```
In [107]:
              ▶ import numpy as np
                  import pandas as pd
                  import matplotlib.pyplot as plt
                  import seaborn as sns
                  import warnings
                  warnings.filterwarnings('ignore')
  In [2]:
                 movies = pd.read csv('tmdb 5000 movies.csv')
                  credits = pd.read csv('tmdb 5000 credits.csv')
  In [3]:
                 movies.head(3)
       Out[3]:
                          budget
                                                                                               id keywords original_language original_title
                                                                              homepage
                                                                                                                                                overview
                                                                                                                                                            popularity
                                       genres
                                                                                                       [{"id":
                                                                                                                                                    In the
                                      [{"id": 28,
                                                                                                       1463.
                                                                                                                                                     22nd
                                       "name":
                                                                                                     "name":
                                                                                                                                                century, a
                   0 237000000
                                                                                                                                                           150.437577
                                      "Action"},
                                                               http://www.avatarmovie.com/
                                                                                           19995
                                                                                                                             en
                                                                                                                                        Avatar
                                                                                                      "culture
                                                                                                                                                paraplegic
                                      {"id": 12,
                                                                                                      clash"},
                                                                                                                                                 Marine is
                                        "nam...
                                                                                                      {"id":...
                                                                                                                                                      di...
                                                                                                                                                  Captain
                                                                                                   [{"id": 270,
                                                                                                                                                Barbossa,
                                                                                                                                  Pirates of the
                                      [{"id": 12,
                                                                                                      "name":
                                                                                                                                                     long
                                       "name":
                                                                                                                                    Caribbean:
                                                  http://disney.go.com/disneypictures/pirates/
                                                                                                    "ocean"},
                   1 300000000
                                                                                              285
                                                                                                                                                 believed
                                                                                                                                                           139.082615
                                  "Adventure"},
                                                                                                                                     At World's
                                                                                                    {"id": 726,
                                                                                                                                                    to be
                                    {"id": 14, "...
                                                                                                                                          End
                                                                                                        "na...
                                                                                                                                                    dead.
                                                                                                                                                     ha...
                                                                                                                                                 A cryptic
                                      [{"id": 28,
                                                                                                   [{"id": 470,
                                                                                                                                                 message
                                                                                                     "name":
                                       "name":
                                                                                                                                                     from
                   2 245000000
                                      "Action"},
                                                http://www.sonypictures.com/movies/spectre/
                                                                                                       "spy"},
                                                                                                                                       Spectre
                                                                                                                                                   Bond's
                                                                                                                                                           107.376788
                                                                                                                             en
                                      {"id": 12,
                                                                                                    {"id": 818,
                                                                                                                                                     past
                                        "nam...
                                                                                                     "name...
                                                                                                                                                sends him
                                                                                                                                                      0...
```

```
In [4]:
   Out[4]:
                movie id
                                                    title
                                                                                        cast
                                                                                                                               crew
                   19995
             0
                                                  Avatar
                                                          [{"cast id": 242, "character": "Jake Sully", "...
                                                                                             [{"credit id": "52fe48009251416c750aca23", "de...
             1
                    285 Pirates of the Caribbean: At World's End [{"cast id": 4, "character": "Captain Jack Spa...
                                                                                             [{"credit id": "52fe4232c3a36847f800b579", "de...
             2
                  206647
                                                         [{"cast id": 1, "character": "James Bond", "cr...
                                                                                             [{"credit id": "54805967c3a36829b5002c41", "de...
             3
                   49026
                                      The Dark Knight Rises
                                                        [{"cast_id": 2, "character": "Bruce Wayne / Ba...
                                                                                              [{"credit id": "52fe4781c3a36847f81398c3", "de...
                                              John Carter
                                                          [{"cast id": 5, "character": "John Carter", "c...
                                                                                             [{"credit id": "52fe479ac3a36847f813eaa3", "de...
                   49529
          credits.iloc[0][2]
In [5]:
    Out[5]: '[{"cast id": 242, "character": "Jake Sully", "credit id": "5602a8a7c3a3685532001c9a", "gender": 2, "id": 6573
            1, "name": "Sam Worthington", "order": 0}, {"cast id": 3, "character": "Neytiri", "credit id": "52fe4800925141
            6c750ac9cb", "gender": 1, "id": 8691, "name": "Zoe Saldana", "order": 1}, {"cast id": 25, "character": "Dr. Gr
            ace Augustine", "credit id": "52fe48009251416c750aca39", "gender": 1, "id": 10205, "name": "Sigourney Weaver",
            "order": 2}, {"cast id": 4, "character": "Col. Quaritch", "credit id": "52fe48009251416c750ac9cf", "gender":
            2, "id": 32747, "name": "Stephen Lang", "order": 3}, {"cast id": 5, "character": "Trudy Chacon", "credit id":
             "52fe48009251416c750ac9d3", "gender": 1, "id": 17647, "name": "Michelle Rodriguez", "order": 4}, {"cast id":
            8, "character": "Selfridge", "credit id": "52fe48009251416c750ac9e1", "gender": 2, "id": 1771, "name": "Giovan
            ni Ribisi", "order": 5}, {"cast id": 7, "character": "Norm Spellman", "credit id": "52fe48009251416c750ac9dd",
             "gender": 2, "id": 59231, "name": "Joel David Moore", "order": 6}, {"cast id": 9, "character": "Moat", "credit
             id": "52fe48009251416c750ac9e5", "gender": 1, "id": 30485, "name": "CCH Pounder", "order": 7}, {"cast id": 1
            1, "character": "Eytukan", "credit id": "52fe48009251416c750ac9ed", "gender": 2, "id": 15853, "name": "Wes Stu
            di", "order": 8}, {"cast id": 10, "character": "Tsu\'Tey", "credit id": "52fe48009251416c750ac9e9", "gender":
            2, "id": 10964, "name": "Laz Alonso", "order": 9}, {"cast id": 12, "character": "Dr. Max Patel", "credit id":
            "52fe48009251416c750ac9f1", "gender": 2, "id": 95697, "name": "Dileep Rao", "order": 10}, {"cast id": 13, "cha
            racter": "Lyle Wainfleet", "credit id": "52fe48009251416c750ac9f5", "gender": 2, "id": 98215, "name": "Matt Ge
            rald", "order": 11}, {"cast id": 32, "character": "Private Fike", "credit id": "52fe48009251416c750aca5b", "ge
            nder": 2, "id": 154153, "name": "Sean Anthony Moran", "order": 12}, {"cast id": 33, "character": "Cryo Vault M
            ed Tech", "credit id": "52fe48009251416c750aca5f", "gender": 2, "id": 397312, "name": "Jason Whyte", "order":
In [6]:
          ▶ movies.shape
    Out[6]: (4803, 20)
```

```
M credits.shape
 In [7]:
     Out[7]: (4803, 4)
 In [8]:
            Movies.columns
     Out[8]: Index(['budget', 'genres', 'homepage', 'id', 'keywords', 'original language',
                        'original title', 'overview', 'popularity', 'production companies',
                        'production countries', 'release date', 'revenue', 'runtime',
                        'spoken languages', 'status', 'tagline', 'title', 'vote_average',
                        'vote count'],
                      dtvpe='object')
            ▶ credits.columns
 In [9]:
     Out[9]: Index(['movie id', 'title', 'cast', 'crew'], dtype='object')
            movies = movies.merge(credits)
In [10]:
               movies.head(2)
    Out[10]:
                      budget
                                                                     homepage
                                                                                   id keywords original_language original_title overview
                                                                                                                                           popularity r
                                   genres
                                                                                           [{"id":
                                                                                                                                   In the
                                 [{"id": 28,
                                                                                           1463,
                                                                                                                                    22nd
                                   "name":
                                                                                         "name":
                                                                                                                                century, a
                0 237000000
                                  "Action"},
                                                      http://www.avatarmovie.com/ 19995
                                                                                                                        Avatar
                                                                                                                                          150.437577
                                                                                                               en
                                                                                         "culture
                                                                                                                                paraplegic
                                  {"id": 12,
                                                                                         clash"},
                                                                                                                                Marine is
                                   "nam...
                                                                                         {"id":...
                                                                                                                                     di...
                                                                                                                                 Captain
                                                                                       [{"id": 270,
                                                                                                                                Barbossa,
                                                                                                                   Pirates of the
                                 [{"id": 12,
                                                                                         "name":
                                                                                                                                    long
                                  "name":
                                                                                                                     Caribbean:
                1 300000000
                                           http://disney.go.com/disneypictures/pirates/
                                                                                        "ocean"},
                                                                                                                                 believed 139.082615
                              "Adventure"},
                                                                                                                     At World's
                                                                                       {"id": 726,
                                                                                                                                    to be
                               {"id": 14, "...
                                                                                                                          End
                                                                                           "na...
                                                                                                                                   dead,
                                                                                                                                    ha...
               2 rows × 23 columns
```

```
    ■ movies.shape

In [11]:
    Out[11]: (4809, 23)
             pd.pandas.set option('display.max columns', None)
In [12]:
In [13]:

    movies.head(2)

    Out[13]:
                                                                                         id keywords original_language original_title
                        budget
                                                                                                                                        overview
                                                                                                                                                     popularity r
                                      genres
                                                                          homepage
                                                                                                 [{"id":
                                                                                                                                             In the
                                    [{"id": 28,
                                                                                                                                             22nd
                                                                                                 1463.
                                     "name":
                                                                                                                                         century, a
                                                                                               "name":
                    237000000
                                                          http://www.avatarmovie.com/ 19995
                                                                                                                                                    150.437577
                                    "Action"},
                                                                                                                       en
                                                                                                                                 Avatar
                                                                                                "culture
                                                                                                                                         paraplegic
                                     {"id": 12,
                                                                                                clash"},
                                                                                                                                          Marine is
                                      "nam...
                                                                                                {"id":...
                                                                                                                                               di...
                                                                                                                                           Captain
                                                                                             [{"id": 270,
                                                                                                                                         Barbossa,
                                    [{"id": 12,
                                                                                                                           Pirates of the
                                                                                                "name":
                                                                                                                                              long
                                     "name":
                                                                                                                             Caribbean:
                 1 300000000
                                              http://disney.go.com/disneypictures/pirates/
                                                                                              "ocean"},
                                                                                                                                          believed 139.082615
                                                                                                                       en
                                "Adventure"},
                                                                                                                              At World's
                                                                                              {"id": 726,
                                                                                                                                             to be
                                 {"id": 14, "...
                                                                                                                                    End
                                                                                                  "na...
                                                                                                                                             dead,
                                                                                                                                              ha...
In [14]:
             # genres
                # movie id
                # Keywords
                # title
                # overview
                # cast
                # crew
 In [ ]:
```

```
In [15]:  M movies = movies[['movie_id', 'title', 'overview', 'genres', 'keywords', 'cast', 'crew']]
    movies.head()
```

crew	cast	keywords	genres	overview	title	movie_id	.5]:
[{"credit_id" "52fe48009251416c750aca23" "de	[{"cast_id": 242, "character": "Jake Sully", "	[{"id": 1463, "name": "culture clash"}, {"id":	[{"id": 28, "name": "Action"}, {"id": 12, "nam	In the 22nd century, a paraplegic Marine is di	Avatar	19995	0
[{"credit_id" "52fe4232c3a36847f800b579" "de	[{"cast_id": 4, "character": "Captain Jack Spa	[{"id": 270, "name": "ocean"}, {"id": 726, "na	[{"id": 12, "name": "Adventure"}, {"id": 14, "	Captain Barbossa, long believed to be dead, ha	Pirates of the Caribbean: At World's End	285	1
[{"credit_id" "54805967c3a36829b5002c41" "de	[{"cast_id": 1, "character": "James Bond", "cr	[{"id": 470, "name": "spy"}, {"id": 818, "name	[{"id": 28, "name": "Action"}, {"id": 12, "nam	A cryptic message from Bond's past sends him o	Spectre	206647	2
[{"credit_id" "52fe4781c3a36847f81398c3" "de	[{"cast_id": 2, "character": "Bruce Wayne / Ba	[{"id": 849, "name": "dc comics"}, {"id": 853,	[{"id": 28, "name": "Action"}, {"id": 80, "nam	Following the death of District Attorney Harve	The Dark Knight Rises	49026	3
[{"credit_id" "52fe479ac3a36847f813eaa3" "de	[{"cast_id": 5, "character": "John Carter", "c	[{"id": 818, "name": "based on novel"}, {"id":	[{"id": 28, "name": "Action"}, {"id": 12, "nam	John Carter is a war-weary, former military ca	John Carter	49529	4

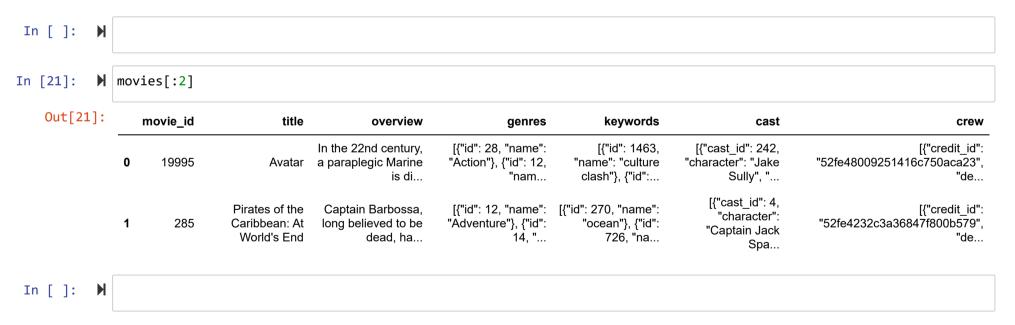
In []: •

```
    movies.info()

In [16]:
             <class 'pandas.core.frame.DataFrame'>
             RangeIndex: 4809 entries, 0 to 4808
             Data columns (total 7 columns):
                  Column
                            Non-Null Count Dtype
                  movie id 4809 non-null
                                            int64
                  title
                            4809 non-null
                                            object
              1
                  overview 4806 non-null
                                            object
                            4809 non-null
                  genres
                                            object
                  keywords 4809 non-null
                                            object
                            4809 non-null
                                            object
                  cast
                            4809 non-null
                                            object
                  crew
             dtypes: int64(1), object(6)
             memory usage: 263.1+ KB
In [17]:
          M movies.isna().sum()
   Out[17]: movie id
                         0
             title
                         0
             overview
                         3
             genres
             keywords
                         0
             cast
                         0
             crew
             dtype: int64
          movies.dropna(inplace=True)
In [18]:
```

```
movies.isna().sum()
In [19]:
   Out[19]: movie id
                          0
             title
                          0
             overview
                          0
                          0
             genres
             keywords
                          0
                          0
             cast
             crew
             dtype: int64
          M movies.duplicated().sum()
In [20]:
   Out[20]: 0
```

Final Outcome - movie_id + title + tags



Checking the records for specific column

```
movies['overview'][0]
In [22]:
   Out[22]: 'In the 22nd century, a paraplegic Marine is dispatched to the moon Pandora on a unique mission, but becomes torn
             between following orders and protecting an alien civilization.'
In [23]:
         movies['genres'][0]
   Out[23]: '[{"id": 28, "name": "Action"}, {"id": 12, "name": "Adventure"}, {"id": 14, "name": "Fantasy"}, {"id": 878, "nam
             e": "Science Fiction"}]'
In [24]:
          ▶ movies['keywords'][0]
   Out[24]: '[{"id": 1463, "name": "culture clash"}, {"id": 2964, "name": "future"}, {"id": 3386, "name": "space war"}, {"i
             d": 3388, "name": "space colony"}, {"id": 3679, "name": "society"}, {"id": 3801, "name": "space travel"}, {"id":
             9685, "name": "futuristic"}, {"id": 9840, "name": "romance"}, {"id": 9882, "name": "space"}, {"id": 9951, "name":
             "alien"}, {"id": 10148, "name": "tribe"}, {"id": 10158, "name": "alien planet"}, {"id": 10987, "name": "cgi"},
             {"id": 11399, "name": "marine"}, {"id": 13065, "name": "soldier"}, {"id": 14643, "name": "battle"}, {"id": 14720,
             "name": "love affair"}, {"id": 165431, "name": "anti war"}, {"id": 193554, "name": "power relations"}, {"id": 206
             690, "name": "mind and soul"}, {"id": 209714, "name": "3d"}]'
```

M movies['cast'][0] In [25]:

Out[25]: '[{"cast id": 242, "character": "Jake Sully", "credit id": "5602a8a7c3a3685532001c9a", "gender": 2, "id": 6573 1, "name": "Sam Worthington", "order": 0}, {"cast id": 3, "character": "Neytiri", "credit id": "52fe4800925141 6c750ac9cb", "gender": 1, "id": 8691, "name": "Zoe Saldana", "order": 1}, {"cast id": 25, "character": "Dr. Gr ace Augustine", "credit id": "52fe48009251416c750aca39", "gender": 1, "id": 10205, "name": "Sigourney Weaver", "order": 2}, {"cast id": 4, "character": "Col. Quaritch", "credit id": "52fe48009251416c750ac9cf", "gender": 2, "id": 32747, "name": "Stephen Lang", "order": 3}, {"cast id": 5, "character": "Trudy Chacon", "credit id": "52fe48009251416c750ac9d3", "gender": 1, "id": 17647, "name": "Michelle Rodriguez", "order": 4}, {"cast id": 8, "character": "Selfridge", "credit id": "52fe48009251416c750ac9e1", "gender": 2, "id": 1771, "name": "Giovan ni Ribisi", "order": 5}, {"cast id": 7, "character": "Norm Spellman", "credit id": "52fe48009251416c750ac9dd", "gender": 2, "id": 59231, "name": "Joel David Moore", "order": 6}, {"cast id": 9, "character": "Moat", "credit id": "52fe48009251416c750ac9e5", "gender": 1, "id": 30485, "name": "CCH Pounder", "order": 7}, {"cast id": 1 1, "character": "Eytukan", "credit id": "52fe48009251416c750ac9ed", "gender": 2, "id": 15853, "name": "Wes Stu di", "order": 8}, {"cast id": 10, "character": "Tsu\'Tey", "credit id": "52fe48009251416c750ac9e9", "gender": 2, "id": 10964, "name": "Laz Alonso", "order": 9}, {"cast id": 12, "character": "Dr. Max Patel", "credit id": "52fe48009251416c750ac9f1", "gender": 2, "id": 95697, "name": "Dileep Rao", "order": 10}, {"cast id": 13, "cha racter": "Lyle Wainfleet", "credit id": "52fe48009251416c750ac9f5", "gender": 2, "id": 98215, "name": "Matt Ge rald", "order": 11}, {"cast id": 32, "character": "Private Fike", "credit id": "52fe48009251416c750aca5b", "ge nder": 2, "id": 154153, "name": "Sean Anthony Moran", "order": 12}, {"cast id": 33, "character": "Cryo Vault M ed Tech", "credit id": "52fe48009251416c750aca5f", "gender": 2, "id": 397312, "name": "Jason Whyte", "order":

```
M movies['crew'][0]
In [26]:
   Out[26]: '[{"credit id": "52fe48009251416c750aca23", "department": "Editing", "gender": 0, "id": 1721, "job": "Editor",
             "name": "Stephen E. Rivkin"}, {"credit id": "539c47ecc3a36810e3001f87", "department": "Art", "gender": 2, "i
             d": 496, "job": "Production Design", "name": "Rick Carter"}, {"credit id": "54491c89c3a3680fb4001cf7", "depart
             ment": "Sound", "gender": 0, "id": 900, "job": "Sound Designer", "name": "Christopher Boyes"}, {"credit id":
             "54491cb70e0a267480001bd0", "department": "Sound", "gender": 0, "id": 900, "job": "Supervising Sound Editor",
             "name": "Christopher Boyes"}, {"credit id": "539c4a4cc3a36810c9002101", "department": "Production", "gender":
             1, "id": 1262, "job": "Casting", "name": "Mali Finn"}, {"credit id": "5544ee3b925141499f0008fc", "department":
             "Sound", "gender": 2, "id": 1729, "job": "Original Music Composer", "name": "James Horner"}, {"credit id": "52
             fe48009251416c750ac9c3", "department": "Directing", "gender": 2, "id": 2710, "job": "Director", "name": "James
             Cameron"}, {"credit id": "52fe48009251416c750ac9d9", "department": "Writing", "gender": 2, "id": 2710, "job":
             "Writer", "name": "James Cameron"}, {"credit_id": "52fe48009251416c750aca17", "department": "Editing", "gende
             r": 2, "id": 2710, "job": "Editor", "name": "James Cameron"}, {"credit id": "52fe48009251416c750aca29", "depar
             tment": "Production", "gender": 2, "id": 2710, "job": "Producer", "name": "James Cameron"}, {"credit_id": "52f
             e48009251416c750aca3f", "department": "Writing", "gender": 2, "id": 2710, "job": "Screenplay", "name": "James
             Cameron"}, {"credit id": "539c4987c3a36810ba0021a4", "department": "Art", "gender": 2, "id": 7236, "job": "Art
             Direction", "name": "Andrew Menzies"}, {"credit id": "549598c3c3a3686ae9004383", "department": "Visual Effect
             s", "gender": 0, "id": 6690, "job": "Visual Effects Producer", "name": "Jill Brooks"}, {"credit id": "52fe4800
             9251416c750aca4b", "department": "Production", "gender": 1, "id": 6347, "job": "Casting", "name": "Margery Sim
             kin"}, {"credit id": "570b6f419251417da70032fe", "department": "Art", "gender": 2, "id": 6878, "job": "Supervi
In [ ]:
```

Now let proceed with Genres column

```
In [29]:
         | import ast
            lst = ast.literal eval(movies['genres'][0])
            lst
   Out[29]: [{'id': 28, 'name': 'Action'},
             {'id': 12, 'name': 'Adventure'},
             {'id': 14, 'name': 'Fantasy'},
             {'id': 878, 'name': 'Science Fiction'}]
In [30]:
         Out[30]: {'id': 28, 'name': 'Action'}
In [31]: | dictt = lst[0]
            dictt.keys()
   Out[31]: dict keys(['id', 'name'])
In [32]:

▶ dictt.values()

   Out[32]: dict_values([28, 'Action'])

  | dictt['id']

In [33]:
   Out[33]: 28
         dictt['name']
In [34]:
   Out[34]: 'Action'
         # for i in movies['genres'][0]:
In [35]:
                  print(i)
```

```
In [36]:
              print(i)
              print(i['name'])
              print('----')
           {'id': 28, 'name': 'Action'}
           Action
           {'id': 12, 'name': 'Adventure'}
           Adventure
           {'id': 14, 'name': 'Fantasy'}
           Fantasy
           {'id': 878, 'name': 'Science Fiction'}
           Science Fiction
In [ ]:
In [37]:

  | def genres_col(dictt):
              List = []
              for i in ast.literal_eval(dictt):
                  List.append(i['name'])
              return List
In [38]:

■ genres_col(movies['genres'][0])
   Out[38]: ['Action', 'Adventure', 'Fantasy', 'Science Fiction']
```

```
movies['genres'].apply(genres col)
In [39]:
    Out[39]: 0
                          [Action, Adventure, Fantasy, Science Fiction]
                1
                                               [Adventure, Fantasy, Action]
                                                  [Action, Adventure, Crime]
                2
                3
                                          [Action, Crime, Drama, Thriller]
                                     [Action, Adventure, Science Fiction]
                4
                4804
                                                   [Action, Crime, Thriller]
                4805
                                                             [Comedy, Romance]
                4806
                                        [Comedy, Drama, Romance, TV Movie]
                4807
                4808
                                                                  [Documentary]
                Name: genres, Length: 4806, dtype: object
In [40]:

    movies[:2]

    Out[40]:
                    movie_id
                                           title
                                                                                              keywords
                                                         overview
                                                                                                                       cast
                                                                               genres
                                                                                                                                                       crew
                                                                      [{"id": 28, "name":
                                                                                                             [{"cast id": 242,
                                                                                                                                                 [{"credit id":
                                                In the 22nd century,
                                                                                             [{"id": 1463,
                                                                      "Action"}, {"id": 12,
                 0
                       19995
                                         Avatar
                                                a paraplegic Marine
                                                                                          "name": "culture
                                                                                                           "character": "Jake
                                                                                                                                 "52fe48009251416c750aca23",
                                                            is di...
                                                                               "nam...
                                                                                           clash"}, {"id":...
                                                                                                                  Sully", "...
                                                                                                                                                       "de...
                                                                                                               [{"cast id": 4,
                                                                                                                                                 [{"credit id":
                                   Pirates of the
                                                 Captain Barbossa.
                                                                      [{"id": 12, "name":
                                                                                       [{"id": 270, "name":
                                                                                                                 "character":
                                                                     "Adventure"}, {"id":
                 1
                         285
                                  Caribbean: At
                                                 long believed to be
                                                                                           "ocean"}, {"id":
                                                                                                                                 "52fe4232c3a36847f800b579",
                                                                                                               "Captain Jack
                                                                               14, "...
                                    World's End
                                                        dead, ha...
                                                                                              726, "na...
                                                                                                                                                       "de...
                                                                                                                     Spa...
In [41]:
             movies['genres'] = movies['genres'].apply(genres col)
```

In [42]: ▶	mov	ies[:5]						
Out[42]:		movie_id	title	overview	genres	keywords	cast	crew
	0	19995	Avatar	In the 22nd century, a paraplegic Marine is di	[Action, Adventure, Fantasy, Science Fiction]	[{"id": 1463, "name": "culture clash"}, {"id":	[{"cast_id": 242, "character": "Jake Sully", "	[{"credit_id": "52fe48009251416c750aca23", "de
	1	285	Pirates of the Caribbean: At World's End	Captain Barbossa, long believed to be dead, ha	[Adventure, Fantasy, Action]	[{"id": 270, "name": "ocean"}, {"id": 726, "na	[{"cast_id": 4, "character": "Captain Jack Spa	[{"credit_id": "52fe4232c3a36847f800b579", "de
	2	206647	Spectre	A cryptic message from Bond's past sends him o	[Action, Adventure, Crime]	[{"id": 470, "name": "spy"}, {"id": 818, "name	[{"cast_id": 1, "character": "James Bond", "cr	[{"credit_id": "54805967c3a36829b5002c41", "de
	3	49026	The Dark Knight Rises	Following the death of District Attorney Harve	[Action, Crime, Drama, Thriller]	[{"id": 849, "name": "dc comics"}, {"id": 853,	[{"cast_id": 2, "character": "Bruce Wayne / Ba	[{"credit_id": "52fe4781c3a36847f81398c3", "de
	4	49529	John Carter	John Carter is a war-weary, former military ca	[Action, Adventure, Science Fiction]	[{"id": 818, "name": "based on novel"}, {"id":	[{"cast_id": 5, "character": "John Carter", "c	[{"credit_id": "52fe479ac3a36847f813eaa3", "de
In []: ▶								

Now same goes for keywords column

```
movies['keywords']
In [43]:
   Out[43]: 0
                     [{"id": 1463, "name": "culture clash"}, {"id":...
                     [{"id": 270, "name": "ocean"}, {"id": 726, "na...
             1
                     [{"id": 470, "name": "spv"}, {"id": 818, "name...
             3
                     [{"id": 849, "name": "dc comics"}, {"id": 853,...
                     [{"id": 818, "name": "based on novel"}, {"id":...
                     [{"id": 5616, "name": "united states\u2013mexi...
             4804
             4805
                     [{"id": 248, "name": "date"}, {"id": 699, "nam...
             4806
             4807
                     [{"id": 1523, "name": "obsession"}, {"id": 224...
             4808
             Name: keywords, Length: 4806, dtype: object
In [44]: | movies['keywords'][0]
   Out[44]: '[{"id": 1463, "name": "culture clash"}, {"id": 2964, "name": "future"}, {"id": 3386, "name": "space war"}, {"i
             d": 3388, "name": "space colony"}, {"id": 3679, "name": "society"}, {"id": 3801, "name": "space travel"}, {"id":
             9685, "name": "futuristic"}, {"id": 9840, "name": "romance"}, {"id": 9882, "name": "space"}, {"id": 9951, "name":
             "alien"}, {"id": 10148, "name": "tribe"}, {"id": 10158, "name": "alien planet"}, {"id": 10987, "name": "cgi"},
             {"id": 11399, "name": "marine"}, {"id": 13065, "name": "soldier"}, {"id": 14643, "name": "battle"}, {"id": 14720,
             "name": "love affair"}, {"id": 165431, "name": "anti war"}, {"id": 193554, "name": "power relations"}, {"id": 206
             690, "name": "mind and soul"}, {"id": 209714, "name": "3d"}]'
 In [ ]:
In [45]: ▶ def keywords col(dictt):
                 List = []
                 for i in ast.literal eval(dictt):
                     List.append(i['name'])
                 return List
```

```
keywords col(movies['keywords'][0])
In [46]:
   Out[46]: ['culture clash',
               'future',
               'space war',
               'space colony',
               'society',
               'space travel',
               'futuristic',
               'romance',
               'space',
               'alien',
               'tribe',
               'alien planet',
               'cgi',
               'marine',
               'soldier',
               'battle',
               'love affair',
               'anti war',
               'power relations',
               'mind and soul',
               '3d'1
          movies['keywords'].apply(keywords col)
In [47]:
   Out[47]: 0
                      [culture clash, future, space war, space colon...
                      [ocean, drug abuse, exotic island, east india ...
             1
                      [spy, based on novel, secret agent, sequel, mi...
             2
             3
                      [dc comics, crime fighter, terrorist, secret i...
             4
                      [based on novel, mars, medallion, space travel...
                      [united states-mexico barrier, legs, arms, pap...
             4804
             4805
                                                                      []
                      [date, love at first sight, narration, investi...
             4806
             4807
             4808
                              [obsession, camcorder, crush, dream girl]
             Name: keywords, Length: 4806, dtype: object
```

```
movies['keywords'] = movies['keywords'].apply(keywords col)
In [48]:
In [49]:
                  movies[:4]
     Out[49]:
                       movie_id
                                                title
                                                                 overview
                                                                                        genres
                                                                                                           keywords
                                                                                                                                       cast
                                                                                                                                                                            crew
                                                                             [Action, Adventure,
                                                                                                                                                                     [{"credit id":
                                                       In the 22nd century,
                                                                                                       [culture clash,
                                                                                                                            [{"cast id": 242,
                                                                                                                          "character": "Jake
                   0
                          19995
                                              Avatar
                                                       a paraplegic Marine
                                                                               Fantasy, Science
                                                                                                   future, space war,
                                                                                                                                                  "52fe48009251416c750aca23",
                                                                                                                                 Sully", "...
                                                                    is di...
                                                                                        Fiction]
                                                                                                       space colon...
                                                                                                                                                                            "de...
                                                                                                                              [{"cast id": 4,
                                       Pirates of the
                                                        Captain Barbossa,
                                                                                                                                                                     [{"credit id":
                                                                                                         [ocean, drug
                                                                                                                                "character":
                                                                                    [Adventure,
                   1
                             285
                                       Caribbean: At
                                                        long believed to be
                                                                                                        abuse, exotic
                                                                                                                                                   "52fe4232c3a36847f800b579",
                                                                                Fantasy, Action]
                                                                                                                             "Captain Jack
                                                                                                  island, east india ...
                                        World's End
                                                                dead, ha...
                                                                                                                                                                            "de...
                                                                                                                                     Spa...
                                                                                                                              [{"cast id": 1,
                                                                                                                                                                     [{"credit id":
                                                         A cryptic message
                                                                                                       [spy, based on
                                                                             [Action, Adventure,
                   2
                         206647
                                             Spectre
                                                          from Bond's past
                                                                                                        novel, secret
                                                                                                                        "character": "James
                                                                                                                                                  "54805967c3a36829b5002c41",
                                                                                         Crime]
                                                             sends him o...
                                                                                                  agent, seguel, mi...
                                                                                                                               Bond", "cr...
                                                                                                                                                                            "de...
                                                                                                                              [{"cast id": 2,
                                                       Following the death
                                                                                                    [dc comics, crime
                                    The Dark Knight
                                                                                 [Action, Crime,
                                                                                                                                                                     [{"credit id":
                                                                                                                        "character": "Bruce
                   3
                          49026
                                                        of District Attorney
                                                                                                     fighter, terrorist,
                                                                                                                                             "52fe4781c3a36847f81398c3", "de...
                                               Rises
                                                                                 Drama, Thriller]
                                                                  Harve...
                                                                                                            secret i...
                                                                                                                              Wayne / Ba...
 In [ ]:
```

Now working with Cast column

```
M movies['cast']
In [50]:
   Out[50]: 0
                     [{"cast id": 242, "character": "Jake Sully", "...
                     [{"cast id": 4, "character": "Captain Jack Spa...
             1
             2
                     [{"cast id": 1, "character": "James Bond", "cr...
                     [{"cast id": 2, "character": "Bruce Wayne / Ba...
             3
                     [{"cast id": 5, "character": "John Carter", "c...
                     [{"cast id": 1, "character": "El Mariachi", "c...
             4804
             4805
                     [{"cast id": 1, "character": "Buzzy", "credit ...
                     [{"cast id": 8, "character": "Oliver 0\u2019To...
             4806
                     [{"cast id": 3, "character": "Sam", "credit id...
             4807
                     [{"cast id": 3, "character": "Herself", "credi...
             4808
             Name: cast, Length: 4806, dtype: object
In [51]: | movies['cast'][0]
   Out[51]: '[{"cast id": 242, "character": "Jake Sully", "credit id": "5602a8a7c3a3685532001c9a", "gender": 2, "id": 6573
             1, "name": "Sam Worthington", "order": 0}, {"cast id": 3, "character": "Neytiri", "credit id": "52fe4800925141
             6c750ac9cb", "gender": 1, "id": 8691, "name": "Zoe Saldana", "order": 1}, {"cast id": 25, "character": "Dr. Gr
             ace Augustine", "credit id": "52fe48009251416c750aca39", "gender": 1, "id": 10205, "name": "Sigourney Weaver",
             "order": 2}, {"cast id": 4, "character": "Col. Quaritch", "credit id": "52fe48009251416c750ac9cf", "gender":
             2, "id": 32747, "name": "Stephen Lang", "order": 3}, {"cast id": 5, "character": "Trudy Chacon", "credit id":
             "52fe48009251416c750ac9d3", "gender": 1, "id": 17647, "name": "Michelle Rodriguez", "order": 4}, {"cast id":
             8, "character": "Selfridge", "credit id": "52fe48009251416c750ac9e1", "gender": 2, "id": 1771, "name": "Giovan
             ni Ribisi", "order": 5}, {"cast id": 7, "character": "Norm Spellman", "credit id": "52fe48009251416c750ac9dd",
             "gender": 2, "id": 59231, "name": "Joel David Moore", "order": 6}, {"cast id": 9, "character": "Moat", "credit
             id": "52fe48009251416c750ac9e5", "gender": 1, "id": 30485, "name": "CCH Pounder", "order": 7}, {"cast id": 1
             1, "character": "Eytukan", "credit id": "52fe48009251416c750ac9ed", "gender": 2, "id": 15853, "name": "Wes Stu
```

di", "order": 8}, {"cast id": 10, "character": "Tsu\'Tey", "credit id": "52fe48009251416c750ac9e9", "gender": 2, "id": 10964, "name": "Laz Alonso", "order": 9}, {"cast id": 12, "character": "Dr. Max Patel", "credit id": "52fe48009251416c750ac9f1", "gender": 2, "id": 95697, "name": "Dileep Rao", "order": 10}, {"cast id": 13, "cha racter": "Lyle Wainfleet", "credit id": "52fe48009251416c750ac9f5", "gender": 2, "id": 98215, "name": "Matt Ge rald", "order": 11}, {"cast id": 32, "character": "Private Fike", "credit id": "52fe48009251416c750aca5b", "ge nder": 2, "id": 154153, "name": "Sean Anthony Moran", "order": 12}, {"cast id": 33, "character": "Cryo Vault M ed Tech", "credit_id": "52fe48009251416c750aca5f", "gender": 2, "id": 397312, "name": "Jason Whyte", "order":

```
ast.literal eval(movies['cast'][0])
In [52]:
   Out[52]: [{'cast id': 242,
               'character': 'Jake Sully',
               'credit id': '5602a8a7c3a3685532001c9a',
               'gender': 2,
               'id': 65731,
               'name': 'Sam Worthington',
               'order': 0},
              {'cast id': 3,
               'character': 'Neytiri',
               'credit id': '52fe48009251416c750ac9cb',
               'gender': 1,
               'id': 8691,
               'name': 'Zoe Saldana',
               'order': 1},
              {'cast id': 25,
               'character': 'Dr. Grace Augustine',
               'credit id': '52fe48009251416c750aca39',
               'gender': 1,
               'id': 10205,
               . . . . . . . . .
          M movies['cast'][0][0:451]
In [53]:
   Out[53]: '[{"cast_id": 242, "character": "Jake Sully", "credit_id": "5602a8a7c3a3685532001c9a", "gender": 2, "id": 65731,
             "name": "Sam Worthington", "order": 0}, {"cast id": 3, "character": "Neytiri", "credit id": "52fe48009251416c750a
             c9cb", "gender": 1, "id": 8691, "name": "Zoe Saldana", "order": 1}, {"cast id": 25, "character": "Dr. Grace Augus
             tine", "credit id": "52fe48009251416c750aca39", "gender": 1, "id": 10205, "name": "Sigourney Weaver", "order":
             2}'
```

```
    def cast col(dictt):

In [54]:
                 List = []
                 counter = 0
                 for i in ast.literal eval(dictt):
                     if counter != 3:
                         List.append(i['name'])
                         counter += 1
                                                             \# counter = counter + 1
                     else:
                         break
                 return List
          movies['cast'].apply(cast col)
In [55]:
   Out[55]: 0
                      [Sam Worthington, Zoe Saldana, Sigourney Weaver]
             1
                         [Johnny Depp, Orlando Bloom, Keira Knightley]
             2
                          [Daniel Craig, Christoph Waltz, Léa Seydoux]
             3
                          [Christian Bale, Michael Caine, Gary Oldman]
             4
                        [Taylor Kitsch, Lynn Collins, Samantha Morton]
                     [Carlos Gallardo, Jaime de Hoyos, Peter Marqua...
             4804
             4805
                          [Edward Burns, Kerry Bishé, Marsha Dietlein]
             4806
                            [Eric Mabius, Kristin Booth, Crystal Lowe]
                             [Daniel Henney, Eliza Coupe, Bill Paxton]
             4807
                     [Drew Barrymore, Brian Herzlinger, Corey Feldman]
             4808
             Name: cast, Length: 4806, dtype: object
```

In [56]:	M	mov:	ies[:3]						
Out[56]:		movie_id	title	overview	genres	keywords	cast	crew
		0	19995	Avatar	In the 22nd century, a paraplegic Marine is di	[Action, Adventure, Fantasy, Science Fiction]	[culture clash, future, space war, space colon	[{"cast_id": 242, "character": "Jake Sully", "	[{"credit_id": "52fe48009251416c750aca23", "de
		1	285	Pirates of the Caribbean: At World's End	Captain Barbossa, long believed to be dead, ha	[Adventure, Fantasy, Action]	[ocean, drug abuse, exotic island, east india	[{"cast_id": 4, "character": "Captain Jack Spa	[{"credit_id": "52fe4232c3a36847f800b579", "de
		2	206647	Spectre	A cryptic message from Bond's past sends him o	[Action, Adventure, Crime]	[spy, based on novel, secret agent, sequel, mi	[{"cast_id": 1, "character": "James Bond", "cr	[{"credit_id": "54805967c3a36829b5002c41", "de
In [57]:	H	mov:	ies['cast	t'] = movies['	<pre>cast'].apply(ca</pre>	st_col)			
In [58]:	M	mov:	ies[:4]						
Out[58]:		movie_id	title	overview	genres	keywords	cast	crew
		0	19995	Avatar	In the 22nd century, a paraplegic Marine is di	[Action, Adventure, Fantasy, Science Fiction]	[culture clash, future, space war, space colon	[Sam Worthington, Zoe Saldana, Sigourney Weaver]	[{"credit_id": "52fe48009251416c750aca23", "de
		1	285	Pirates of the Caribbean: At World's End	Captain Barbossa, long believed to be dead, ha	[Adventure, Fantasy, Action]	[ocean, drug abuse, exotic island, east india 	[Johnny Depp, Orlando Bloom, Keira Knightley]	[{"credit_id": "52fe4232c3a36847f800b579", "de
		2	206647	Spectre	A cryptic message from Bond's past sends him o	[Action, Adventure, Crime]	[spy, based on novel, secret agent, sequel, mi	[Daniel Craig, Christoph Waltz, Léa Seydoux]	[{"credit_id": "54805967c3a36829b5002c41", "de
		3	49026	The Dark Knight Rises	Following the death of District Attorney Harve	[Action, Crime, Drama, Thriller]	[dc comics, crime fighter, terrorist, secret i	[Christian Bale, Michael Caine, Gary Oldman]	[{"credit_id": "52fe4781c3a36847f81398c3", "de

```
In [ ]: N
```

Crew Column

```
In [59]:
          movies['crew']
   Out[59]: 0
                     [{"credit id": "52fe48009251416c750aca23", "de...
                    [{"credit id": "52fe4232c3a36847f800b579", "de...
                     [{"credit id": "54805967c3a36829b5002c41", "de...
                     [{"credit_id": "52fe4781c3a36847f81398c3", "de...
                     [{"credit_id": "52fe479ac3a36847f813eaa3", "de...
                     [{"credit id": "52fe44eec3a36847f80b280b", "de...
             4804
                     [{"credit id": "52fe487dc3a368484e0fb013", "de...
             4805
                     [{"credit id": "52fe4df3c3a36847f8275ecf", "de...
             4806
                     [{"credit id": "52fe4ad9c3a368484e16a36b", "de...
             4807
                     [{"credit id": "58ce021b9251415a390165d9", "de...
             4808
             Name: crew, Length: 4806, dtype: object
```

```
M movies['crew'][0]
In [60]:
   Out[60]: '[{"credit id": "52fe48009251416c750aca23", "department": "Editing", "gender": 0, "id": 1721, "job": "Editor",
             "name": "Stephen E. Rivkin"}, {"credit id": "539c47ecc3a36810e3001f87", "department": "Art", "gender": 2, "i
             d": 496, "job": "Production Design", "name": "Rick Carter"}, {"credit id": "54491c89c3a3680fb4001cf7", "depart
             ment": "Sound", "gender": 0, "id": 900, "job": "Sound Designer", "name": "Christopher Boyes"}, {"credit id":
             "54491cb70e0a267480001bd0", "department": "Sound", "gender": 0, "id": 900, "job": "Supervising Sound Editor",
             "name": "Christopher Boyes"}, {"credit id": "539c4a4cc3a36810c9002101", "department": "Production", "gender":
             1, "id": 1262, "job": "Casting", "name": "Mali Finn"}, {"credit id": "5544ee3b925141499f0008fc", "department":
             "Sound", "gender": 2, "id": 1729, "job": "Original Music Composer", "name": "James Horner"}, {"credit id": "52
             fe48009251416c750ac9c3", "department": "Directing", "gender": 2, "id": 2710, "job": "Director", "name": "James
             Cameron"}, {"credit id": "52fe48009251416c750ac9d9", "department": "Writing", "gender": 2, "id": 2710, "job":
             "Writer", "name": "James Cameron"}, {"credit_id": "52fe48009251416c750aca17", "department": "Editing", "gende
             r": 2, "id": 2710, "job": "Editor", "name": "James Cameron"}, {"credit id": "52fe48009251416c750aca29", "depar
             tment": "Production", "gender": 2, "id": 2710, "job": "Producer", "name": "James Cameron"}, {"credit_id": "52f
             e48009251416c750aca3f", "department": "Writing", "gender": 2, "id": 2710, "job": "Screenplay", "name": "James
             Cameron"}, {"credit id": "539c4987c3a36810ba0021a4", "department": "Art", "gender": 2, "id": 7236, "job": "Art
             Direction", "name": "Andrew Menzies"}, {"credit id": "549598c3c3a3686ae9004383", "department": "Visual Effect
             s", "gender": 0, "id": 6690, "job": "Visual Effects Producer", "name": "Jill Brooks"}, {"credit id": "52fe4800
             9251416c750aca4b", "department": "Production", "gender": 1, "id": 6347, "job": "Casting", "name": "Margery Sim
             kin"}, {"credit id": "570b6f419251417da70032fe", "department": "Art", "gender": 2, "id": 6878, "job": "Supervi
In [74]:  def fetch director(data):
                 L = []
                 for i in ast.literal eval(data):
                         print(i['job'])
```

```
fetch director(movies['crew'][0])
In [75]:
             Editor
             Production Design
             Sound Designer
             Supervising Sound Editor
             Casting
             Original Music Composer
             Director
             Writer
             Editor
             Producer
             Screenplay
             Art Direction
             Visual Effects Producer
             Casting
             Supervising Art Director
             Music Editor
             Sound Effects Editor
             Foley
             Foley
 In [ ]:

    def fetch_director(data):

In [76]:
                 L = []
                 for i in ast.literal eval(data):
                     if i['job'] == 'Director':
                        L.append(i['name'])
                 return L
          fetch_director(movies['crew'][0])
In [77]:
   Out[77]: ['James Cameron']
```

```
movies['crew'].apply(fetch director)
In [78]:
    Out[78]: 0
                                                           [James Cameron]
                                                          [Gore Verbinski]
                 1
                 2
                                                               [Sam Mendes]
                 3
                                                      [Christopher Nolan]
                 4
                                                          [Andrew Stanton]
                 4804
                                                        [Robert Rodriguez]
                 4805
                                                             [Edward Burns]
                 4806
                                                              [Scott Smith]
                 4807
                                                              [Daniel Hsia]
                           [Brian Herzlinger, Jon Gunn, Brett Winn]
                 4808
                 Name: crew, Length: 4806, dtype: object
In [79]:
                movies[:4]
    Out[79]:
                     movie_id
                                            title
                                                           overview
                                                                                                 keywords
                                                                                genres
                                                                                                                            cast
                                                                                                                                                             crew
                                                                      [Action, Adventure,
                                                                                                               [Sam Worthington,
                                                                                                                                                       [{"credit id":
                                                  In the 22nd century,
                                                                                             [culture clash,
                  0
                        19995
                                          Avatar a paraplegic Marine
                                                                       Fantasy, Science
                                                                                                                   Zoe Saldana,
                                                                                                                                      "52fe48009251416c750aca23",
                                                                                          future, space war,
                                                                                Fiction]
                                                                                                              Sigourney Weaver]
                                                                                                                                                              "de...
                                                              is di...
                                                                                             space colon...
                                                                                               [ocean, drug
                                    Pirates of the
                                                   Captain Barbossa,
                                                                                                                  [Johnny Depp,
                                                                                                                                                       [{"credit id":
                                                                            [Adventure,
                                                                                              abuse, exotic
                  1
                          285
                                    Caribbean: At
                                                   long believed to be
                                                                                                             Orlando Bloom, Keira
                                                                                                                                       "52fe4232c3a36847f800b579",
                                                                         Fantasy, Action]
                                                                                           island, east india
                                     World's End
                                                                                                                       Knightley]
                                                                                                                                                              "de...
                                                          dead, ha...
                                                   A cryptic message
                                                                                             [spy, based on
                                                                                                                   [Daniel Craig,
                                                                                                                                                       [{"credit id":
                                                                      [Action, Adventure,
                  2
                       206647
                                                    from Bond's past
                                                                                                             Christoph Waltz, Léa
                                                                                                                                      "54805967c3a36829b5002c41",
                                         Spectre
                                                                                              novel, secret
                                                                                Crime]
                                                       sends him o...
                                                                                         agent, sequel, mi...
                                                                                                                       Seydoux1
                                                                                                                                                              "de...
                                                                                                                  [Christian Bale,
                                                  Following the death
                                                                                          [dc comics, crime
                                                                                                                                                       [{"credit id":
                                                                         [Action, Crime,
                                 The Dark Knight
                  3
                        49026
                                                                                                             Michael Caine, Gary
                                                                                                                                       "52fe4781c3a36847f81398c3",
                                                   of District Attorney
                                                                                            fighter, terrorist,
                                           Rises
                                                                         Drama, Thriller]
                                                            Harve...
                                                                                                 secret i...
                                                                                                                        Oldman]
                                                                                                                                                              "de...
In [80]:
             movies['crew'] = movies['crew'].apply(fetch director)
```

In [81]: ▶	movi	ies[:4]						
Out[81]:	İ	movie_id	title	overview	genres	keywords	cast	crew
	0	19995	Avatar	In the 22nd century, a paraplegic Marine is di	[Action, Adventure, Fantasy, Science Fiction]	[culture clash, future, space war, space colon	[Sam Worthington, Zoe Saldana, Sigourney Weaver]	[James Cameron]
	1	285	Pirates of the Caribbean: At World's End	Captain Barbossa, long believed to be dead, ha	[Adventure, Fantasy, Action]	[ocean, drug abuse, exotic island, east india	[Johnny Depp, Orlando Bloom, Keira Knightley]	[Gore Verbinski]
	2	206647	Spectre	A cryptic message from Bond's past sends him o	[Action, Adventure, Crime]	[spy, based on novel, secret agent, sequel, mi	[Daniel Craig, Christoph Waltz, Léa Seydoux]	[Sam Mendes]
	3	49026	The Dark Knight Rises	Following the death of District Attorney Harve	[Action, Crime, Drama, Thriller]	[dc comics, crime fighter, terrorist, secret i	[Christian Bale, Michael Caine, Gary Oldman]	[Christopher Nolan]
In []: 🕨								

Overview Column

```
movies['overview']
In [82]:
   Out[82]: 0
                     In the 22nd century, a paraplegic Marine is di...
                     Captain Barbossa, long believed to be dead, ha...
             1
             2
                     A cryptic message from Bond's past sends him o...
             3
                     Following the death of District Attorney Harve...
             4
                     John Carter is a war-weary, former military ca...
                     El Mariachi just wants to play his guitar and ...
             4804
                     A newlywed couple's honeymoon is upended by th...
             4805
             4806
                     "Signed, Sealed, Delivered" introduces a dedic...
                     When ambitious New York attorney Sam is sent t...
             4807
                     Ever since the second grade when he first saw ...
             4808
             Name: overview, Length: 4806, dtype: object
```

```
movies['overview'][0].split()
In [86]:
   Out[86]: ['In',
               'the',
               '22nd',
               'century,',
               'a',
               'paraplegic',
               'Marine',
               'is',
               'dispatched',
               'to',
               'the',
               'moon',
               'Pandora',
               'on',
               'a',
               'unique',
               'mission,',
               'but',
               'becomes',
               'torn',
               'between',
               'following',
               'orders',
               'and',
               'protecting',
               'an',
               'alien',
               'civilization.']
In [87]:
           movies['overview'] = movies['overview'].apply(lambda x: x.split())
```

```
M movies[:5]
In [88]:
    Out[88]:
                     movie_id
                                               title
                                                                  overview
                                                                                                                keywords
                                                                                                                                               cast
                                                                                           genres
                                                                                                                                                              crew
                                                                                                                               [Sam Worthington, Zoe
                                                                                 [Action, Adventure,
                                                                                                       [culture clash, future,
                                                      [In, the, 22nd, century.,
                                                                                                                                                            [James
                  0
                        19995
                                             Avatar
                                                                                  Fantasy, Science
                                                                                                                                  Saldana, Sigourney
                                                                                                          space war, space
                                                       a, paraplegic, Marin...
                                                                                                                                                          Cameron]
                                                                                           Fiction]
                                                                                                                  colon...
                                                                                                                                            Weaverl
                                       Pirates of the
                                                        [Captain, Barbossa,,
                                                                                                       [ocean, drug abuse,
                                                                               [Adventure, Fantasy,
                                                                                                                               [Johnny Depp. Orlando
                                                                                                                                                             [Gore
                  1
                          285
                                       Caribbean: At
                                                        long, believed, to, be,
                                                                                                    exotic island, east india
                                                                                            Action1
                                                                                                                              Bloom, Keira Knightley]
                                                                                                                                                          Verbinski]
                                        World's End
                                                                        d...
                                                        [A, cryptic, message,
                                                                                                      [spy, based on novel,
                                                                                 [Action, Adventure,
                                                                                                                              [Daniel Craig, Christoph
                  2
                       206647
                                                          from, Bond's, past,
                                                                                                                                                      [Sam Mendes]
                                            Spectre
                                                                                                      secret agent, sequel,
                                                                                                                                 Waltz, Léa Sevdouxl
                                                                                            Crime]
                                                                     send...
                                                                                                         [dc comics, crime
                                    The Dark Knight
                                                       [Following, the, death,
                                                                             [Action, Crime, Drama,
                                                                                                                              [Christian Bale, Michael
                                                                                                                                                        [Christopher
                  3
                        49026
                                                                                                     fighter, terrorist, secret
                                              Rises
                                                        of, District, Attorney...
                                                                                           Thriller]
                                                                                                                                Caine, Gary Oldman]
                                                                                                                                                             Nolan]
                                                                                                     [based on novel, mars,
                                                     [John, Carter, is, a, war-
                                                                                [Action, Adventure,
                                                                                                                                  [Taylor Kitsch, Lynn
                                                                                                                                                           [Andrew
                  4
                        49529
                                        John Carter
                                                                                                          medallion, space
                                                                                                                           Collins, Samantha Morton]
                                                        weary,, former, mili...
                                                                                   Science Fiction]
                                                                                                                                                           Stanton]
                                                                                                                  travel...
 In [ ]:
             H
             movies['genres'] = movies['genres'].apply(lambda x: [i.replace(" ", "") for i in x])
In [89]:
             movies['keywords'] = movies['keywords'].apply(lambda x: [i.replace(" ", "") for i in x])
In [91]:
                 movies['cast'] = movies['cast'].apply(lambda x: [i.replace(" ", "") for i in x])
                 movies['crew'] = movies['crew'].apply(lambda x: [i.replace(" ", "") for i in x])
```

In [92]: ▶	mov	ies.head())					
Out[92]:		movie_id	title	overview	genres	keywords	cast	crew
	0	19995	Avatar	[In, the, 22nd, century,, a, paraplegic, Marin	[Action, Adventure, Fantasy, ScienceFiction]	[cultureclash, future, spacewar, spacecolony,	[SamWorthington, ZoeSaldana, SigourneyWeaver]	[JamesCameron]
	1	285	Pirates of the Caribbean: At World's End	[Captain, Barbossa,, long, believed, to, be, d	[Adventure, Fantasy, Action]	[ocean, drugabuse, exoticisland, eastindiatrad	[JohnnyDepp, OrlandoBloom, KeiraKnightley]	[GoreVerbinski]
	2	206647	Spectre	[A, cryptic, message, from, Bond's, past, send	[Action, Adventure, Crime]	[spy, basedonnovel, secretagent, sequel, mi6,	[DanielCraig, ChristophWaltz, LéaSeydoux]	[SamMendes]
	3	49026	The Dark Knight Rises	[Following, the, death, of, District, Attorney	[Action, Crime, Drama, Thriller]	[dccomics, crimefighter, terrorist, secretiden	[ChristianBale, MichaelCaine, GaryOldman]	[ChristopherNolan]
	4	49529	John Carter	[John, Carter, is, a, war-weary,, former, mili	[Action, Adventure, ScienceFiction]	[basedonnovel, mars, medallion, spacetravel, p	[TaylorKitsch, LynnCollins, SamanthaMorton]	[AndrewStanton]
In []: 🕨								

Now create a 'Tags' column

```
In [93]:  M movies['tags'] = movies['overview'] + movies['genres'] + movies['keywords'] + movies['cast'] + movies['crew']
```

```
    movies[:3]

In [94]:
    Out[94]:
                    movie id
                                        title
                                                   overview
                                                                                         keywords
                                                                                                                    cast
                                                                        genres
                                                                                                                                   crew
                                                                                                                                                    tags
                                                [In, the, 22nd,
                                                                                                                                            [In, the, 22nd,
                                                               [Action, Adventure,
                                                                                [cultureclash, future.
                                                                                                        [SamWorthington,
                                                  century,, a,
                                                                                                                                               century,, a,
                      19995
                                                                                                             ZoeSaldana,
                                                                                                                         [JamesCameron]
                 0
                                      Avatar
                                                                       Fantasy,
                                                                                         spacewar,
                                                  paraplegic,
                                                                                                                                               paraplegic,
                                                                                    spacecolony, ...
                                                                 ScienceFiction]
                                                                                                        SigourneyWeaver]
                                                     Marin...
                                                                                                                                                 Marin...
                                                    [Captain,
                                                                                                                                                [Captain,
                                Pirates of the
                                                                                                           [JohnnyDepp.
                                                                                 [ocean, drugabuse,
                                              Barbossa,, long,
                                                                    [Adventure,
                                                                                                                                          Barbossa,, long,
                 1
                        285
                               Caribbean: At
                                                                                       exoticisland.
                                                                                                           OrlandoBloom.
                                                                                                                           [GoreVerbinski]
                                              believed, to, be,
                                                                 Fantasy, Action]
                                                                                                                                           believed, to, be,
                                 World's End
                                                                                     eastindiatrad...
                                                                                                           KeiraKnightley]
                                                         d...
                                                                                                                                                     d...
                                                  [A, cryptic,
                                                                                                                                               [A, cryptic,
                                                                                 [spy, basedonnovel,
                                                                                                            [DanielCraig,
                                              message, from,
                                                              [Action, Adventure,
                                                                                                                                           message, from,
                 2
                     206647
                                                                                                          ChristophWaltz,
                                    Spectre
                                                                                secretagent, sequel,
                                                                                                                            [SamMendes]
                                                 Bond's, past,
                                                                         Crime]
                                                                                                                                             Bond's, past,
                                                                                                            LéaSeydoux]
                                                                                           mi6, ...
                                                      send...
                                                                                                                                                  send...
            print(movies['tags'][0])
In [96]:
               ['In', 'the', '22nd', 'century,', 'a', 'paraplegic', 'Marine', 'is', 'dispatched', 'to', 'the', 'moon', 'Pandor
               a', 'on', 'a', 'unique', 'mission,', 'but', 'becomes', 'torn', 'between', 'following', 'orders', 'and', 'protecti
               ng', 'an', 'alien', 'civilization.', 'Action', 'Adventure', 'Fantasy', 'ScienceFiction', 'cultureclash', 'futur
               e', 'spacewar', 'spacecolony', 'society', 'spacetravel', 'futuristic', 'romance', 'space', 'alien', 'tribe', 'ali
               enplanet', 'cgi', 'marine', 'soldier', 'battle', 'loveaffair', 'antiwar', 'powerrelations', 'mindandsoul', '3d',
                'SamWorthington', 'ZoeSaldana', 'SigourneyWeaver', 'JamesCameron']
In [ ]:
            H
```

Create new dataframe

```
new df = movies[['movie id', 'title', 'tags']]
 In [97]:
                 new df.head()
     Out[97]:
                                                              title
                     movie id
                                                                                                        tags
                  0
                        19995
                                                            Avatar
                                                                      [In, the, 22nd, century,, a, paraplegic, Marin...
                  1
                          285
                               Pirates of the Caribbean: At World's End
                                                                     [Captain, Barbossa,, long, believed, to, be, d...
                  2
                       206647
                                                           Spectre [A, cryptic, message, from, Bond's, past, send...
                                              The Dark Knight Rises
                  3
                        49026
                                                                      [Following, the, death, of, District, Attorney...
                  4
                        49529
                                                       John Carter
                                                                      [John, Carter, is, a, war-weary,, former, mili...
              ▶ | new df['tags'] = new df['tags'].apply(lambda x: " ".join(x))
 In [99]:
                 new df.head()
In [100]:
    Out[100]:
                     movie_id
                                                              title
                                                                                                          tags
                  0
                        19995
                                                            Avatar
                                                                      In the 22nd century, a paraplegic Marine is di...
                               Pirates of the Caribbean: At World's End
                                                                    Captain Barbossa, long believed to be dead, ha...
                  2
                       206647
                                                           Spectre A cryptic message from Bond's past sends him o...
                  3
                        49026
                                              The Dark Knight Rises
                                                                      Following the death of District Attorney Harve...
                        49529
                                                       John Carter
                                                                      John Carter is a war-weary, former military ca...
In [101]:
                new_df['tags'][0]
    Out[101]: 'In the 22nd century, a paraplegic Marine is dispatched to the moon Pandora on a unique mission, but becomes torn
                 between following orders and protecting an alien civilization. Action Adventure Fantasy ScienceFiction culturecla
                 sh future spacewar spacecolony society spacetravel futuristic romance space alien tribe alienplanet cgi marine so
```

ldier battle loveaffair antiwar powerrelations mindandsoul 3d SamWorthington ZoeSaldana SigourneyWeaver JamesCame

ron'

In [103]: ► new_df

Out[103]:

	movie_id	title	tags
	0 19995	Avatar	in the 22nd century, a paraplegic marine is di
	1 285	Pirates of the Caribbean: At World's End	captain barbossa, long believed to be dead, ha
	2 206647	Spectre	a cryptic message from bond's past sends him o
	3 49026	The Dark Knight Rises	following the death of district attorney harve
	4 49529	John Carter	john carter is a war-weary, former military ca
480	9367	El Mariachi	el mariachi just wants to play his guitar and
480	72766	Newlyweds	a newlywed couple's honeymoon is upended by th
480	23 1617	Signed, Sealed, Delivered	"signed, sealed, delivered" introduces a dedic
480	126186	Shanghai Calling	when ambitious new york attorney sam is sent t
480	25975	My Date with Drew	ever since the second grade when he first saw

4806 rows × 3 columns

In [105]: ▶ new_df.tags[0]

Out[105]: 'in the 22nd century, a paraplegic marine is dispatched to the moon pandora on a unique mission, but becomes torn between following orders and protecting an alien civilization. action adventure fantasy sciencefiction culturecla sh future spacewar spacecolony society spacetravel futuristic romance space alien tribe alienplanet cgi marine so ldier battle loveaffair antiwar powerrelations mindandsoul 3d samworthington zoesaldana sigourneyweaver jamescame ron'

Now will work using some NLP tasks

Out[128]: 'in the 22nd century, a parapleg marin is dispatch to the moon pandora on a uniqu mission, but becom torn between follow order and protect an alien civilization. action adventur fantasi sciencefict cultureclash futur spacewar s pacecoloni societi spacetravel futurist romanc space alien tribe alienplanet cgi marin soldier battl loveaffair a ntiwar powerrel mindandsoul 3d samworthington zoesaldana sigourneyweav jamescameron'

Convert text into vectors using CountVectorizer

```
▶ | from sklearn.feature extraction.text import CountVectorizer
In [111]:

  | cv = CountVectorizer(max_features=5000, stop words='english')

In [129]:
In [130]: N vectors = cv.fit transform(new df['tags']).toarray()
              vectors
   Out[130]: array([[0, 0, 0, ..., 0, 0, 0],
                     [0, 0, 0, \ldots, 0, 0, 0],
                     [0, 0, 0, ..., 0, 0, 0],
                     [0, 0, 0, \ldots, 0, 0, 0],
                     [0, 0, 0, ..., 0, 0, 0],
                     [0, 0, 0, ..., 0, 0, 0]], dtype=int64)
In [131]:
           vectors[0]
   Out[131]: array([0, 0, 0, ..., 0, 0, 0], dtype=int64)

    len(vectors)

In [132]:
   Out[132]: 4806
```

```
N cv.get feature names out()
In [133]:
   Out[133]: array(['000', '007', '10', ..., 'zone', 'zoo', 'zooeydeschanel'],
                    dtype=object)
           ▶ len(cv.get_feature_names_out())
In [134]:
   Out[134]: 5000
           for i in cv.get_feature_names_out():
In [135]:
                  print(i)
              ador
              adrienbrodi
              adult
              adultanim
              adulteri
              adulthood
              advanc
              adventur
              adventure
              adventures
              advertis
              advic
              advis
              affair
              affect
              afghanistan
              africa
              african
              africanamerican
              aftercreditecting
 In [ ]:
```

Cosine Similarity

```
In [136]: ▶ from sklearn.metrics.pairwise import cosine similarity
In [137]: ▶ | similarity = cosine similarity(vectors)
In [138]: ▶ similarity
  Out[138]: array([[1.
                           , 0.08346223, 0.0860309 , ..., 0.04499213, 0.
                  0.
                  [0.08346223, 1. , 0.06063391, ..., 0.02378257, 0.
                  0.02615329],
                  [0.0860309, 0.06063391, 1., 0.02451452, 0.
                  0.
                  [0.04499213, 0.02378257, 0.02451452, ..., 1. , 0.03962144,
                  0.04229549],
                  [0., 0., 0., 0.03962144, 1.
                  0.08714204],
                          , 0.02615329, 0. , ..., 0.04229549, 0.08714204,
                  [0.
                           11)
                  1.
In [139]: ▶ len(similarity)
  Out[139]: 4806
Out[140]: array([1.
                         , 0.08346223, 0.0860309 , ..., 0.04499213, 0.
                         1)
```

```
    for i in similarity[0]:

In [141]:
                  print(i)
              0.08346223261119858
              0.08603090020146065
              0.0734718358370645
              0.1892994097121204
              0.10838874619051501
              0.04024218182927669
              0.14673479641335554
              0.05923488777590923
              0.0967301666813349
              0.10259783520851541
              0.09464970485606021
              0.09037128496931669
              0.04499212706658476
              0.12824729401064427
              0.06282808624375433
              0.07894736842105264
              0.13977653617040256
              0.09493290614465533
              0 0000043004704500
In [142]:
           new df[new df['title'] == 'Avatar']
   Out[142]:
                 movie_id
                           title
                                                              tags
               0
                    19995 Avatar in the 22nd century, a parapleg marin is dispa...
In [143]:
           new_df[new_df['title'] == 'Avatar'].index
   Out[143]: Index([0], dtype='int64')
In [144]:
           new_df[new_df['title'] == 'Avatar'].index[0]
   Out[144]: 0
```

Finding the distances between vectors

```
    distances = similarity[0]

In [147]:
               distances
   Out[147]: array([1.
                                 , 0.08346223, 0.0860309 , ..., 0.04499213, 0.
                                 1)
In [148]:

■ sorted(similarity[0])

   Out[148]: [0.0,
                0.0,
                0.0,
                0.0,
                0.0,
                0.0,
                0.0,
                0.0,
                0.0,
                0.0,
                0.0,
                0.0,
                0.0,
                0.0,
                0.0,
                0.0,
                0.0,
                0.0,
                0.0,
```

```
■ sorted(similarity[0])[-10:]

In [151]:
   Out[151]: [0.23174488732966073,
             0.23179316248638276,
             0.24455799402225925,
             0.24511108480187255,
             0.25038669783359574,
             0.255608593705383,
             0.2605130246476754,
             0.26901379342448517,
             0.28676966733820225,
             In [150]:
   0.28676966733820225,
             0.26901379342448517,
             0.2605130246476754,
             0.255608593705383,
             0.25038669783359574,
             0.24511108480187255,
             0.24455799402225925,
             0.23179316248638276,
             0.23174488732966073,
             0.2278389747471728,
             0.2252817784447915,
             0.22269966704152225,
             0.21853668936906193,
             0.21239769762143662,
             0.2108663315950723,
             0.2105263157894737,
             0.20443988269091456,
             0.20437977982832192,
             0 20205070426402276
```

```
▶ list(enumerate(similarity[0]))
In [154]:
   Out[154]: [(0, 1.0000000000000000),
               (1, 0.08346223261119858),
               (2, 0.08603090020146065),
               (3, 0.0734718358370645),
               (4, 0.1892994097121204),
               (5, 0.10838874619051501),
               (6, 0.04024218182927669),
               (7, 0.14673479641335554),
               (8, 0.05923488777590923),
               (9, 0.0967301666813349),
               (10, 0.10259783520851541),
               (11, 0.09464970485606021),
               (12, 0.09037128496931669),
               (13, 0.04499212706658476),
               (14, 0.12824729401064427),
               (15, 0.06282808624375433),
               (16, 0.07894736842105264),
               (17, 0.13977653617040256),
               (18, 0.09493290614465533),
```

```
▶ | sorted(list(enumerate(similarity[0])), reverse=True)
In [157]:
   Out[157]: [(4805, 0.0),
               (4804, 0.0),
                (4803, 0.04499212706658476),
                (4802, 0.046829290579084706),
                (4801, 0.019252140716412975),
                (4800, 0.0),
                (4799, 0.052631578947368425),
                (4798, 0.04223886030955117),
                (4797, 0.0),
                (4796, 0.0),
                (4795, 0.0),
               (4794, 0.0),
                (4793, 0.05407380704358751),
                (4792, 0.0),
                (4791, 0.0),
               (4790, 0.0582716546748065),
               (4789, 0.060833032924035954),
                (4788, 0.0),
                (4787, 0.019117977822546817),
           # sorted(list(enumerate(similarity[0])), reverse=True)[0][1]
In [161]:
```

```
▶ | sorted(list(enumerate(similarity[0])), reverse=True, key=lambda x: x[1])
In [163]:
   Out[163]: [(0, 1.0000000000000000),
               (1216, 0.28676966733820225),
                (2409, 0.26901379342448517),
                (3730, 0.2605130246476754),
               (507, 0.255608593705383),
               (539, 0.25038669783359574),
               (582, 0.24511108480187255),
               (1204, 0.24455799402225925),
               (1194, 0.23179316248638276),
               (778, 0.23174488732966073),
                (4048, 0.2278389747471728),
               (1920, 0.2252817784447915),
                (61, 0.22269966704152225),
               (2786, 0.21853668936906193),
               (172, 0.21239769762143662),
               (972, 0.2108663315950723),
               (322, 0.2105263157894737),
               (2333, 0.20443988269091456),
                (3608, 0.20437977982832192),
           ▶ | sorted(list(enumerate(similarity[0])), reverse=True, key=lambda x: x[1])[1:6]
In [165]:
   Out[165]: [(1216, 0.28676966733820225),
               (2409, 0.26901379342448517),
               (3730, 0.2605130246476754),
               (507, 0.255608593705383),
               (539, 0.25038669783359574)]
           ▶ new df.iloc[1216]
In [166]:
   Out[166]: movie id
                                                                          440
              title
                                                 Aliens vs Predator: Requiem
                          a sequel to 2004' alien vs. predator, the icon...
              tags
              Name: 1216, dtype: object
```

```
    new df.iloc[2409]

In [168]:
   Out[168]: movie id
                                                                          679
              title
                                                                      Aliens
                          when ripley' lifepod is found by a salvag crew...
              tags
              Name: 2409, dtype: object
           ▶ new df.iloc[3730]
In [169]:
   Out[169]: movie id
                                                                      270938
                                                               Falcon Rising
              title
              tags
                          chapman is an ex-marin in brazil' slums, battl...
              Name: 3731, dtype: object
 In [ ]:
```

Now make a final function to return top 5 similar movies

```
recommend('Avatar')
In [187]:
              Aliens vs Predator: Requiem
              Aliens
              Falcon Rising
              Independence Day
              Titan A.E.

    recommend('Iron Man')

In [188]:
              Iron Man 3
              Iron Man 2
              Avengers: Age of Ultron
              The Avengers
              Captain America: Civil War
           recommend('Iron Man 3')
In [191]:
              Iron Man
              Iron Man 2
              Avengers: Age of Ultron
              Captain America: Civil War
              X-Men

    recommend('Batman Begins')

In [192]:
              The Dark Knight
              Batman
              Batman
              The Dark Knight Rises
              10th & Wolf
 In [ ]:
```

In []:	M	
In []:	M	
In [175]:	M	new_df.iloc[1216]
Out[175	5]:	movie_id 440 title Aliens vs Predator: Requiem tags a sequel to 2004' alien vs. predator, the icon Name: 1216, dtype: object
In [178]:	H	new_df.iloc[1216].title
Out[178	8]:	'Aliens vs Predator: Requiem'
In []:	H	
In []:	M	
In []:	H	
In []:	M	
In []:	M	
In []:	H	
In []:	M	

In []: N