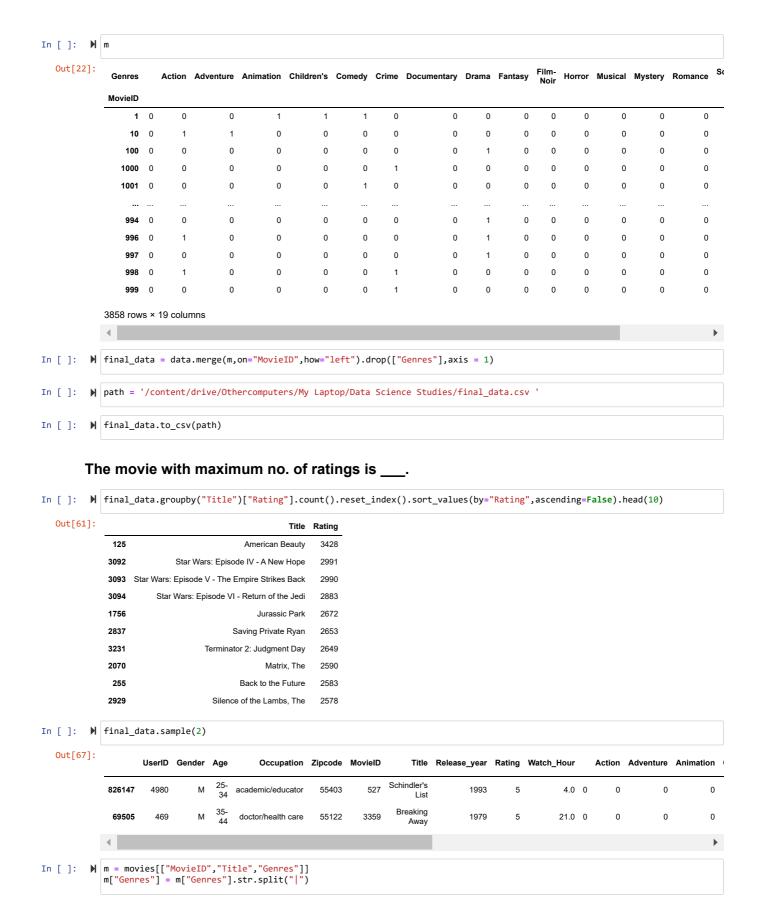
## optional work

```
In [ ]: ▶
            import numpy as np
            import pandas as pd
            import matplotlib.pyplot as plt
            plt.rcParams["figure.figsize"] = (12,8)
            import warnings
            warnings.filterwarnings("ignore")
            pd.set_option('display.max_rows', 500)
            pd.set_option('display.max_columns', 500)
            pd.set_option('display.width', 10000)
            from google.colab import drive
            drive.mount('/content/drive')
            Mounted at /content/drive
In []: M movies = pd.read_fwf("/content/drive/Othercomputers/My Laptop/Data Science Studies/GitHub_Desktop/BusinessCase_Data_Explo
            ratings =pd.read_fwf("/content/drive/Othercomputers/My Laptop/Data Science Studies/GitHub_Desktop/BusinessCase_Data_Explo
            users = pd.read_fwf("/content/drive/Othercomputers/My Laptop/Data Science Studies/GitHub_Desktop/BusinessCase_Data_Explor
In [ ]: ▶ | delimiter ="::"
            users = users["UserID::Gender::Age::Occupation::Zip-code"].str.split(delimiter,expand = True)
            users.columns = ["UserID", "Gender", "Age", "Occupation", "Zipcode"]
            users["Occupation"] = users["Occupation"].astype(int).replace({0: "other",1: "academic/educator",2: "artist",
                                                                          3: "clerical/admin",4: "college/grad student",
                                                        5: "customer service",6: "doctor/health care",7: "executive/managerial",8: "farmer",9: "homemaker",10: "K-12 student",11: "lawyer",
                                                        12: "programmer",13: "retired",14: "sales/marketing",15: "scientist",
                                                        16: "self-employed",17: "technician/engineer"
                                                        18: "tradesman/craftsman",19: "unemployed",20: "writer"},
            delimiter ="::"
            ratings = ratings["UserID::MovieID::Rating::Timestamp"].str.split(delimiter,expand = True)
            ratings.columns = ["UserID", "MovieID", "Rating", "Timestamp"]
            movies.drop(["Unnamed: 1","Unnamed: 2"],axis = 1,inplace=True)
            delimiter ="::"
            movies = movies["Movie ID::Title::Genres"].str.split(delimiter,expand = True)
            movies.columns = ["MovieID","Title","Genres"]
            movies.shape,ratings.shape,users.shape
  Out[11]: ((3883, 3), (1000209, 4), (6040, 5))
In [ ]: M movies["Release_year"] = movies["Title"].str.extract('^(.+)\s\(([0-9]*)\)$',expand = True)[1]
            movies["Title"] = movies["Title"].str.split("(").apply(lambda x:x[0])
In [ ]: ▶ from datetime import datetime
            ratings["Watch\_Hour"] = ratings["Timestamp"].apply(lambda x:datetime.fromtimestamp(int(x)).hour)
            ratings.drop(["Timestamp"],axis = 1,inplace=True)
In [ ]:  ▶ | movies.shape,ratings.shape,users.shape
  Out[14]: ((3883, 4), (1000209, 4), (6040, 5))
In [ ]: M df = users.merge(movies.merge(ratings,on="MovieID",how="outer"),on="UserID",how="outer")
```

```
In [ ]: | (df.isna().sum())/len(df) * 100
   Out[9]: UserID
                              0.017693
             Gender
                              0.017693
                              0.017693
             Age
                              0.017693
             Occupation
             Zipcode
                              0.017693
             MovieTD
                              9.999999
             Title
                              0.000000
             Genres
                              0.406443
             Release_year
                              0.377854
             Rating
                              0.017693
             Watch_Hour
                              0.017693
             dtype: float64
In []: M data = df.copy()
             data.dropna(inplace= True)
In [ ]: ▶ data
  Out[11]:
                      UserID Gender
                                            Occupation Zipcode MovielD
                                                                             Title
                                                                                                          Genres Release_year Rating Watch_Hour
                                      Age
                   0
                                            K-12 student
                                                         48067
                                                                         Toy Story
                                                                                         Animation|Children's|Comedy
                                                                                                                        1995
                                                                                                                                  5
                                                                                                                                            23.0
                                        18
                                    Under
                   1
                                  F
                                           K-12 student
                                                         48067
                                                                       Pocahontas Animation|Children's|Musical|Romance
                                                                                                                        1995
                                                                                                                                  5
                                                                                                                                            23.0
                                                                    48
                                        18
                                    Under
                                  F
                   2
                           1
                                           K-12 student
                                                                   150
                                                                                                                                  5
                                                                                                                                            22.0
                                                         48067
                                                                         Apollo 13
                                                                                                          Drama
                                                                                                                        1995
                                        18
                                                                         Star Wars:
                                     Under
                                                                         Episode IV
                                           K-12 student
                                                         48067
                                                                   260
                                                                                             Action|Adventure|Fantas
                                                                                                                         1977
                                                                                                                                            22.0
                                                                           - A New
                                        18
                                                                            Hope
                                    Under
                                                                        Schindler's
                                  F
                   4
                                           K-12 student
                                                         48067
                                                                   527
                                                                                                       Drama|War
                                                                                                                         1993
                                                                                                                                  5
                                                                                                                                            23.0
                                           doctor/health
                                                                            Blood
              1000204
                        6040
                                  M 25-34
                                                         11106
                                                                  3683
                                                                                                   DramalFilm-Noir
                                                                                                                         1984
                                                                                                                                            8.0
                                                                                                                                  4
                                                                           Simple
                                           doctor/health
              1000205
                        6040
                                  M 25-34
                                                         11106
                                                                  3703
                                                                        Mad Max 2
                                                                                                      Action|Sci-Fi
                                                                                                                         1981
                                                                                                                                            23.0
                                           doctor/health
              1000206
                        6040
                                  M 25-34
                                                         11106
                                                                  3735
                                                                           Serpico
                                                                                                      Crime|Drama
                                                                                                                         1973
                                                                                                                                  4
                                                                                                                                            8.0
                                           doctor/health
                                                                          Chicken
              1000207
                        6040
                                     25-34
                                                         11106
                                                                  3751
                                                                                         Animation|Children's|Comedy
                                                                                                                        2000
                                                                                                                                  4
                                                                                                                                            23.0
                                                                             Run
                                                  care
                                  M 25-34 doctor/health
              1000208
                        6040
                                                         11106
                                                                  3819
                                                                          Tampopo
                                                                                                         Comedy
                                                                                                                         1986
                                                                                                                                  5
                                                                                                                                            23.0
             996144 rows × 11 columns
<class 'pandas.core.frame.DataFrame'>
             Int64Index: 996144 entries, 0 to 1000208
             Data columns (total 11 columns):
                                 Non-Null Count
              #
                  Column
                                                    Dtype
              0
                                  996144 non-null
                  UserID
                                                    object
                                  996144 non-null
                  Gender
              1
                                                    object
                                  996144 non-null
              2
                  Age
                                                    object
                  Occupation
              3
                                  996144 non-null
                                                    object
                                  996144 non-null
              4
                  Zipcode
                                                    obiect
              5
                  MovieID
                                  996144 non-null
                                                    object
              6
                  Title
                                  996144 non-null
                                                    object
              7
                  Genres
                                  996144 non-null
                                                    object
              8
                  Release_year
                                  996144 non-null
                                                    object
                  Rating
                                  996144 non-null
                                                    object
              10 Watch_Hour
                                 996144 non-null
                                                    float64
             dtypes: float64(1), object(10)
             memory usage: 91.2+ MB
Out[13]: UserID
                              6040
             Gender
                                 2
             Age
                                  7
             Occupation
                                21
             Zipcode
                              3439
             MovieID
                              3682
             Title
                               3633
             Genres
                                81
             Release_year
             Rating
             Watch_Hour
                                 24
             dtype: int64
```

```
In [ ]: ▶ # 6040 unique UserID
         # 7 different age groups
         # 21 occupations
         # 3493 different locations of users
         # 3682 unique movies
\# Thats is the reason we have lots of NaN values in our final dataset.
In [ ]: ► data.shape
 Out[16]: (996144, 11)
In [ ]:  M m["Genres"] = m["Genres"].str.split("|")
In [ ]: ▶
         m = m.explode("Genres")
         })
In [ ]: ► M m
  Out[20]:
             MovielD
                            Title
                                  Genres
           0
                  1
                          Toy Story Animation
            0
                          Toy Story Children's
            0
                  1
                          Toy Story
                                 Comedy
                 2
            1
                           Jumanii Adventure
            1
                 2
                           Jumanji Children's
          3879
               3949 Requiem for a Dream
                                  Drama
          3880
                3950
                          Tigerland
                3951
                     Two Family House
          3881
                                  Drama
          3882
                3952
                       Contender, The
                                  Drama
          3882
                3952
                       Contender, The
                                  Thriller
         6366 rows × 3 columns
In [ ]: ▶
         m = pd.crosstab(m["MovieID"],
                                     m["Genres"])
         m = pd.DataFrame(np.where(m>=1,1,0),index = m.index,columns=m.columns)
```



```
"Documenta": "Documentary",
                                "Wester": "Western",
                                "Fant": "Fantasy", "Chil": "Children's", "R": "Romance", "D": "Drama", "Rom": "Romance",
                                "Animati": "Animation",
                                "Childr": "Children's", "Childre": "Children's",
                               "Fantasy", "Come": "Comedy", "Dram": "Drama", "S": "Sci-Fi", "Roma": "Romance", "A": "Adventure", "Children": "Children's",
                                "Adventu": "Adventure",
                                 "Adv":"Adventure",
                                "Wa":"War",
"Thrille" :"Thriller"
                                com" :"Thril
Com" :"Comedy"
"Comed" :"C
                                          :"Comedy",
                                 "Acti" :"Action",

"Advent" :"Adve
                                  Adventur" :"Adventure",
                                                 :"Adventure",
                                  "Thri": "Thriller",
"Chi": "Children's",
                                  "Ro": "Romance",
                                  "F":"Fantasy",
                                  "We":"Western",
                                   "Documen": "Documentary"
                                   "Music": "Musical"
                                   "Children":"Children's",
                                "Horr":"Horror"
"Children'":"Children's",
                                "Roman": "Romance", "Docu": "Documentary", "Th": "Thriller", "Document": "Documentary"
                               })
Out[139]:
                       Title UserID
                Genres
                          8
                             1344
                 Action
                        497
                             6011
                             5894
              Adventure
              Animation
                        104
                             4794
                             5280
              Children's
                        244
                Comedy 1182
                             6031
                 Crime
                        209
                             5662
                             2237
                        124
            Documentary
                             6037
                 Drama 1569
                Fantasy
                             4617
               Film-Noir
                             4150
                 Horror
                        332
                             5299
                        113
                             4752
                Musical
                Mystery
               Romance
                        461
                             5910
                 Sci-Fi
                        261
                             5809
                 Thriller
                        481
                             5989
                   War
                        139
                             5649
                Western
                         68
                             4100
```

# Number of Movie Titles been raten as per Genre by each type of use Occupation :

```
In []: N Occupation_genre_count = merged_data.groupby(["Occupation","Genres"])["Title"].nunique().sort_values(ascending=False).res
```

```
In []: M Occupation_genre_count.pivot(index="Occupation",columns="Genres",values="Title")
  Out[135]:
                                                                                                                                Film-
Marie Horror Musical Mystery F
                                       Action Adventure Animation Children's Comedy Crime Documentary Drama Fantasy
                            Genres
                        Occupation
                       K-12 student 5
                                                     234
                                                                            229
                                                                                     855
                                                                                             135
                                                                                                           35
                                                                                                                  852
                                                                                                                                   30
                                                                                                                                         228
                                                                                                                                                             74
                  academic/educator 6
                                          458
                                                     262
                                                                  96
                                                                            231
                                                                                    1033
                                                                                             175
                                                                                                           74
                                                                                                                 1239
                                                                                                                            62
                                                                                                                                   40
                                                                                                                                         282
                                                                                                                                                   104
                                                                                                                                                             94
                                          457
                                                     260
                                                                  99
                                                                            224
                                                                                            176
                                                                                                           69
                                                                                                                 1203
                                                                                                                            60
                                                                                                                                   39
                                                                                                                                         284
                              artist 6
                                                                                    1002
                                                                                                                                                   107
                                                                                                                                                             93
                      clerical/admin 6
                                          440
                                                     245
                                                                  98
                                                                            212
                                                                                     969
                                                                                            169
                                                                                                           56
                                                                                                                 1122
                                                                                                                            61
                                                                                                                                   36
                                                                                                                                         246
                                                                                                                                                   104
                                                                                                                                                             96
                college/grad student 7
                                          464
                                                     266
                                                                 103
                                                                            237
                                                                                    1072
                                                                                            185
                                                                                                           83
                                                                                                                 1294
                                                                                                                            61
                                                                                                                                   41
                                                                                                                                         297
                                                                                                                                                   107
                                                                                                                                                             98
                   customer service 6
                                                                            202
                                                                                     883
                                                                                            160
                                                                                                            40
                                                                                                                            59
                                                                                                                                                   98
                  doctor/health care 7
                                          450
                                                     253
                                                                 102
                                                                            228
                                                                                     946
                                                                                            165
                                                                                                           68
                                                                                                                 1137
                                                                                                                            60
                                                                                                                                   41
                                                                                                                                         285
                                                                                                                                                   111
                                                                                                                                                             89
                executive/managerial 6
                                          466
                                                                 101
                                                                            234
                                                                                    1038
                                                                                                           70
                                                                                                                            62
                                                                                                                                         301
                                                                                                                                                   106
                                                     269
                                                                                            180
                                                                                                                 1243
                                                                                                                                   40
                                                                                                                                                             94
                                                                  81
                                                                            151
                                                                                     519
                                                                                                            8
                                                                                                                            44
                                                                                                                                                   57
                        homemaker 5
                                          358
                                                     211
                                                                  87
                                                                            211
                                                                                     808
                                                                                            127
                                                                                                           21
                                                                                                                  789
                                                                                                                            55
                                                                                                                                   26
                                                                                                                                          159
                                                                                                                                                   90
                                                                                                                                                             79
                                          425
                                                                  91
                                                                                                                  933
                                                                                                                            55
                                                                                                                                                   96
                             lawyer 5
                                                     238
                                                                            200
                                                                                     853
                                                                                            154
                                                                                                           44
                                                                                                                                   37
                                                                                                                                         224
                                                                                                                                                             87
                              other
                                          472
                                                     273
                                                                 103
                                                                            240
                                                                                     1081
                                                                                            195
                                                                                                            86
                                                                                                                 1330
                                                                                                                            63
                                                                                                                                   43
                                                                                                                                         317
                                                                                                                                                   108
                                                                                                                                                             97
                                          455
                                                     259
                                                                            219
                                                                                     943
                                                                                                           59
                                                                                                                 1089
                                                                                                                            60
                                                                                                                                   36
                                                                                                                                                   98
                       programmer 6
                                                                 100
                                                                                            163
                                                                                                                                         264
                                                                                                                                                             94
                             retired 3
                                          366
                                                     213
                                                                  63
                                                                            152
                                                                                     750
                                                                                             139
                                                                                                           41
                                                                                                                  925
                                                                                                                            45
                                                                                                                                   34
                                                                                                                                          177
                                                                                                                                                   89
                                                                                                                                                             86
                    sales/marketing 6
                                          459
                                                     251
                                                                  97
                                                                            217
                                                                                     970
                                                                                             175
                                                                                                           59
                                                                                                                 1100
                                                                                                                            59
                                                                                                                                   36
                                                                                                                                         251
                                                                                                                                                   102
                                                                                                                                                             93
                           scientist 6
                                          424
                                                     230
                                                                  86
                                                                            176
                                                                                     804
                                                                                            142
                                                                                                           42
                                                                                                                  971
                                                                                                                            54
                                                                                                                                   31
                                                                                                                                          194
                                                                                                                                                   86
                                                                                                                                                            87
                      self-employed 7
                                                     258
                                                                                     994
                                                                                            179
                                                                                                                 1182
                                                                                                                            57
                                                                                                                                         281
                                                                                                                                                   99
                 technician/engineer 6
                                          466
                                                     262
                                                                 101
                                                                            233
                                                                                     990
                                                                                            172
                                                                                                           64
                                                                                                                 1146
                                                                                                                            62
                                                                                                                                   38
                                                                                                                                         295
                                                                                                                                                   104
                                                                                                                                                             96
                                                                  76
                                                                                                                            57
                                                                                                                                                   87
                                                                                                                                                             83
               tradesman/craftsman 4
                                          414
                                                     226
                                                                            182
                                                                                     789
                                                                                            155
                                                                                                           36
                                                                                                                  853
                                                                                                                                   33
                                                                                                                                         279
                       unemployed 6
                                          432
                                                     235
                                                                  93
                                                                                             148
                                                                                                                                                   95
                                                                                                                                                             83
                                                                            208
                                                                                                                            59
                                                                                                                                          266
                             writer 8
                                          461
                                                     263
                                                                 102
                                                                            233
                                                                                    1053
                                                                                            182
                                                                                                           80
                                                                                                                 1294
                                                                                                                            62
                                                                                                                                   42
                                                                                                                                         295
                                                                                                                                                   107
                                                                                                                                                             99
In []: M pd.set_option("max_colwidth", None)
```

# **Co-occurance | Frequency Based Recommender System (Apriory)**

```
In []: M rules.groupby(["antecedents"])["lift"].max().reset_index().merge(rules,on=["antecedents","lift"])
                     (Star vvars: Episode VI - Return or
                 20
                              the Jedi . Men in Black .
                                                     2.047691
                                                                              (Matrix, The.)
                                                                                               0.229470
                                                                                                             0.428808 0.201490
                                                                                                                                   0.878066 0.103091
                                                                                                                                                          4.684451
                          Terminator 2: Judgment Day )
                       (Raiders of the Lost Ark . Back to
                                                                   (Star Wars: Episode IV - A
                     the Future , Star Wars: Episode V 1.832652
- The Empire Strikes Back )
                                                                                               0.230960
                                                                                                             0.495199 0.209603
                                                                                                                                   0.907527 0.095231
                                                                                                                                                          5.458898
                                                                               New Hope )
                       (Raiders of the Lost Ark , Back to
                                                                 (Star Wars: Episode V - The
                 22 the Future , Star Wars: Episode IV 1.879063
- A New Hope )
                                                                                               0.225331
                                                                                                             0.495033 0.209603
                                                                                                                                   0.930198 0.098056
                                                                                                                                                          7.234315
                        (Total Recall , Terminator, The ) 2.090372 (Terminator 2: Judgment Day
                                                                                               0.228808
                                                                                                             0.438576 0.209768
                                                                                                                                   0.916787 0.109419
                                                                                                                                                          6.746850
                 24
                                             (Babe) 1.260843
                                                                                               0.289901
                                                                                                             0.567550 0.207450
                                                                                                                                   0.715591 0.042917
                                                                                                                                                          1.520523
                                                                         (American Beauty)
                                                                   (Star Wars: Episode IV - A
                 25
                                      (Aliens , Alien ) 1.826383
                                                                                               0.232119
                                                                                                             0.495199 0.209934
                                                                                                                                   0.904422 0.094989
                                                                                                                                                          5.281578
                                                                               New Hope )
                     (Star Wars: Episode VI - Return of the Jedi , Aliens , Star Wars: 1.941061
                                                                 (Star Wars: Episode V - The
                                                                                               0.215894
                                                                                                             0.495033 0.207450
                                                                                                                                   0.960890 0.100576 12.911310
                                                                      Empire Strikes Back )
                            Episode IV - A New Hope )
                    (Aliens , Star Wars: Episode IV - A
                                                                  (Star Wars: Episode VI -
In []: M rules[rules["antecedents"] == rules.loc[4606]["antecedents"]].sort_values(by="lift",ascending=False).head(5)
   Out[44]:
                      antecedents
                                         consequents antecedent support consequent support support confidence
                                                                                                                            lift leverage conviction
                       (Fight Club ) (American Beauty )
                                                                 0.240232
                                                                                       0.56755
                                                                                                     0.2
                                                                                                           0.832529 1.466884 0.063657
```

#### **Item-Item Similarity Based Rec System**

## Name the top 3 movies similar to 'Liar Liar' on the item-based approach.

```
In [ ]:  M movies[movies["Title"].str.contains("Liar Liar")]
   Out[46]:
                              Title Genres Release year
              1455
                      1485 Liar Liar Comedy
In [ ]: | m = movies[["MovieID","Title","Genres"]]
             m = m.explode("Genres")
             "Dr":"Drama",
"Documenta":"Documentary",
                                     "Wester": "Western",
                                     "Fant": "Fantasy", "Chil": "Children's", "R": "Romance", "D": "Drama", "Rom": "Romance",
                                    "Animati": "Animation",
"Childre": "Children's", "Childre": "Children's",
                                     "Fantas": "Fantasy", "Come": "Comedy", "Dram": "Drama", "S": "Sci-Fi", "Roma": "Romance", "A": "Adventure", "Children": "Children's",
                                     "Adventu":"Adventure",
                                      "Adv":"Adventure",
                                      "Wa":"War",
"Thrille" :"Thriller"
                                      "Com"
                                                   :"Comedy"
                                      "Comed"
                                                      :"Comedy",
                                      'Acti"
                                                       :"Action"
                                        "Advent"
                                                       :"Adventure"
                                        "Adventur"
                                                         :"Adventure",
                                       "Thri":"Thriller",
"Chi":"Children's",
                                       "Ro": "Romance",
"F": "Fantasy",
                                        "We": "Western",
                                        "Documen": "Documentary"
                                        "Music": "Musical"
                                        "Children":"Children's"
                                     "Horr":"Horror"
"Children'":"Children's",
                                     "Roman": "Romance", "Docu": "Documentary", "Th": "Thriller", "Document": "Documentary"
                                    })
             m = pd.crosstab(m["MovieID"],m["Genres"])
               = pd.DataFrame(np.where(m>=1,1,0),index = m.index,columns=m.columns)
```

```
return np.sum(abs(x1-x2))
```

```
In [ ]: N Ranks = []
               Query = "1485"
               for candidate in m.index:
                 if candidate == Query:
                  Ranks.append([Query,candidate,Hamming_distance(m.loc[Query],m.loc[candidate])])
In [ ]:  M Ranks = pd.DataFrame(Ranks,columns=["Query","Candidate","Hamming_distance"])
               Ranks = Ranks.merge(movies[['MovieID', 'Title']], left_on='Query', right_on='MovieID').rename(columns={'Title': 'query_ti
Ranks = Ranks.merge(movies[['MovieID', 'Title']], left_on='Candidate', right_on='MovieID').rename(columns={'Title': 'candidate', right_on='MovieID').rename(columns={'Title': 'candidate', right_on='MovieID').rename(columns={'Title': 'candidate', right_on='MovieID').rename(columns={'Title': 'candidate', right_on='MovieID').rename(columns={'Title': 'candidate', right_on='MovieID').rename(columns={'Title': 'query_ti
               Ranks = Ranks.sort_values(by=['Query', 'Hamming_distance'])
In [ ]: ▶
               Ranks.head()
   Out[52]:
                    Query Candidate Hamming_distance query_tittle candidate_tittle
                 4
                    1485
                                 1001
                                                               Liar Liar
                                                                         Associate, The
                     1485
                                 1002
                 5
                                                        0
                                                              Liar Liar Ed's Next Move
                     1485
                                  101
                                                                           Bottle Rocket
                24
                     1485
                                  102
                                                        0
                                                               Liar Liar
                                                                              Mr. Wrong
                25
                     1485
                                 1020
                                                        0
                                                               Liar Liar Cool Runnings
           Collaborative Filtering:
In [ ]: | user_movie_ratings = ratings.pivot(index ="UserID",
                                columns = "MovieID"
                                 values ="Rating").fillna(0)
In [ ]: | user_movie_ratings.shape
   Out[54]: (6040, 3706)
In []: ▶ # 6040 users # 3706 movies
In [ ]: ▶ ratings.shape, users.shape
   Out[56]: ((1000209, 4), (6040, 5))
In [ ]: ▶
               rm_raw = ratings[['UserID', 'MovieID', 'Rating']].copy()
rm_raw.columns = ['UserId', 'ItemId', 'Rating'] # Lib requires specific column names
In [ ]: | rm_raw.Rating = rm_raw.Rating.astype(int)
               rm_raw.UserId = rm_raw.UserId.astype(int)
               rm_raw.ItemId = rm_raw.ItemId.astype(int)
In [ ]: | rm_raw.nunique()
 Out[133]: UserId
                           6040
                           3706
               ItemId
               Rating
               dtype: int64
In [ ]: ▶ # !pip install cmfrec
In []: ▶ from cmfrec import CMF
               model = CMF(k=5, lambda_=0.1, user_bias=False, item_bias=False, verbose=False)
               model.fit(rm_raw)
 Out[111]: Collective matrix factorization model
               (explicit-feedback variant)
Out[112]: ((1000209, 3), (6040, 5), (3706, 5))
```

```
In [ ]:  M model.A_.shape,model.B_.T.shape
 Out[113]: ((6040, 5), (5, 3706))
In [ ]: ▶ model.topN(user=8, n=10)
 Out[126]: array([2323, 296, 1136, 1089, 3421, 858, 1198, 3552, 260, 2776])
            movies_to_recommend = model.topN(user=1, n=10)
In [ ]: ▶
             movies_to_recommend = movies_to_recommend[movies_to_recommend<3706]</pre>
             movies_to_recommend
             movies.MovieID = movies.MovieID.astype(int)
            movies.loc[movies to recommend]
 Out[143]:
                  MovieID
                                     Title
                                                              Genres Release vear
             1421
                     1446
                                    Kolva
                                                              Comedy
                                                                            1996
             2099
                    2168
                             Dance with Me
                                                        Drama|Romance
                                                                            1998
             2776
                    2845
                                White Boys
                                                               Drama
                                                                            1999
             2211
                    2280
                               Clay Pigeons
                                                               Crime
                                                                            1998
             1022
                    1035 Sound of Music, The
                                                              Musical
                                                                            1965
                                 Peter Pan Animation|Children's|Fantasy|Musical
                    2087
                                                                            1953
             2018
             1951
                    2020 Dangerous Liaisons
                                                        DramalRomance
                                                                            1988
             1028
                    1041
                              Secrets & Lies
                                                               Drama
                                                                            1996
             2035
                    2104
                                     Tex
                                                               Drama
                                                                            1982
                     599
             595
                            Wild Bunch, The
                                                              Western
                                                                            1969
In [ ]: ▶ !pip install scikit-surprise
            Looking in indexes: https://pypi.org/simple, (https://pypi.org/simple,) https://us-python.pkg.dev/colab-wheels/public/s
            imple/ (https://us-python.pkg.dev/colab-wheels/public/simple/)
            Collecting scikit-surprise
              Downloading scikit-surprise-1.1.3.tar.gz (771 kB)
                                                        772.0/772.0 KB 14.2 MB/s eta 0:00:00
              Preparing metadata (setup.py) \dots done
            Requirement already satisfied: joblib>=1.0.0 in /usr/local/lib/python3.8/dist-packages (from scikit-surprise) (1.2.0)
            Requirement already satisfied: numpy>=1.17.3 in /usr/local/lib/python3.8/dist-packages (from scikit-surprise) (1.21.6)
            Requirement already satisfied: scipy>=1.3.2 in /usr/local/lib/python3.8/dist-packages (from scikit-surprise) (1.7.3)
            Building wheels for collected packages: scikit-surprise
              Building wheel for scikit-surprise (setup.py) ... done
              Created wheel for scikit-surprise: filename=scikit_surprise-1.1.3-cp38-cp38-linux_x86_64.whl size=3366465 sha256=a3d8
            dd2a966666053cdeab5b9aa99d54786198cded3c182220c579d2074fab3b
              Stored in directory: /root/.cache/pip/wheels/af/db/86/2c18183a80ba05da35bf0fb7417aac5cddbd93bcb1b92fd3ea
            Successfully built scikit-surprise
            Installing collected packages: scikit-surprise
            Successfully installed scikit-surprise-1.1.3
In []: M final_data= pd.read_csv("/content/drive/Othercomputers/My Laptop/Data Science Studies/GitHub_Desktop/BusinessCase_Data_Ex
            from surprise import KNNWithMeans
            from surprise import Dataset
            from surprise import accuracy
            from surprise.model_selection import train_test_split
            from surprise import Reader
            ## The Reader class is used to parse a file containing ratings.It orders the data in format of (userid,title,rating) and
            reader = Reader(rating_scale=(0.5, 5))
            # The columns must correspond to user id, item id and ratings (in that order).
            data = Dataset.load_from_df(final_data[["UserID", "MovieID", "Rating"]], reader) # loading the data as per the format
                                                                                                                                    Þ
In [ ]: ▶ data
  Out[26]: <surprise.dataset.DatasetAutoFolds at 0x7f814324beb0>
In [ ]: M trainset, testset = train_test_split(data, test_size=.15) # Splitting the data
```

## **User Based Collaborative Filtering:**

print("Item-based Model : Test Set")
accuracy.rmse(test pred, verbose=True)

Item-based Model : Test Set

RMSE: 0.8926
Out[156]: 0.8926413172784623

```
In [ ]: M algo = KNNWithMeans(k = 50, sim_options={'name': 'cosine', 'user_based': True})
           # K value represents the (max) number of neighbors to take into account for aggregation. Example for every item it gives
           # There are many similarity options to calculate the similarity between the neighbors. Here, we have used the cosine simi
           # when user_based = True then it performs user based collaborative filtering
           algo.fit(trainset) #fitting the train dataset
           Computing the cosine similarity matrix...
           Done computing similarity matrix.
 Out[145]: <surprise.prediction algorithms.knns.KNNWithMeans at 0x7f2ed8392df0>
In []: ▶ # run the trained model against the testset
           test_pred = algo.test(testset)
In [ ]: | test_pred[0]
 Out[147]: Prediction(uid='770', iid='2383', r_ui=1.0, est=0.8775986827947615, details={'actual_k': 50, 'was_impossible': False})
In [ ]: ▶ # get RMSE on test set
           print("User-based Model : Test Set")
           accuracy.rmse(test_pred, verbose=True)
           User-based Model : Test Set
           RMSE: 0.9345
 Out[148]: 0.9344753997438834
In [ ]: | accuracy.mae(test_pred, verbose=True)
           MAE: 0.7433
 Out[149]: 0.7432764090468391
In []: ▶ # we can query for specific predictions
           uid = str(1) # raw user id
           iid = str(1) # raw item id
r_ui = None est = 4.88 {'actual_k': 50, 'was_impossible': False}
           user: 1
                          item: 1
In [ ]: # anti_pre = algo.test(anti_set)
           # pred_df = pd.DataFrame(anti_pre).merge(movies , Left_on = ['iid'], right_on = ['MovieID'])
           \# pred_df = pd.DataFrame(pred_df).merge(users , left_on = ['uid'], right_on = ['UserID'])
        Item Based Collaborative Filtering:
In []: M # K value represents the (max) number of neighbors to take into account for aggregation. Example for every item it gives
           # There are many similarity options to calculate the similarity between the neighbors . Here, we have used the cosine sim
           # when user_based = False then it performs item based collaborative filtering
           algo_i = KNNWithMeans(k=30, sim_options={'name': 'cosine', 'user_based': False})
           algo_i.fit(trainset)
           4
           Computing the cosine similarity matrix...
           Done computing similarity matrix.
 Out[153]: <surprise.prediction_algorithms.knns.KNNWithMeans at 0x7f2ed36d0be0>
Out[155]: Prediction(uid='770', iid='2383', r_ui=1.0, est=1.1141077593830107, details={'actual_k': 30, 'was_impossible': False})
In [ ]: ▶ # get RMSE on test set
```

```
In [ ]: m{M} # we can query for specific predictions
            uid = str(196) # raw user id
            iid = str(303) # raw item id
In [ ]: | pred = algo_i.predict(uid, iid, verbose=True)
                                               r ui = None est = 2.90 {'actual k': 29, 'was impossible': False}
            user: 196
                             item: 303
In [ ]: | # final_data[final_data["MovieID"]=="984"]
In [ ]: ▶
            tsr inner id = algo i.trainset.to inner iid("1485") #Considering the movieId 1485 : Liar Liar
            tsr_neighbors = algo_i.get_neighbors(tsr_inner_id, k=5) #Getting the 5 nearest neighbors for movieId 1
In [ ]: ▶
            movies[movies.MovieID.isin([algo.trainset.to raw iid(inner id)
                                    for inner_id in tsr_neighbors])] #Displaying the 5 nearest neighbors to the Liar Liar
 Out[161]:
                  MovielD
                                                 Title
                                                              Genres Release_year
              127
                      129
                                           Pie in the Sky Comedy|Romance
                                                                            1995
              968
                     980
                                       In the Line of Duty 2
                                                                            1987
                                                               Action
             1645
                     1692
                                           Alien Escape
                                                           Horror|Sci-Fi
                                                                            1995
                     2158 Henry: Portrait of a Serial Killer, Part 2
                                                          CrimelHorror
                                                                            1996
             2089
             2500
                     2569
                                           Among Giants DramalRomance
                                                                            1998
         Matrix Factorisation:
In [ ]: ▶ from surprise import SVD
            from surprise.model_selection import cross_validate
In []: N svd = SVD() #Suprise library uses the SVD algorithm to perform the matrix factorisation where as other libraries uses ALS
            cross_validate(svd, data, measures=['rmse','mae'], cv = 5 , return_train_measures=True,verbose=True)
            ##The dataset is divided into train and test and with 5 folds the rmse has been calculated
            Evaluating RMSE, MAE of algorithm SVD on 5 split(s).
                               Fold 1 Fold 2 Fold 3 Fold 4 Fold 5 Mean
                                                                                Std
            RMSE (testset)
                               0.8768 0.8736 0.8753 0.8726 0.8729 0.8742 0.0016
                               0.6885 0.6858 0.6883
                                                       0.6852 0.6862
                                                                       0.6868
                                                                                0.0014
            MAE (testset)
            RMSE (trainset)
                               0.6691 0.6702 0.6712
                                                       0.6692 0.6702
                                                                       9.6799
                                                                                0.0008
            MAE (trainset)
                               0.5298 0.5303 0.5310 0.5297
                                                               0.5302
                                                                       0.5302
                                                                                0.0005
            Fit time
                               13.61
                                       14.77
                                               14.17
                                                       14.88
                                                               14.20
                                                                       14.33
                                                                                0.46
            Test time
                               2.63
                                       1.64
                                               4.00
                                                       3.98
                                                               1.65
                                                                        2.78
                                                                                1.05
   Out[30]: {'test_rmse': array([0.87679055, 0.87363147, 0.87525437, 0.87256188, 0.8728605]),
              'train_rmse': array([0.6691445 , 0.67017315, 0.67120119, 0.6692182 , 0.67018469]),
              'test_mae': array([0.68851652, 0.68582647, 0.68829271, 0.68518788, 0.68615884]),
              'train_mae': array([0.52975006, 0.53025507, 0.53101726, 0.52965381, 0.53016147]),
              'fit_time': (13.613670587539673,
              14.765673160552979,
              14.170423984527588,
              14.876917362213135,
              14.199043035507202)
              'test time': (2.6307239532470703,
              1.639423131942749,
              3.9973740577697754,
              3.9794485569000244
              1,650454044342041)}
In [ ]: ▶
            import pandas as pd
            final_data= pd.read_csv("/content/drive/Othercomputers/My Laptop/Data Science Studies/GitHub_Desktop/BusinessCase_Data_Ex
            from surprise import Dataset
            from surprise import SVD
            from surprise import Reader
            ## The Reader class is used to parse a file containing ratings. It orders the data in format of (userid, title, rating) and
            reader = Reader(rating_scale=(0.5 , 5))
# The columns must correspond to user id, item id and ratings (in that order).
```

data = Dataset.load\_from\_df(final\_data[["UserID","MovieID","Rating"]], reader) # Loading the data as per the format

```
trainset = data.build_full_trainset()
             svd.fit(trainset) ##Fitting the trainset with the help of svd
    Out[6]: <surprise.prediction algorithms.matrix factorization.SVD at 0x7fc169469d30>
In []: N svd.pu.shape , svd.qi.shape #pu gives the embeddings of Users and qi gives the embeddings of Items.
   Out[7]: ((6040, 10), (3682, 10))
In [ ]: ▶ #Storing all the movie titles in items
            items = movies['Title'].unique()
             ##Considering the user '662
             test = [[662, iid, 4] for iid in items]
             ##Finding the user predictions(ratings) for all the movies
             predictions = svd.test(test)
             pred = pd.DataFrame(predictions)
In []: Ma = pred.sort_values(by='est', ascending=False) ##Sorting the values based on the estimated predictions
             a[0:10] ##TOP 10
  Out[16]:
                   uid
                                              iid r_ui
                                                                           details
                                                          est
                0 662
                                         Toy Story
                                                   4 3.385119 {'was_impossible': False}
             2560 662
                              It Came from Hollywood 4 3.385119 {'was_impossible': False}
                               House of Frankenstein 4 3.385119 {'was_impossible': False}
             2548 662
             2549 662
                                      Frankenstein 4 3.385119 {'was_impossible': False}
             2550 662
                                 Son of Frankenstein 4 3.385119 {'was_impossible': False}
             2551 662
                            Ghost of Frankenstein, The 4 3.385119 {'was_impossible': False}
             2552 662 Frankenstein Meets the Wolf Man
                                                  4 3.385119 {'was_impossible': False}
                            Curse of Frankenstein, The 4 3.385119 {'was_impossible': False}
             2553 662
             2554 662
                                     Son of Dracula 4 3.385119 {'was_impossible': False}
                                     Wolf Man, The 4 3.385119 {'was_impossible': False}
             2555 662
predictions_svd = svd.test(testset) #Predicting for the test set
In [ ]: ▶ from surprise import accuracy
In []: N print('SVD - RMSE:', accuracy.rmse(predictions_svd, verbose=False))
print('SVD - MAE:', accuracy.mae(predictions_svd, verbose=False))
```

SVD - RMSE: 0.7034125702508484 SVD - MAE: 0.5440720263794291