

```
In [2]: import pandas as pd
movies = pd.read_csv(r'C:\Users\Ravi\Downloads\archive\movie.csv')
```

```
In [3]: movies.shape
```

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

```
In [4]: movies.head(5)
```

```
Out[4]:
```

	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 [5]: ratings = pd.read_csv(r'C:\Users\Ravi\Downloads\archive\rating.csv')
```

```
In [6]: ratings.shape
```

```
Out[6]: (20000263, 4)
```

```
In [7]: tags = pd.read_csv(r'C:\Users\Ravi\Downloads\archive>tag.csv')
```

```
In [8]: tags.shape
```

```
Out[8]: (465564, 4)
```

```
In [9]: print(movies.shape)
print(ratings.shape)
print(tags.shape)
```

```
(27278, 3)
```

```
(20000263, 4)
```

```
(465564, 4)
```

```
In [10]: print(movies.columns)
print(ratings.columns)
print(tags.columns)
```

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

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

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

```
In [11]: del ratings['timestamp']
del tags['timestamp']
```

```
In [12]: print(movies.columns)
print(ratings.columns)
print(tags.columns)
```

```
Index(['movieId', 'title', 'genres'], dtype='object')
Index(['userId', 'movieId', 'rating'], dtype='object')
Index(['userId', 'movieId', 'tag'], dtype='object')
```

```
In [13]: tags.head()
```

```
Out[13]:
```

	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 [18]: row_0=tags.iloc[0]
         type(row_0)
```

```
Out[18]: pandas.core.series.Series
```

```
In [ ]:
```

```
In [19]: print(row_0)
```

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

```
In [20]: row_0.index
```

```
Out[20]: Index(['userId', 'movieId', 'tag'], dtype='object')
```

```
In [21]: row_0['userId']
```

```
Out[21]: np.int64(18)
```

```
In [22]: 'rating' in row_0
```

```
Out[22]: False
```

```
In [23]: row_0.name
```

```
Out[23]: 0
```

```
In [24]: row_0 = row_0.rename('firstRow')
         row_0.name
```

```
Out[24]: 'firstRow'
```

```
In [25]: tags.head()
```

```
Out[25]:
```

	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 [27]: tags.index
```

```
Out[27]: RangeIndex(start=0, stop=465564, step=1)
```

```
In [28]: tags.columns
```

```
Out[28]: Index(['userId', 'movieId', 'tag'], dtype='object')
```

```
In [29]: tags.iloc[[0,11,500]]
```

```
Out[29]:
```

	userId	movieId	tag
0	18	4141	Mark Waters
11	65	1783	noir thriller
500	342	55908	entirely dialogue

```
In [31]: ratings['rating'].describe()
```

```
Out[31]: 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 [32]: ratings.describe()
```

Out[32]:

	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 [33]: ratings['rating'].mean()

Out[33]: np.float64(3.5255285642993797)

In [34]: ratings.mean()

Out[34]:

userId	69045.872583
movieId	9041.567330
rating	3.525529
dtype:	float64

In [35]: ratings['rating'].min()

Out[35]: 0.5

In [36]: ratings['rating'].max()

Out[36]: 5.0

In [37]: ratings['rating'].std()

Out[37]: 1.051988919275684

In [38]: ratings['rating'].mode()

Out[38]:

0	4.0
---	-----

Name: rating, dtype: float64

In [39]: ratings.corr()

Out[39]:

	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 [40]:

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

```
filter1.any()
```

```
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
```

```
Out[40]: np.False_
```

```
In [41]: filter2 = ratings['rating'] > 0
filter2.all()
```

```
Out[41]: np.True_
```

```
In [42]: movies.shape
```

```
Out[42]: (27278, 3)
```

```
In [43]: movies.isnull().any().any()
```

```
Out[43]: np.False_
```

```
In [44]: ratings.shape
```

```
Out[44]: (20000263, 3)
```

```
In [45]: ratings.isnull().any().any()
```

```
Out[45]: np.False_
```

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

```
Out[46]: (465564, 3)
```

```
In [47]: tags.isnull().any().any()
```

```
Out[47]: np.True_
```

```
In [48]: tags=tags.dropna()
```

```
In [49]: tags.isnull().any().any()
```

```
Out[49]: np.False_
```

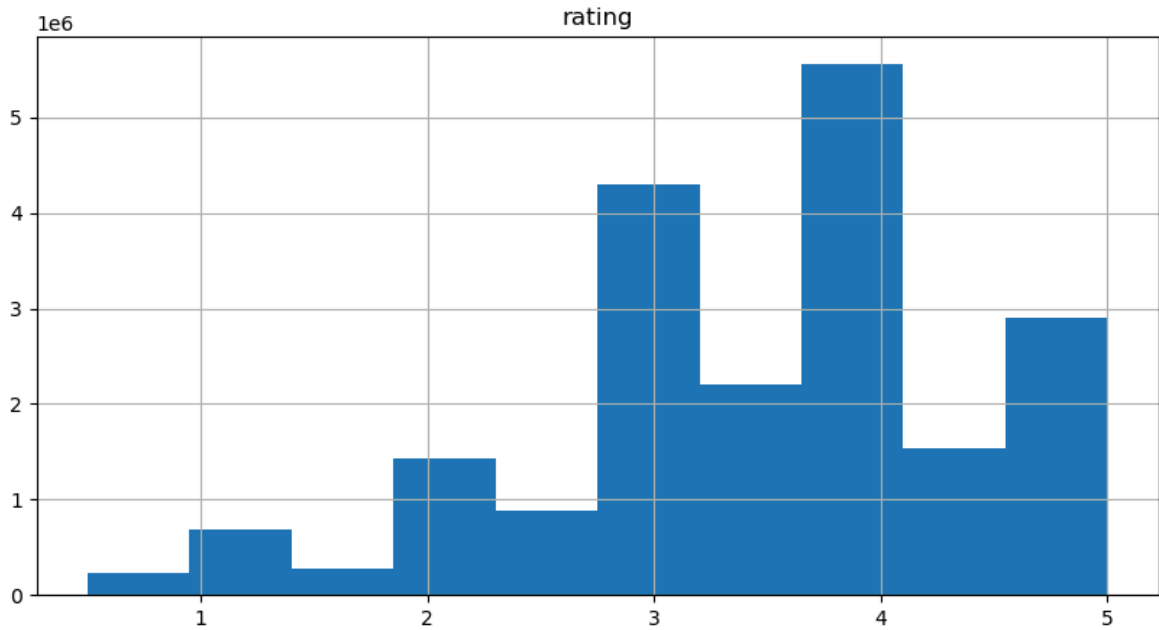
```
In [50]: tags.shape
```

```
Out[50]: (465548, 3)
```

```
In [51]: %matplotlib inline
```

```
ratings.hist(column='rating', figsize=(10,5))
```

```
Out[51]: array([[<Axes: title={'center': 'rating'}>]], dtype=object)
```



```
In [52]: tags['tag'].head()
```

```
Out[52]: 0    Mark Waters
1    dark hero
2    dark hero
3    noir thriller
4    dark hero
Name: tag, dtype: object
```

```
In [53]: movies[['title', 'genres']].head()
```

```
Out[53]:
```

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

Out[54]:

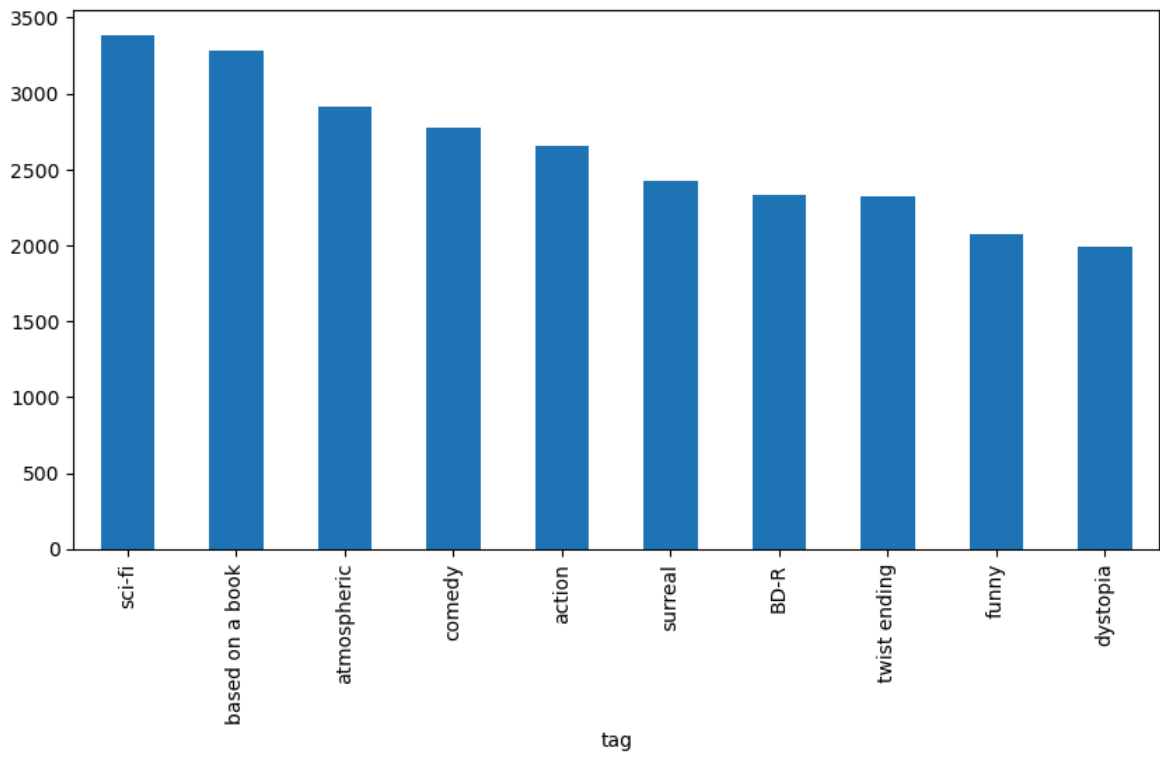
	userId	movieId	rating
20000253	138493	60816	4.5
20000254	138493	61160	4.0
20000255	138493	65682	4.5
20000256	138493	66762	4.5
20000257	138493	68319	4.5
20000258	138493	68954	4.5
20000259	138493	69526	4.5
20000260	138493	69644	3.0
20000261	138493	70286	5.0
20000262	138493	71619	2.5

```
In [55]: tag_counts = tags['tag'].value_counts()
tag_counts[-10:]
```

```
Out[55]: 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 [56]: tag_counts[:10].plot(kind='bar', figsize=(10,5))
```

```
Out[56]: <Axes: xlabel='tag'>
```



In []: