

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

```
In [2]: moive=pd.read_csv(r'C:\Users\soham\OneDrive\Desktop\movie.csv')
```

moive

```
In [6]: print(type(moive))
```

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

```
In [8]: moive.head(30)
```

Out[8]:

|    | 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                                      |
| 5  | 6       | Heat (1995)                                       | Action Crime Thriller                       |
| 6  | 7       | Sabrina (1995)                                    | Comedy Romance                              |
| 7  | 8       | Tom and Huck (1995)                               | Adventure Children                          |
| 8  | 9       | Sudden Death (1995)                               | Action                                      |
| 9  | 10      | GoldenEye (1995)                                  | Action Adventure Thriller                   |
| 10 | 11      | American President, The (1995)                    | Comedy Drama Romance                        |
| 11 | 12      | Dracula: Dead and Loving It (1995)                | Comedy Horror                               |
| 12 | 13      | Balto (1995)                                      | Adventure Animation Children                |
| 13 | 14      | Nixon (1995)                                      | Drama                                       |
| 14 | 15      | Cutthroat Island (1995)                           | Action Adventure Romance                    |
| 15 | 16      | Casino (1995)                                     | Crime Drama                                 |
| 16 | 17      | Sense and Sensibility (1995)                      | Drama Romance                               |
| 17 | 18      | Four Rooms (1995)                                 | Comedy                                      |
| 18 | 19      | Ace Ventura: When Nature Calls (1995)             | Comedy                                      |
| 19 | 20      | Money Train (1995)                                | Action Comedy Crime Drama Thriller          |
| 20 | 21      | Get Shorty (1995)                                 | Comedy Crime Thriller                       |
| 21 | 22      | Copycat (1995)                                    | Crime Drama Horror Mystery Thriller         |
| 22 | 23      | Assassins (1995)                                  | Action Crime Thriller                       |
| 23 | 24      | Powder (1995)                                     | Drama Sci-Fi                                |
| 24 | 25      | Leaving Las Vegas (1995)                          | Drama Romance                               |
| 25 | 26      | Othello (1995)                                    | Drama                                       |
| 26 | 27      | Now and Then (1995)                               | Children Drama                              |
| 27 | 28      | Persuasion (1995)                                 | Drama Romance                               |
| 28 | 29      | City of Lost Children, The (Cit  des enfants p... | Adventure Drama Fantasy Mystery Sci-Fi      |

|    | movielid | title   | genres      |
|----|----------|---|-------------|
| 29 | 30       | Shanghai Triad (Yao a yao yao dao waipo qiao) ... | Crime Drama |

```
In [10]: tags = pd.read_csv(r'C:\Users\soham\OneDrive\Desktop\28th- Kaggle Workshop\imbd dat
```

```
In [12]: tags.head()
```

```
Out[12]:
```

|   | userId | movielid | 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 [17]: rating = pd.read_csv(r'C:\Users\soham\OneDrive\Desktop\28th- Kaggle Workshop\imbd d
```

```
In [19]: #ratings = pd.read_csv(r'C:\Users\soham\OneDrive\Desktop\28th- Kaggle Workshop\imbd
#ratings.head()
```

```
Out[19]:
```

|   | userId | movielid | rating | timestamp           |
|---|--------|----------|--------|---------------------|
| 0 | 1      | 2        | 3.5    | 2005-04-02 23:53:47 |
| 1 | 1      | 29       | 3.5    | 2005-04-02 23:31:16 |
| 2 | 1      | 32       | 3.5    | 2005-04-02 23:33:39 |
| 3 | 1      | 47       | 3.5    | 2005-04-02 23:32:07 |
| 4 | 1      | 50       | 3.5    | 2005-04-02 23:29:40 |

```
In [21]: rating.head()
```

```
Out[21]:
```

|   | userId | movielid | rating | timestamp           |
|---|--------|----------|--------|---------------------|
| 0 | 1      | 2        | 3.5    | 2005-04-02 23:53:47 |
| 1 | 1      | 29       | 3.5    | 2005-04-02 23:31:16 |
| 2 | 1      | 32       | 3.5    | 2005-04-02 23:33:39 |
| 3 | 1      | 47       | 3.5    | 2005-04-02 23:32:07 |
| 4 | 1      | 50       | 3.5    | 2005-04-02 23:29:40 |

```
In [24]: del rating['timestamp']      #For current analysis, we will remove timestamp
del tags['timestamp']
```

```
In [26]: row = tags.iloc[0]    #Series  
type(row)
```

```
Out[26]: pandas.core.series.Series
```

```
In [28]: print(row)
```

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

```
In [30]: row.index
```

```
Out[30]: Index(['userId', 'movieId', 'tag'], dtype='object')
```

```
In [32]: row['userId']
```

```
Out[32]: 18
```

```
In [34]: 'rating' in row
```

```
Out[34]: False
```

```
In [36]: row.name
```

```
Out[36]: 0
```

```
In [38]: row = row.rename('firstRow')  
row.name
```

```
Out[38]: 'firstRow'
```

```
In [42]: #dataframe#
```

```
In [46]: tags.head()
```

```
Out[46]:
```

|   | 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 [48]: tags.index
```

```
Out[48]: RangeIndex(start=0, stop=465564, step=1)
```

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

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

```
In [56]: tags.iloc [[0,11,500]]
```

```
Out[56]:
```

|            | userId | movieId | tag               |
|------------|--------|---------|-------------------|
| <b>0</b>   | 18     | 4141    | Mark Waters       |
| <b>11</b>  | 65     | 1783    | noir thriller     |
| <b>500</b> | 342    | 55908   | entirely dialogue |

## Descriptive Statistics

### Let's look how the ratings are distributed!

```
In [60]: rating['rating'].describe()
```

```
Out[60]: 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 [62]: rating.describe()
```

```
Out[62]:
```

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

```
In [70]: rating['rating'].mean()
```

Out[70]: 3.5255285642993797

In [72]: `rating.mean()`

Out[72]:

|         |              |
|---------|--------------|
| userId  | 69045.872583 |
| movieId | 9041.567330  |
| rating  | 3.525529     |
| dtype:  | float64      |

In [74]: `rating['rating'].min()`

Out[74]: 0.5

In [76]: `rating['rating'].std()`

Out[76]: 1.051988919275684

In [78]: `rating['rating'].mode()`

Out[78]:

|   |     |
|---|-----|
| 0 | 4.0 |
|---|-----|

Name: rating, dtype: float64

In [80]: `rating.corr()`

Out[80]:

|         | 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 [86]:

```

filter = rating['rating'] > 10
print(filter)
filter.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[86]: False

In [84]: `filter.any()`

Out[84]: False

```
In [88]: filter1 = rating['rating'] > 0  
filter1.all()
```

```
Out[88]: True
```

## Data Cleaning: Handling Missing Data

```
In [96]: moive.shape
```

```
Out[96]: (27278, 3)
```

```
In [100... moive.isnull().any()
```

```
Out[100... movieId    False  
title         False  
genres        False  
dtype: bool
```

```
In [102... moive.isnull().any().any()
```

```
Out[102... False
```

```
In [104... tags.shape
```

```
Out[104... (465564, 3)
```

```
In [106... tags.isnull().any().any()
```

```
Out[106... True
```

```
In [ ]: #We have some tags which are NULL.
```

```
In [114... tags=tags.dropna()
```

```
In [116... tags.isnull().any().any()
```

```
Out[116... False
```

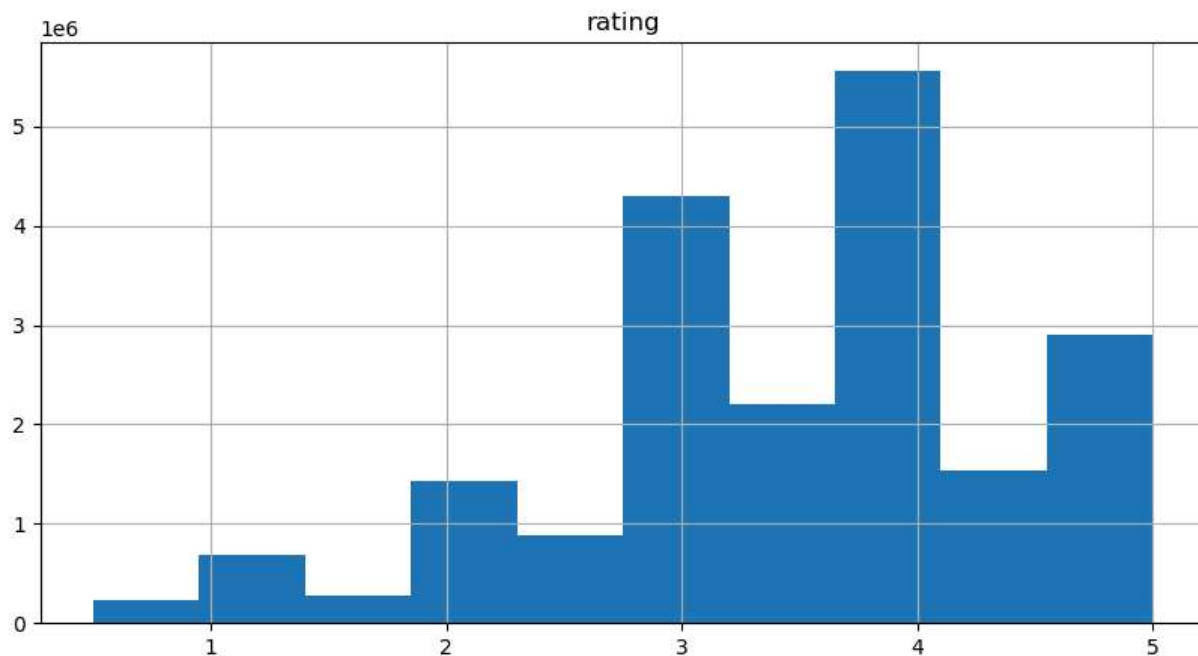
```
In [118... tags.shape
```

```
Out[118... (465548, 3)
```

```
In [ ]: #Data visulazation
```

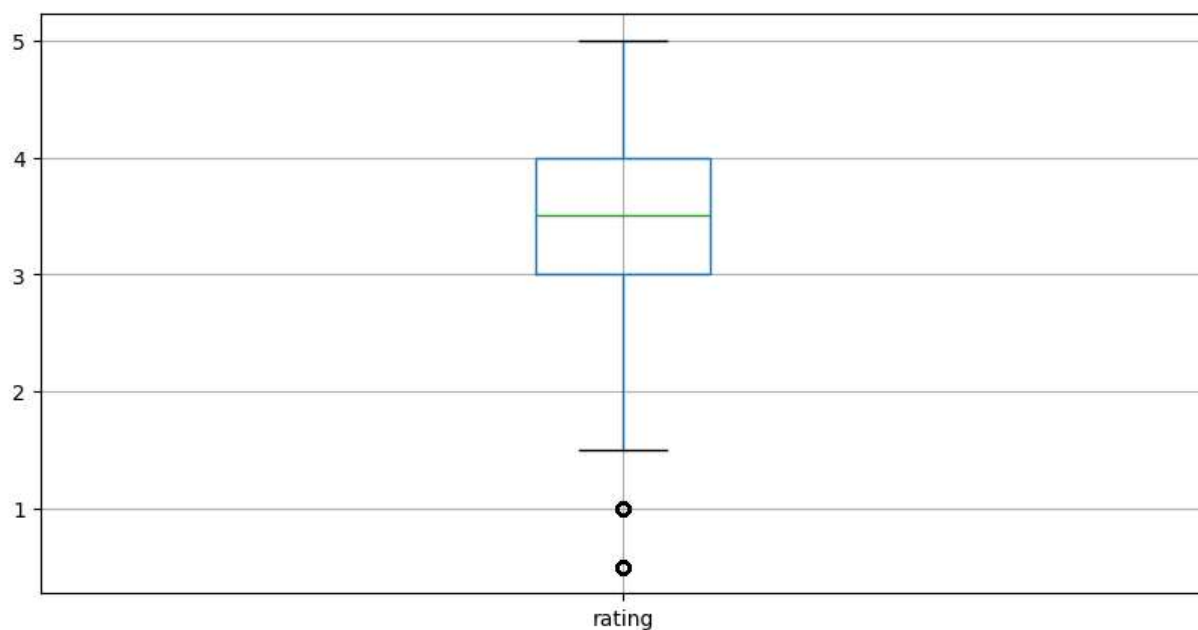
```
In [122... %matplotlib inline  
  
rating.hist(column='rating',figsize=(10,5))
```

```
Out[122... array([[<Axes: title={'center': 'rating'}>]], dtype=object)
```



```
In [128...] rating.boxplot(column='rating', figsize=(10,5))
```

```
Out[128...] <Axes: >
```



```
In [ ]: #slicing Out Columns
```

```
In [132...] tags['tag'].head()
```

```
Out[132...] 0    Mark Waters  
            1    dark hero  
            2    dark hero  
            3    noir thriller  
            4    dark hero  
            Name: tag, dtype: object
```



In [134... `moive[['title', 'genres']].head()`

Out[134... 

|   | 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 [136... `rating[-10:]`

Out[136... 

|          | 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 [140... `tags_count = tags['tag'].value_counts()  
tags_count[-10:]`

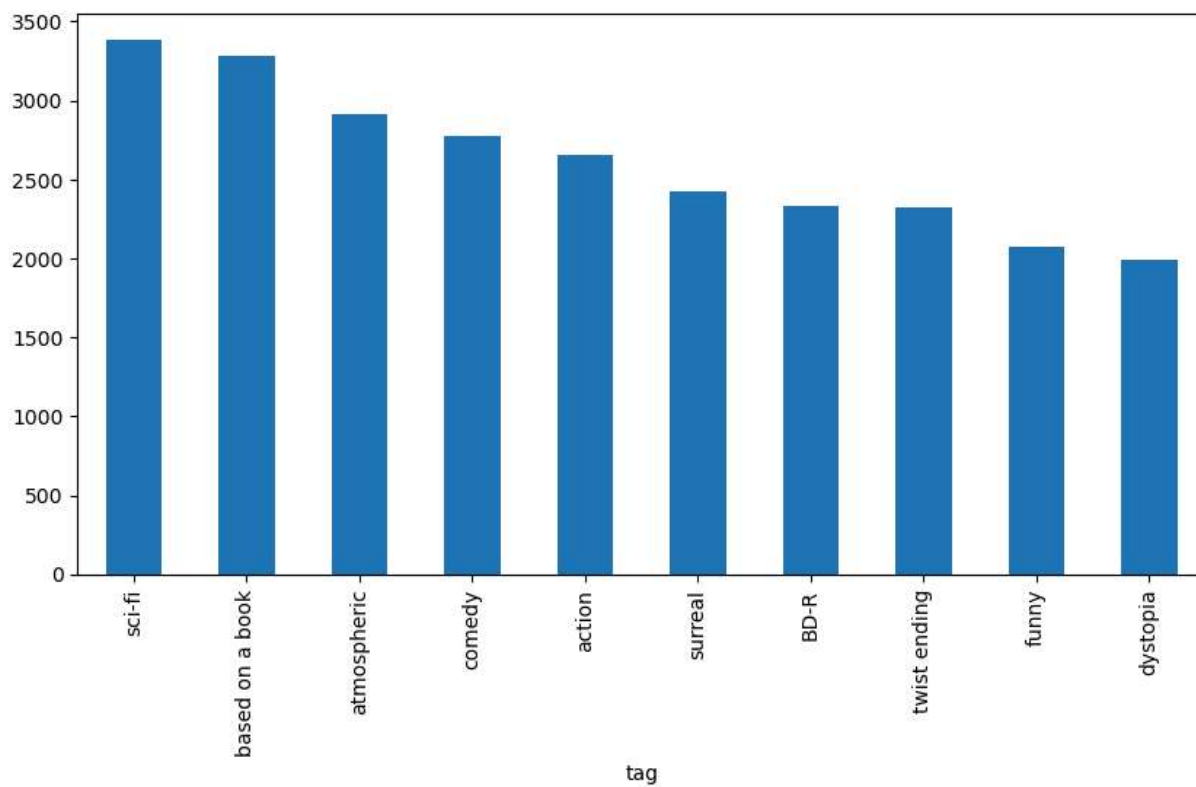
Out[140... 

| 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 [142... `tags_count[:10].plot(kind='bar', figsize=(10, 5))`

Out[142... &lt;Axes: xlabel='tag'&gt;



In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]: