

```
In [1]: import pandas as pd
```

```
In [5]: movies = pd.read_csv(r'C:\Users\DELL\Documents\chaitu DS\senapathi material\clas
```

```
In [19]: movies.head(2)
```

```
Out[19]:
```

	movieId	title	genres
0	1	Toy Story (1995)	Adventure Animation Children Comedy Fantasy
1	2	Jumanji (1995)	Adventure Children Fantasy

```
In [11]: print(type(movies))
```

```
<class 'pandas.core.frame.DataFrame'>
```

```
In [13]: movies.shape
```

```
Out[13]: (27278, 3)
```

```
In [37]: ratings=pd.read_csv(r'C:\Users\DELL\Documents\chaitu DS\senapathi material\class
```

```
In [40]: ratings.shape
```

```
Out[40]: (20000263, 4)
```

```
In [44]: tags=pd.read_csv(r'C:\Users\DELL\Documents\chaitu DS\senapathi material\class ro
```

```
In [46]: tags.shape
```

```
Out[46]: (465564, 4)
```

```
In [48]: tags.columns
```

```
Out[48]: Index(['userId', 'movieId', 'tag', 'timestamp'], dtype='object')
```

```
In [50]: ratings.columns
```

```
Out[50]: Index(['userId', 'movieId', 'rating', 'timestamp'], dtype='object')
```

```
In [52]: movies.columns
```

```
Out[52]: Index(['movieId', 'title', 'genres'], dtype='object')
```

```
In [54]: del ratings['timestamp']  
del tags['timestamp']
```

```
In [56]: tags.columns
```

```
Out[56]: Index(['userId', 'movieId', 'tag'], dtype='object')
```

```
In [58]: ratings.columns
```

```
Out[58]: Index(['userId', 'movieId', 'rating'], dtype='object')
```

```
In [62]: tags.head()
```

```
Out[62]:
```

	userId	movieId	tag
0	18	4141	Mark Waters
1	65	208	dark hero
2	65	353	dark hero
3	65	521	noir thriller
4	65	592	dark hero

```
In [64]: tags.iloc[0]
```

```
Out[64]:
```

userId	18
movieId	4141
tag	Mark Waters

Name: 0, dtype: object

```
In [66]: tags.iloc[2]
```

```
Out[66]:
```

userId	65
movieId	353
tag	dark hero

Name: 2, dtype: object

```
In [68]: row_0=tags.iloc[0]  
print(row_0)
```

```
userId      18  
movieId     4141  
tag         Mark Waters  
Name: 0, dtype: object
```

```
In [70]: row_0.index
```

```
Out[70]: Index(['userId', 'movieId', 'tag'], dtype='object')
```

```
In [78]: row_0['userId']
```

```
Out[78]: 18
```

```
In [104... 'rating' in row_0
```

```
Out[104... False
```

```
In [106... row_0.name
```

```
Out[106... 0
```

```
In [108... row_0 = row_0.rename('firstRow')  
row_0.name
```

```
Out[108... 'firstRow'
```

```
In [110... ratings.head()
```

Out[110...

	userId	movieId	rating
0	1	2	3.5
1	1	29	3.5
2	1	32	3.5
3	1	47	3.5
4	1	50	3.5

In [112... ratings['rating'].describe()

Out[112... count 2.000026e+07
mean 3.525529e+00
std 1.051989e+00
min 5.000000e-01
25% 3.000000e+00
50% 3.500000e+00
75% 4.000000e+00
max 5.000000e+00
Name: rating, dtype: float64

In [114... ratings.describe()

Out[114...

	userId	movieId	rating
count	2.000026e+07	2.000026e+07	2.000026e+07
mean	6.904587e+04	9.041567e+03	3.525529e+00
std	4.003863e+04	1.978948e+04	1.051989e+00
min	1.000000e+00	1.000000e+00	5.000000e-01
25%	3.439500e+04	9.020000e+02	3.000000e+00
50%	6.914100e+04	2.167000e+03	3.500000e+00
75%	1.036370e+05	4.770000e+03	4.000000e+00
max	1.384930e+05	1.312620e+05	5.000000e+00

In [116... ratings['rating'].mean()

Out[116... 3.5255285642993797

In [118... ratings.mean()

Out[118... userId 69045.872583
movieId 9041.567330
rating 3.525529
dtype: float64

In [120... ratings['rating'].min()

Out[120... 0.5

```
In [122... ratings.min()
```

```
Out[122...  userId      1.0  
      movieId    1.0  
      rating     0.5  
      dtype: float64
```

```
In [124... ratings['rating'].max()
```

```
Out[124... 5.0
```

```
In [128... ratings['rating'].std()
```

```
Out[128... 1.051988919275684
```

```
In [130... ratings['rating'].mode()
```

```
Out[130... 0      4.0  
      Name: rating, dtype: float64
```

```
In [132... ratings.corr()
```

```
Out[132...      userId  movieId  rating  
userId  1.000000 -0.000850  0.001175  
movieId -0.000850  1.000000  0.002606  
rating  0.001175  0.002606  1.000000
```

```
In [134... filter1 = ratings['rating'] > 10  
print(filter1)
```

```
0      False  
1      False  
2      False  
3      False  
4      False  
...  
20000258  False  
20000259  False  
20000260  False  
20000261  False  
20000262  False  
Name: rating, Length: 20000263, dtype: bool
```

```
In [136... filter1.any()
```

```
Out[136... False
```

```
In [138... filter1 = ratings['rating'] > 0  
print(filter1)
```

```
0          True
1          True
2          True
3          True
4          True
...
20000258   True
20000259   True
20000260   True
20000261   True
20000262   True
Name: rating, Length: 20000263, dtype: bool
```

```
In [140... movies.shape
```

```
Out[140... (27278, 3)
```

```
In [142... movies.isnull().any().any()
```

```
Out[142... False
```

```
In [144... ratings.shape
```

```
Out[144... (20000263, 3)
```

```
In [146... ratings.isnull().any().any()
```

```
Out[146... False
```

```
In [148... tags.shape
```

```
Out[148... (465564, 3)
```

```
In [150... tags.isnull().any().any()
```

```
Out[150... True
```

```
In [152... tags=tags.dropna()
```

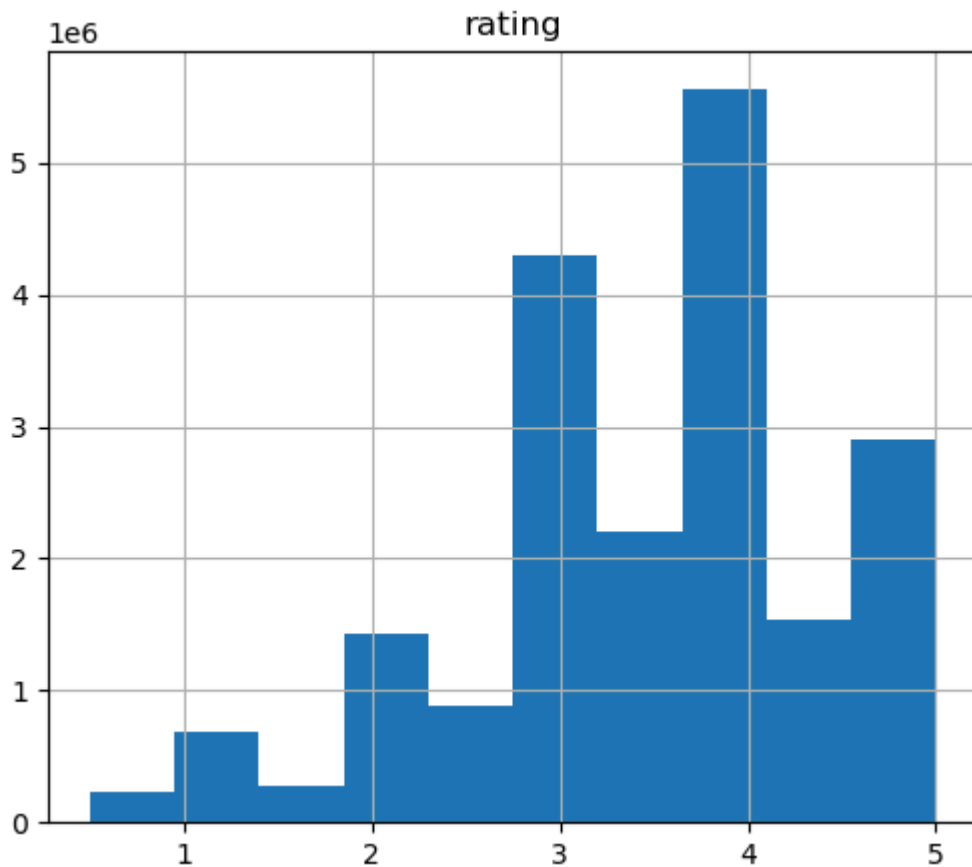
```
In [162... tags.isnull().any().any()
```

```
Out[162... False
```

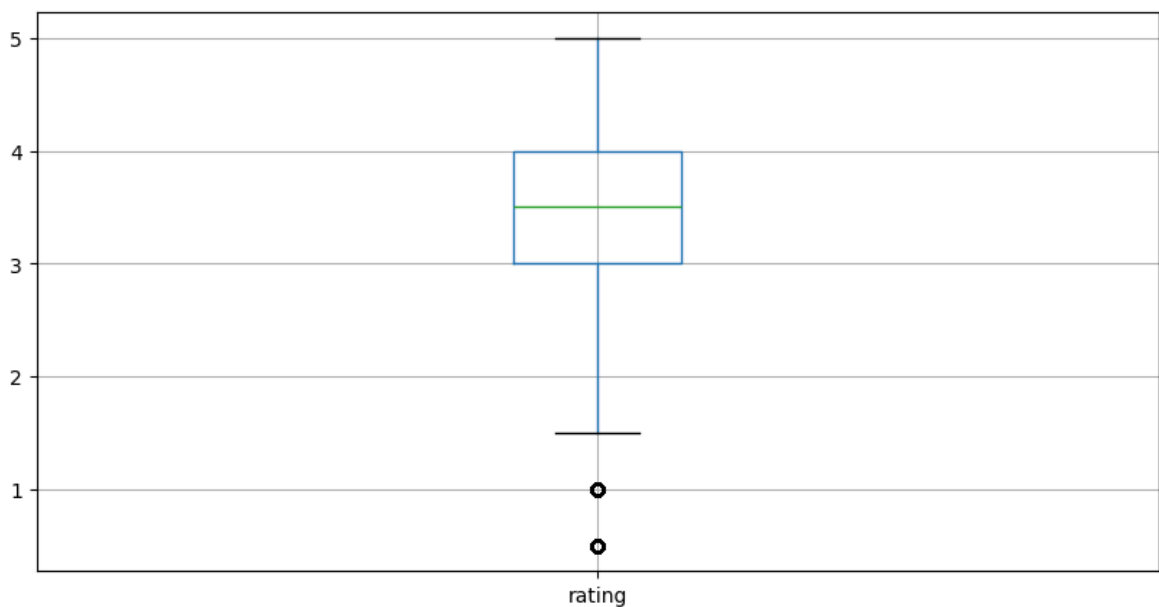
```
In [164... tags.shape
```

```
Out[164... (465548, 3)
```

```
In [176... import matplotlib.pyplot as plt
%matplotlib inline
ratings.hist(column='rating', figsize=(6,5))
plt.show()
```



```
In [198... ratings.boxplot(column='rating', figsize=(10,5))
plt.show()
```



```
In [186... tags['tag'].head()
```

```
Out[186... 0    Mark Waters
1    dark hero
2    dark hero
3    noir thriller
4    dark hero
Name: tag, dtype: object
```

```
In [188... movies[['title', 'genres']].head()
```

Out[188...

	title	genres
0	Toy Story (1995)	Adventure Animation Children Comedy Fantasy
1	Jumanji (1995)	Adventure Children Fantasy
2	Grumpier Old Men (1995)	Comedy Romance
3	Waiting to Exhale (1995)	Comedy Drama Romance
4	Father of the Bride Part II (1995)	Comedy

In [190...

```
ratings[-10:]
```

Out[190...

	userId	movieId	rating
20000253	138493	60816	4.5
20000254	138493	61160	4.0
20000255	138493	65682	4.5
20000256	138493	66762	4.5
20000257	138493	68319	4.5
20000258	138493	68954	4.5
20000259	138493	69526	4.5
20000260	138493	69644	3.0
20000261	138493	70286	5.0
20000262	138493	71619	2.5

In [192...

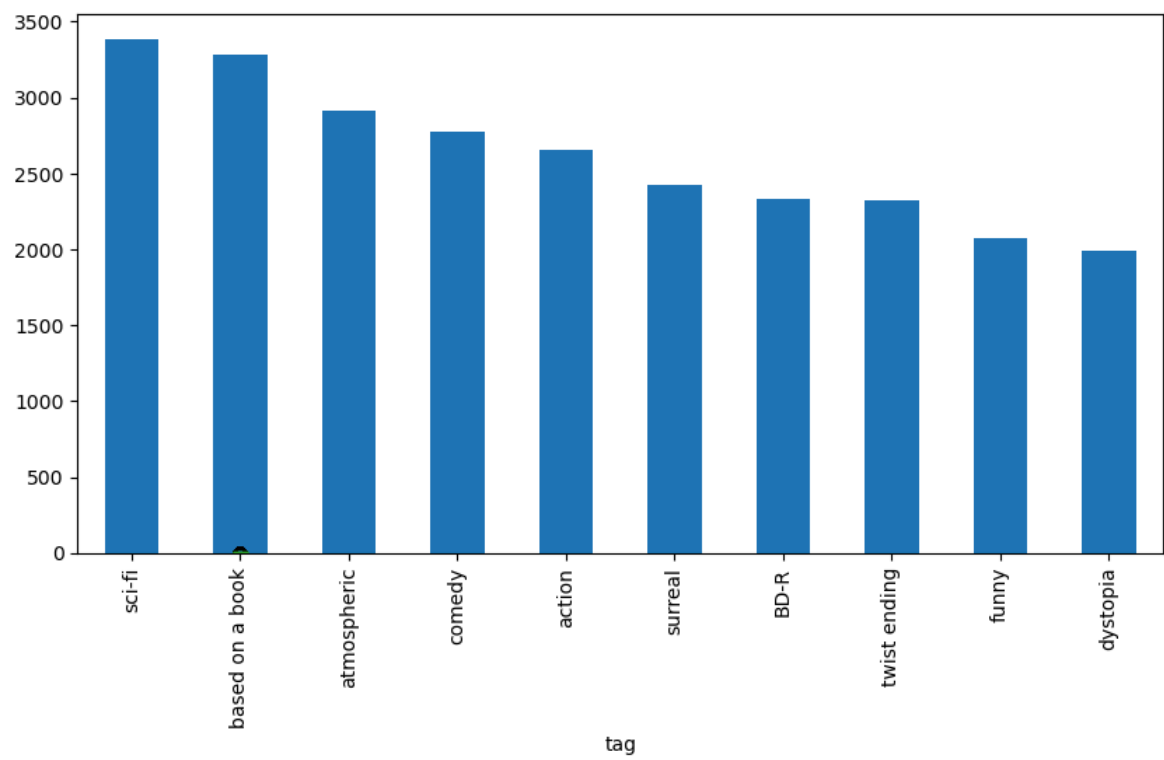
```
tag_counts = tags['tag'].value_counts()
tag_counts[-10:]
```

Out[192...

```
tag
missing child          1
Ron Moore              1
Citizen Kane          1
mullet                1
biker gang            1
Paul Adelstein        1
the wig               1
killer fish           1
genetically modified monsters 1
topless scene         1
Name: count, dtype: int64
```

In [196...

```
tag_counts[:10].plot(kind='bar', figsize=(10,5))
plt.show()
```



```
In [200... is_highly_rated = ratings['rating'] >= 5.0  
ratings[is_highly_rated][30:50]
```


Out[200...

	userId	movieId	rating
239	3	50	5.0
242	3	175	5.0
244	3	223	5.0
245	3	260	5.0
246	3	316	5.0
247	3	318	5.0
248	3	329	5.0
252	3	457	5.0
253	3	480	5.0
254	3	490	5.0
256	3	541	5.0
258	3	593	5.0
263	3	858	5.0
264	3	904	5.0
267	3	924	5.0
268	3	953	5.0
271	3	1060	5.0
272	3	1073	5.0
275	3	1084	5.0
276	3	1089	5.0

In [209...

```
is_action= movies['genres'].str.contains('Action')
movies[is_action][5:15]
```

Out[209...

	movieId	title	genres
22	23	Assassins (1995)	Action Crime Thriller
41	42	Dead Presidents (1995)	Action Crime Drama
43	44	Mortal Kombat (1995)	Action Adventure Fantasy
50	51	Guardian Angel (1994)	Action Drama Thriller
65	66	Lawnmower Man 2: Beyond Cyberspace (1996)	Action Sci-Fi Thriller
69	70	From Dusk Till Dawn (1996)	Action Comedy Horror Thriller
70	71	Fair Game (1995)	Action
75	76	Screamers (1995)	Action Sci-Fi Thriller
77	78	Crossing Guard, The (1995)	Action Crime Drama Thriller
85	86	White Squall (1996)	Action Adventure Drama

In [211...

```
movies[is_action].head(15)
```

Out[211...

	movieId	title	genres
5	6	Heat (1995)	Action Crime Thriller
8	9	Sudden Death (1995)	Action
9	10	GoldenEye (1995)	Action Adventure Thriller
14	15	Cutthroat Island (1995)	Action Adventure Romance
19	20	Money Train (1995)	Action Comedy Crime Drama Thriller
22	23	Assassins (1995)	Action Crime Thriller
41	42	Dead Presidents (1995)	Action Crime Drama
43	44	Mortal Kombat (1995)	Action Adventure Fantasy
50	51	Guardian Angel (1994)	Action Drama Thriller
65	66	Lawnmower Man 2: Beyond Cyberspace (1996)	Action Sci-Fi Thriller
69	70	From Dusk Till Dawn (1996)	Action Comedy Horror Thriller
70	71	Fair Game (1995)	Action
75	76	Screamers (1995)	Action Sci-Fi Thriller
77	78	Crossing Guard, The (1995)	Action Crime Drama Thriller
85	86	White Squall (1996)	Action Adventure Drama

In [213...

```
ratings_count = ratings[['movieId', 'rating']].groupby('rating').count()
ratings_count
```

Out[213...

movieId	
rating	
0.5	239125
1.0	680732
1.5	279252
2.0	1430997
2.5	883398
3.0	4291193
3.5	2200156
4.0	5561926
4.5	1534824
5.0	2898660

In [215...

```
average_rating = ratings[['movieId','rating']].groupby('movieId').mean()  
average_rating.head()
```

Out[215...

rating	
movieId	
1	3.921240
2	3.211977
3	3.151040
4	2.861393
5	3.064592

In [217...

```
movie_count = ratings[['movieId','rating']].groupby('movieId').count()  
movie_count.head()
```

Out[217...

rating	
movieId	
1	49695
2	22243
3	12735
4	2756
5	12161

In [219...

```
movie_count = ratings[['movieId','rating']].groupby('movieId').count()  
movie_count.tail()
```

Out[219...

rating	
movieId	
131254	1
131256	1
131258	1
131260	1
131262	1

In [231...

```
tags.head()
```

Out[231...

	userId	movieId	tag
0	18	4141	Mark Waters
1	65	208	dark hero
2	65	353	dark hero
3	65	521	noir thriller
4	65	592	dark hero

In [233...

```
movies.head()
```

Out[233...

	movieId	title	genres
0	1	Toy Story (1995)	Adventure Animation Children Comedy Fantasy
1	2	Jumanji (1995)	Adventure Children Fantasy
2	3	Grumpier Old Men (1995)	Comedy Romance
3	4	Waiting to Exhale (1995)	Comedy Drama Romance
4	5	Father of the Bride Part II (1995)	Comedy

In [235...

```
t = movies.merge(tags, on='movieId', how='inner')
t.head()
```

Out[235...

	movieId	title	genres	userId	tag
0	1	Toy Story (1995)	Adventure Animation Children Comedy Fantasy	1644	Watched
1	1	Toy Story (1995)	Adventure Animation Children Comedy Fantasy	1741	computer animation
2	1	Toy Story (1995)	Adventure Animation Children Comedy Fantasy	1741	Disney animated feature
3	1	Toy Story (1995)	Adventure Animation Children Comedy Fantasy	1741	Pixar animation
4	1	Toy Story (1995)	Adventure Animation Children Comedy Fantasy	1741	TÃ©a Leoni does not star in this movie

In [237...

```
avg_ratings= ratings.groupby('movieId', as_index=False).mean()
del avg_ratings['userId']
avg_ratings.head()
```

Out[237...

	movieId	rating
0	1	3.921240
1	2	3.211977
2	3	3.151040
3	4	2.861393
4	5	3.064592

In [239...

```
box_office = movies.merge(avg_ratings, on='movieId', how='inner')
box_office.tail()
```

Out[239...

	movieId	title	genres	rating
26739	131254	Kein Bund für's Leben (2007)	Comedy	4.0
26740	131256	Feuer, Eis & Dosenbier (2002)	Comedy	4.0
26741	131258	The Pirates (2014)	Adventure	2.5
26742	131260	Rentun Ruusu (2001)	(no genres listed)	3.0
26743	131262	Innocence (2014)	Adventure Fantasy Horror	4.0

In [241...

```
is_highly Rated = box_office['rating'] >= 4.0
box_office[is_highly Rated][-5:]
```

Out[241...

	movieid	title	genres	rating
26737	131250	No More School (2000)	Comedy	4.0
26738	131252	Forklift Driver Klaus: The First Day on the Jo...	Comedy Horror	4.0
26739	131254	Kein Bund für's Leben (2007)	Comedy	4.0
26740	131256	Feuer, Eis & Dosenbier (2002)	Comedy	4.0
26743	131262	Innocence (2014)	Adventure Fantasy Horror	4.0

In [243...

```
is_Adventure = box_office['genres'].str.contains('Adventure')
box_office[is_Adventure][:5]
```

Out[243...

	movieid	title	genres	rating
0	1	Toy Story (1995)	Adventure Animation Children Comedy Fantasy	3.921240
1	2	Jumanji (1995)	Adventure Children Fantasy	3.211977
7	8	Tom and Huck (1995)	Adventure Children	3.142049
9	10	GoldenEye (1995)	Action Adventure Thriller	3.430029
12	13	Balto (1995)	Adventure Animation Children	3.272416

In [245...

```
box_office[is_Adventure & is_highly Rated][-5:]
```

Out[245...

	movieid	title	genres	rating
26611	130586	Itinerary of a Spoiled Child (1988)	Adventure Drama	4.5
26655	130996	The Beautiful Story (1992)	Adventure Drama Fantasy	5.0
26667	131050	Stargate SG-1 Children of the Gods - Final Cut...	Adventure Sci-Fi Thriller	5.0
26736	131248	Brother Bear 2 (2006)	Adventure Animation Children Comedy Fantasy	4.0
26743	131262	Innocence (2014)	Adventure Fantasy Horror	4.0

In [249...

```
movies.head()
```

Out[249...

	movieId	title	genres
0	1	Toy Story (1995)	Adventure Animation Children Comedy Fantasy
1	2	Jumanji (1995)	Adventure Children Fantasy
2	3	Grumpier Old Men (1995)	Comedy Romance
3	4	Waiting to Exhale (1995)	Comedy Drama Romance
4	5	Father of the Bride Part II (1995)	Comedy

In [255...

```
movie_genres = movies['genres'].str.split('|', expand=True)
```

In [256...

```
movie_genres[:10]
```

Out[256...

	0	1	2	3	4	5	6	7	8	9
0	Adventure	Animation	Children	Comedy	Fantasy	None	None	None	None	None
1	Adventure	Children	Fantasy	None	None	None	None	None	None	None
2	Comedy	Romance	None	None	None	None	None	None	None	None
3	Comedy	Drama	Romance	None	None	None	None	None	None	None
4	Comedy	None	None	None	None	None	None	None	None	None
5	Action	Crime	Thriller	None	None	None	None	None	None	None
6	Comedy	Romance	None	None	None	None	None	None	None	None
7	Adventure	Children	None	None	None	None	None	None	None	None
8	Action	None	None	None	None	None	None	None	None	None
9	Action	Adventure	Thriller	None	None	None	None	None	None	None

In [263...

```
movie_genres['isComedy'] = movies['genres'].str.contains('Comedy')
```

In [265...

```
movie_genres[:10]
```

Out[265...

	0	1	2	3	4	5	6	7	8	9
0	Adventure	Animation	Children	Comedy	Fantasy	None	None	None	None	None
1	Adventure	Children	Fantasy	None	None	None	None	None	None	None
2	Comedy	Romance	None	None	None	None	None	None	None	None
3	Comedy	Drama	Romance	None	None	None	None	None	None	None
4	Comedy	None	None	None	None	None	None	None	None	None
5	Action	Crime	Thriller	None	None	None	None	None	None	None
6	Comedy	Romance	None	None	None	None	None	None	None	None
7	Adventure	Children	None	None	None	None	None	None	None	None
8	Action	None	None	None	None	None	None	None	None	None
9	Action	Adventure	Thriller	None	None	None	None	None	None	None

In [271...

```
movies['year'] = movies['title'].str.extract('.*\((.*)\).*', expand=True)
```

```
<>:1: SyntaxWarning: invalid escape sequence '\('
<>:1: SyntaxWarning: invalid escape sequence '\('
C:\Users\DELL\AppData\Local\Temp\ipykernel_57480\275227335.py:1: SyntaxWarning: i
nvalid escape sequence '\('
  movies['year'] = movies['title'].str.extract('.*\((.*)\).*', expand=True)
```

In [273...

```
movies.tail()
```

Out[273...

	movieId	title	genres	year
27273	131254	Kein Bund für's Leben (2007)	Comedy	2007
27274	131256	Feuer, Eis & Dosenbier (2002)	Comedy	2002
27275	131258	The Pirates (2014)	Adventure	2014
27276	131260	Rentun Ruusu (2001)	(no genres listed)	2001
27277	131262	Innocence (2014)	Adventure Fantasy Horror	2014

In [277...

```
tags = pd.read_csv(r'C:\Users\DELL\Documents\chaitu DS\senapathi material\class
```

In [281...

```
tags.dtypes
```

Out[281...

```
userId      int64
movieId     int64
tag         object
timestamp   object
dtype: object
```

In [283...

```
tags.head(5)
```


Out[283...

	userId	movieId	tag	timestamp
0	18	4141	Mark Waters	2009-04-24 18:19:40
1	65	208	dark hero	2013-05-10 01:41:18
2	65	353	dark hero	2013-05-10 01:41:19
3	65	521	noir thriller	2013-05-10 01:39:43
4	65	592	dark hero	2013-05-10 01:41:18

In [293...

```
average_rating = ratings[['movieId', 'rating']].groupby('movieId', as_index=False)
average_rating.tail()
```

Out[293...

	movieId	rating
26739	131254	4.0
26740	131256	4.0
26741	131258	2.5
26742	131260	3.0
26743	131262	4.0

In [297...

```
joined = movies.merge(average_rating, on='movieId', how='inner')
joined.head()
joined.corr()
```

```

-----
ValueError                                Traceback (most recent call last)
Cell In[297], line 3
      1 joined = movies.merge(average_rating, on='movieId', how='inner')
      2 joined.head()
----> 3 joined.corr()

File ~\anaconda3\Lib\site-packages\pandas\core\frame.py:11049, in DataFrame.corr
(self, method, min_periods, numeric_only)
    11047 cols = data.columns
    11048 idx = cols.copy()
> 11049 mat = data.to_numpy(dtype=float, na_value=np.nan, copy=False)
    11051 if method == "pearson":
    11052     correl = libalgos.nancorr(mat, minp=min_periods)

File ~\anaconda3\Lib\site-packages\pandas\core\frame.py:1993, in DataFrame.to_numpy(self, dtype, copy, na_value)
    1991 if dtype is not None:
    1992     dtype = np.dtype(dtype)
-> 1993 result = self._mgr.as_array(dtype=dtype, copy=copy, na_value=na_value)
    1994 if result.dtype is not dtype:
    1995     result = np.asarray(result, dtype=dtype)

File ~\anaconda3\Lib\site-packages\pandas\core\internals\managers.py:1694, in BlockManager.as_array(self, dtype, copy, na_value)
    1692     arr.flags.writeable = False
    1693 else:
-> 1694     arr = self._interleave(dtype=dtype, na_value=na_value)
    1695     # The underlying data was copied within _interleave, so no need
    1696     # to further copy if copy=True or setting na_value
    1698 if na_value is lib.no_default:

File ~\anaconda3\Lib\site-packages\pandas\core\internals\managers.py:1753, in BlockManager._interleave(self, dtype, na_value)
    1751     else:
    1752         arr = blk.get_values(dtype)
-> 1753     result[rl.indexer] = arr
    1754     itemmask[rl.indexer] = 1
    1756 if not itemmask.all():

ValueError: could not convert string to float: 'Toy Story (1995)'

```

In []: