

In []:

IMDB movie data analysis using pandas

```
In [ ]: numpy - image to array (1d, 2d, 3d)
        pytorch -- python + torch
        pytorch - image to tensor
                1d tensor , 2d tensor, 3d tensor , nd tensor

        kaggle.com
        world - (kaggle.com )

        kaggle.com -- data anlysis, data scientist, business analyst
        kaggle.com -- multiple dataset
                banking, healthcare, insure, supplyIF
```

```
In [1]: #Import Libraries
        import pandas as pd
```

```
In [3]: # read the dataset
        movies = pd.read_csv(r'E:\Data Science & AI\Dataset files\archive\movie.csv')
```

```
In [5]: ratings = pd.read_csv(r'E:\Data Science & AI\Dataset files\archive\rating.csv')
```

```
In [6]: tags = pd.read_csv(r'E:\Data Science & AI\Dataset files\archive>tag.csv')
```

```
In [7]: print(type(movies))
        movies.head(20)
```

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

Out[7]:	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

```
In [8]: print(type(ratings))
ratings.head(10)
```

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

Out[8]:

	userId	movieId	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
5	1	112	3.5	2004-09-10 03:09:00
6	1	151	4.0	2004-09-10 03:08:54
7	1	223	4.0	2005-04-02 23:46:13
8	1	253	4.0	2005-04-02 23:35:40
9	1	260	4.0	2005-04-02 23:33:46

In [9]: `print(type(tags))`
`tags.head(10)`

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

Out[9]:

	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
5	65	668	bollywood	2013-05-10 01:37:56
6	65	898	screwball comedy	2013-05-10 01:42:40
7	65	1248	noir thriller	2013-05-10 01:39:43
8	65	1391	mars	2013-05-10 01:40:55
9	65	1617	neo-noir	2013-05-10 01:43:37

In [10]: `tags.head()`

Out[10]:

	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 [11]: ratings.head()
```

```
Out[11]:
```

	userId	movieId	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 [12]: #For current analysis, we will remove timestamp  
del ratings['timestamp']  
del tags['timestamp']
```

Data Structures:

Series

```
In [14]: row_0 = tags.iloc[0]  
type(row_0)
```

```
Out[14]: pandas.core.series.Series
```

```
In [15]: print(row_0)  
  
userId          18  
movieId         4141  
tag            Mark Waters  
Name: 0, dtype: object
```

```
In [16]: row_0.index
```

```
Out[16]: Index(['userId', 'movieId', 'tag'], dtype='object')
```

```
In [17]: row_0['userId']
```

```
Out[17]: 18
```

```
In [18]: 'rating' in row_0
```

```
Out[18]: False
```

```
In [19]: row_0.name
```

```
Out[19]: 0
```

```
In [20]: row_0 = row_0.rename('firstRow')  
row_0.name
```

```
Out[20]: 'firstRow'
```

```
In [21]: #DataFrames
tags.head()
```

```
Out[21]:
```

	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 [22]: tags.index
```

```
Out[22]: RangeIndex(start=0, stop=465564, step=1)
```

```
In [23]: tags.columns
```

```
Out[23]: Index(['userId', 'movieId', 'tag'], dtype='object')
```

```
In [24]: tags.iloc[ [0,11,500] ]
```

```
Out[24]:
```

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

```
In [25]: # 📊 📉 Descriptive Statistics
ratings['rating'].describe()
```

```
Out[25]: 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 [26]: ratings.describe()
```

```
Out[26]:
```

	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 [27]: ratings['rating'].mean()
```

```
Out[27]: 3.5255285642993797
```

```
In [28]: ratings.mean()
```

```
Out[28]:
```

userId	69045.872583
movieId	9041.567330
rating	3.525529
dtype:	float64

```
In [29]: ratings['rating'].min()
```

```
Out[29]: 0.5
```

```
In [30]: ratings['rating'].max()
```

```
Out[30]: 5.0
```

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

```
Out[31]: 1.051988919275684
```

```
In [32]: ratings['rating'].mode()
```

```
Out[32]:
```

0	4.0
Name: rating, dtype: float64	

```
In [33]: ratings.corr()
```

```
Out[33]:
```

	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 [34]: 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[34]: False

```
In [35]: filter2 = ratings['rating'] > 0
filter2.all()
```

Out[35]: True

```
In [36]: # 🐞 Data Cleaning: Handling Missing Data
movies.shape
```

Out[36]: (27278, 3)

```
In [37]: movies.isnull().any().any()
```

Out[37]: False

```
In [38]: ratings.shape
```

Out[38]: (20000263, 3)

```
In [39]: ratings.isnull().any().any()
```

Out[39]: False

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

Out[40]: (465564, 3)

```
In [41]: tags.isnull().any().any()
```

Out[41]: True

```
In [43]: tags=tags.dropna()
```

```
In [44]: tags.isnull().any().any()
```

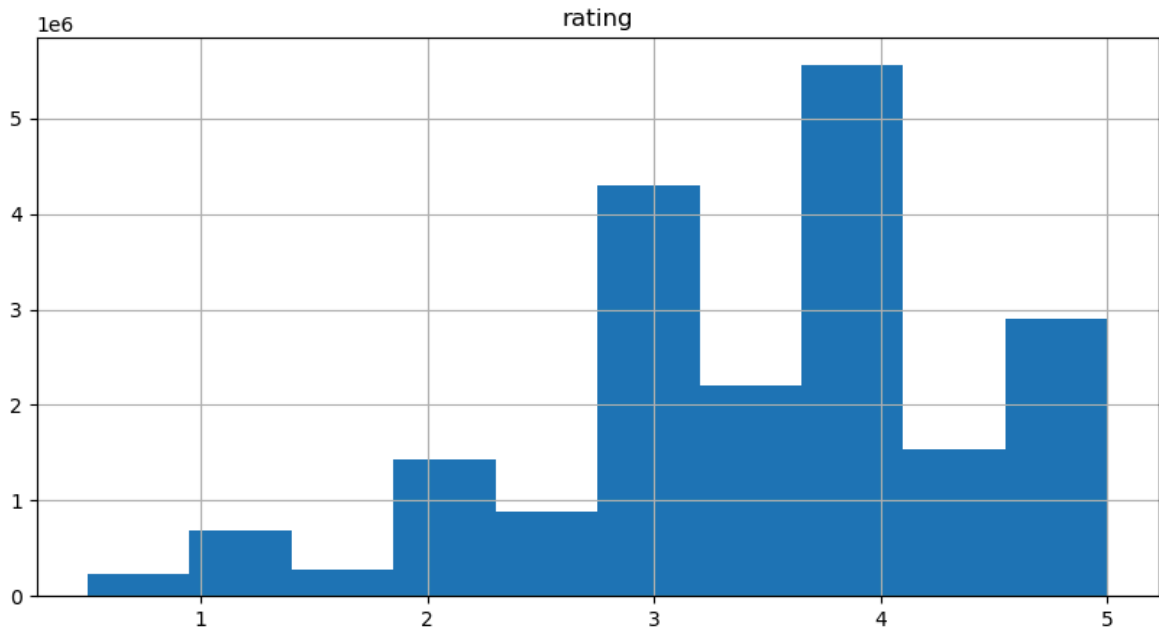
Out[44]: False

```
In [53]: tags.shape
```

Out[53]: (465548, 3)

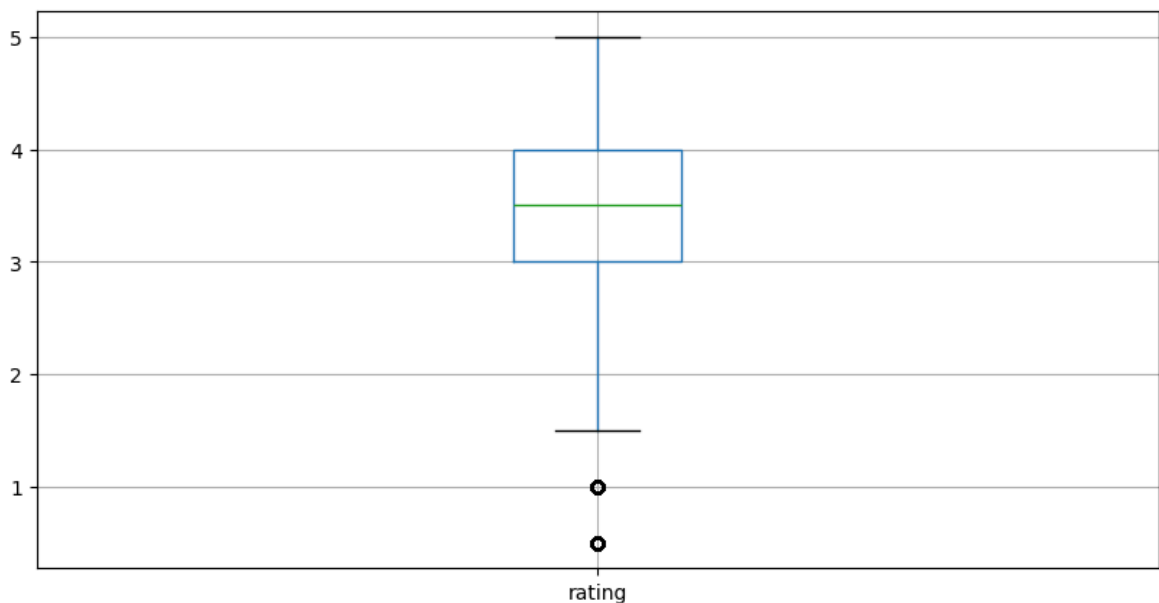
```
In [64]: # 📊 Data Visualization
%matplotlib inline
ratings.hist(column='rating', figsize=(10,5))
```

```
Out[64]: array([[<Axes: title={'center': 'rating'}>]], dtype=object)
```



```
In [66]: ratings.boxplot(column='rating', figsize=(10,5))
```

```
Out[66]: <Axes: >
```



```
In [67]: # 🗑️ Slicing Out Columns
tags['tag'].head()
```

```
Out[67]: 0    Mark Waters
1    dark hero
2    dark hero
3    noir thriller
4    dark hero
Name: tag, dtype: object
```

```
In [74]: movies[['title', 'genres']].head()
```


	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 [76]: ratings[-10:]

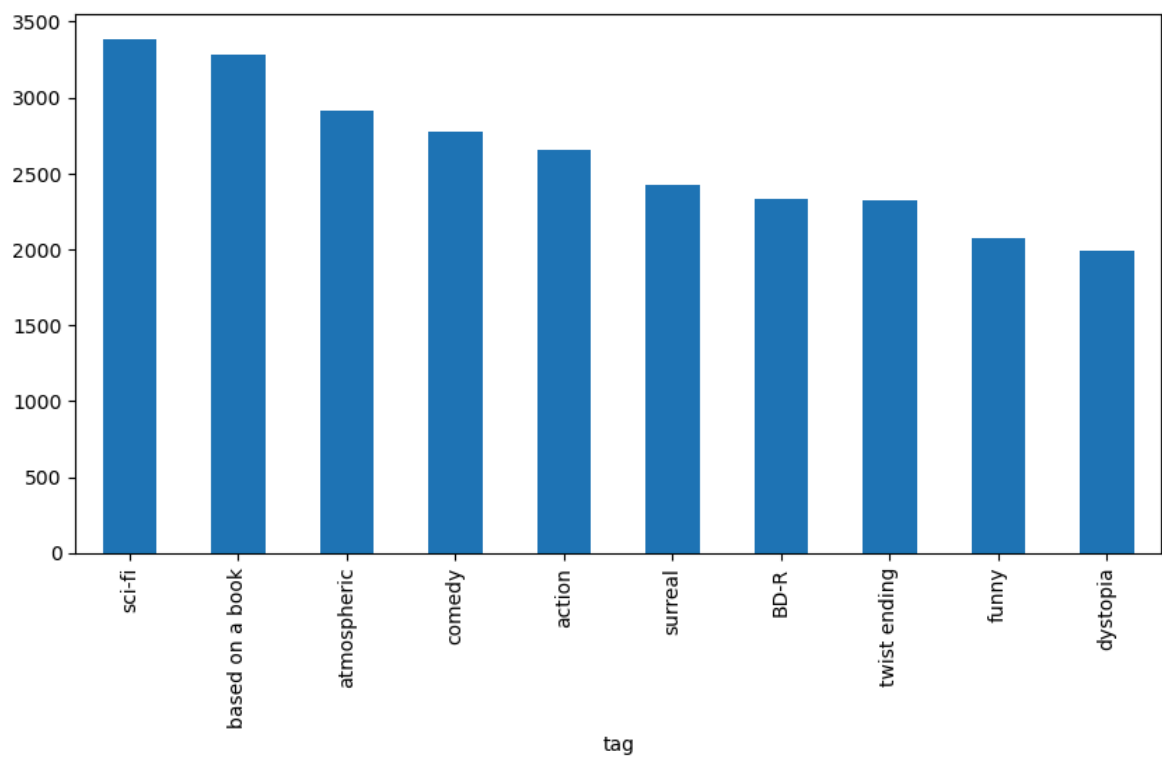
	userId	movieId	rating
20000253	138493	60816	4.5
20000254	138493	61160	4.0
20000255	138493	65682	4.5
20000256	138493	66762	4.5
20000257	138493	68319	4.5
20000258	138493	68954	4.5
20000259	138493	69526	4.5
20000260	138493	69644	3.0
20000261	138493	70286	5.0
20000262	138493	71619	2.5

In [89]: tag_counts = tags['tag'].value_counts()
tag_counts[-10:]

Out[89]: 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 [90]: tag_counts[:10].plot(kind='bar', figsize=(10,5))

Out[90]: <Axes: xlabel='tag'>



In []: