

PROJECT-4 : KAGGLE INTRODUCTION & IMDB Movie Rating Analysis

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

```
In [2]: ratings = pd.read_csv(r'E:\archive\rating.csv')
```

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

```
Out[3]: (20000263, 4)
```

```
In [4]: ratings.head(1)
```

```
Out[4]:
```

	userId	movieId	rating	timestamp
0	1	2	3.5	2005-04-02 23:53:47

```
In [5]: tags = pd.read_csv(r'E:\archive\tag.csv')
```

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

```
Out[6]: (465564, 4)
```

```
In [7]: tags.head(1)
```

```
Out[7]:
```

	userId	movieId	tag	timestamp
0	18	4141	Mark Waters	2009-04-24 18:19:40

```
In [8]: movies = pd.read_csv(r'E:\archive\movie.csv')
movies.head(1)
```

```
Out[8]:
```

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

```
In [9]: # For current analysis, we will remove timestamp

del ratings['timestamp']
del tags['timestamp']
```

```
In [10]: ratings.columns
```

```
Out[10]: Index(['userId', 'movieId', 'rating'], dtype='object')
```

```
In [11]: tags.columns
```

```
Out[11]: Index(['userId', 'movieId', 'tag'], dtype='object')
```

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

```
Out[12]:
```

	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

Data Structures

Series

```
In [13]: row_0 = tags.iloc[0]  
type(row_0)
```

```
Out[13]: pandas.core.series.Series
```

```
In [14]: print(row_0)
```

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

```
In [15]: row_0.index
```

```
Out[15]: Index(['userId', 'movieId', 'tag'], dtype='object')
```

```
In [16]: row_0['userId']
```

```
Out[16]: 18
```

```
In [17]: 'rating' in row_0
```

```
Out[17]: False
```

```
In [18]: row_0.name
```

```
Out[18]: 0
```

```
In [19]: row_0 = row_0.rename('firstRow')  
row_0.name
```

```
Out[19]: 'firstRow'
```

DataFrames

```
In [20]: tags.head()
```

```
Out[20]:
```

	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 [21]: tags.index
```

```
Out[21]: RangeIndex(start=0, stop=465564, step=1)
```

```
In [22]: tags.columns
```

```
Out[22]: Index(['userId', 'movieId', 'tag'], dtype='object')
```

```
In [23]: tags.iloc[ [0,11,500] ]
```

```
Out[23]:
```

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

Descriptive Statistics

how the ratings are distributed!

```
In [24]: ratings['rating'].describe()
```

```
Out[24]: 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 [25]: ratings.describe()
```

```
Out[25]:
```

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

```
Out[26]: 3.5255285642993797
```

```
In [27]: ratings.mean()
```

```
Out[27]: userId      69045.872583
         movieId     9041.567330
         rating        3.525529
         dtype: float64
```

```
In [28]: ratings['rating'].min()
```

```
Out[28]: 0.5
```

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

```
Out[29]: 5.0
```

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

```
Out[30]: 1.051988919275684
```

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

```
Out[31]: 0    4.0
         Name: rating, dtype: float64
```

```
In [32]: ratings.corr()
```

```
Out[32]:
```

	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 [33]: 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[33]: False
```

```
In [34]: filter2 = ratings['rating'] > 0
         filter2.all()
```

```
Out[34]: True
```

Data Cleaning: Handling Missing Data

```
In [35]: movies.shape
```

```
Out[35]: (27278, 3)
```

```
In [36]: movies.isnull().any().any()
```

```
Out[36]: False
```

No NULL values !

```
In [37]: ratings.shape
```

```
Out[37]: (20000263, 3)
```

```
In [38]: ratings.isnull().any().any()
```

```
Out[38]: False
```

No NULL values in Tags !

```
In [39]: tags.shape
```

```
Out[39]: (465564, 3)
```

```
In [40]: tags.isnull().any().any()
```

```
Out[40]: True
```

We have some tags which are NULL.

```
In [41]: tags=tags.dropna()
```

```
In [42]: tags.isnull().any().any()
```

```
Out[42]: False
```

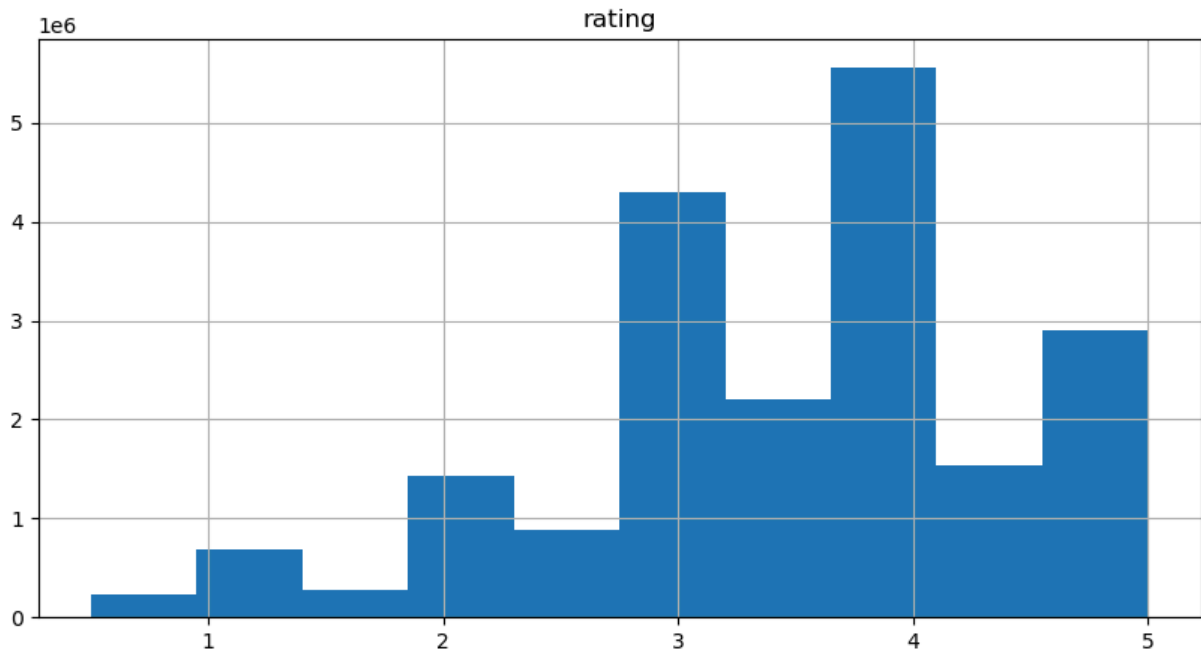
```
In [43]: tags.shape
```

```
Out[43]: (465548, 3)
```

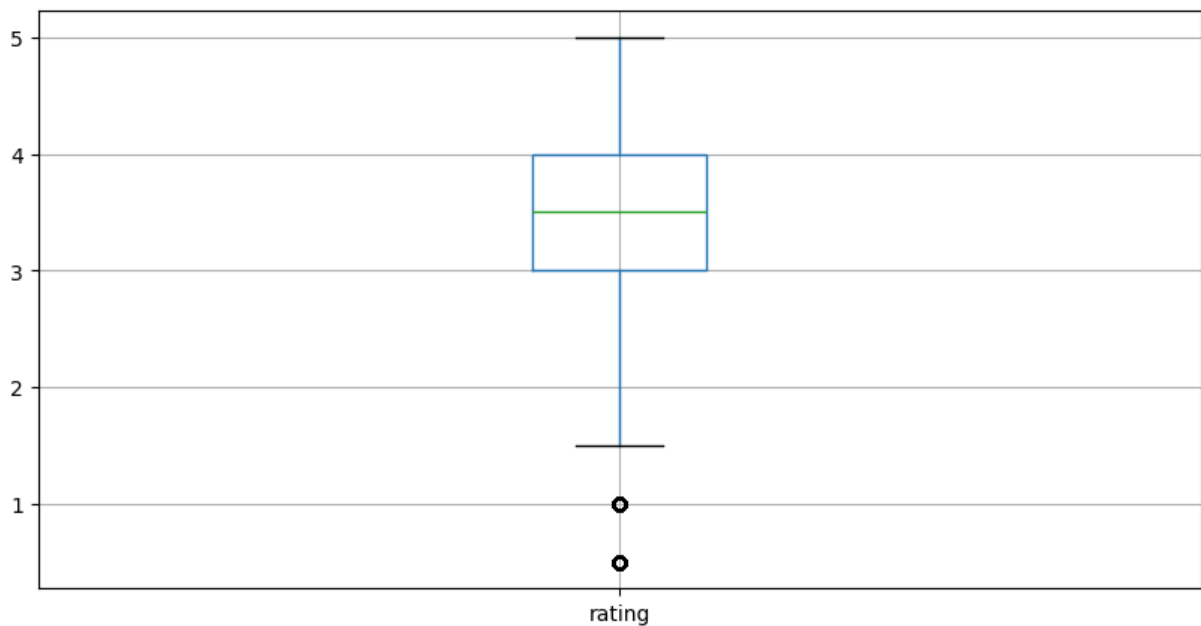
Thats nice ! No NULL values ! Notice the number of lines have reduced.

Data Visualization

```
In [45]: %matplotlib inline
import matplotlib.pyplot as plt
ratings.hist(column='rating', figsize=(10,5))
plt.show()
```



```
In [46]: ratings.boxplot(column='rating', figsize=(10,5))  
plt.show()
```



Slicing Out Columns

```
In [47]: tags['tag'].head()
```

```
Out[47]: 0    Mark Waters  
1    dark hero  
2    dark hero  
3    noir thriller  
4    dark hero  
Name: tag, dtype: object
```

```
In [48]: movies[['title', 'genres']].head()
```

```
Out[48]:
```

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

```
Out[49]:
```

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

```
Out[50]: 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 [51]: colors = plt.cm.Paired.colors
tag_counts[:10].plot(kind='bar', figsize=(10,5), color=colors)
```



```
plt.show()
```

