [71]: print('y-intercept: ', linear_reg.intercept_ print('Beta coefficients: ', linear_reg.coef_ print('Mean Abs Error MAE: ', metrics.mean_absolute_error(y_test, y_pred) print('Mean Sq Error MSE: ', metrics.mean_squared_error(y_test, y_pred) print('Root Mean Sq Error RMSE:', np.sqrt(metrics.mean_squared_error(y_test, y_pred)) print('r2 value: ', metrics.r2_score(y_test, y_pred) y-intercept: 3.3714137555159684 Beta coefficients: [0.00406322 0.00098825 -0.0933231 0.00822898 0.41190314 -0.32536968 -0.00937548 0.07845926 0.43311855 0.22781148 0.07368389 0.3951835 $-0.29085584 \quad 0.12523149 \quad 0.02288591 \quad 0.00234758 \quad -0.01347635 \quad 0.06128953$ 0.30880281 0.14777492 0.01440465 -0.01440465] Mean Abs Error MAE: 0.8978299534841195 Mean Sq Error MSE: 1.1977731707567232 Root Mean Sq Error RMSE: 1.0944282391992282 r2 value: 0.03795269985311833 In [75]: ###Age & Occupation are the main features that affect the movie ratings ###Develop an appropriate model to predict the movie ratings In [73]: x_train.dtypes Out[73]: Age int64 **Occupation** int64 Genres_Action int64 Genres_Adventure int64 Genres_Animation int64 Genres_Children's int64 Genres_Comedy int64 Genres_Crime int64 Genres_Documentary int64 Genres_Drama int64 Genres_Fantasy int64 Genres_Film-Noir int64 Genres_Horror int64 Genres_Musical int64 Genres_Mystery int64 Genres_Romance int64 Genres_Sci-Fi int64 Genres_Thriller int64 Genres_War int64 Genres_Western int64 uint8 Gender_F Gender_M uint8 dtype: object In [74]: prediction_df = pd.DataFrame({'Test': y_test, 'Prediction': y_pred}) prediction_df.head() Out[74]: **Test Prediction** 187446 4 4.322363 69421 4 3.439548 941725 3 3.408593 841836 4 3.652663 869012 4 3.559433

In [56]: import numpy as np

import pandas as pd

import seaborn as sns
%matplotlib inline

engine='python',
header=None

users_df = pd.read_csv(

names=['MovieID','Title','Genres'],

from pandas import Series
from pandas import DataFrame
import matplotlib.pyplot as plt
from matplotlib import style