

```
In [3]: import pandas as pd
```

```
In [4]: movies = pd.read_csv('C:/Users/abhin/Desktop/Full Stack DS &
AI/Works/CSV/movie.csv', sep = ',')
```

```
In [5]: print(type(movies))
```

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

```
In [6]: movies.head(20)
```

```
Out[6]:
```

	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 [7]: tags = pd.read_csv('C:/Users/abhin/Desktop/Full Stack DS &
AI/Works/CSV/tag.csv', sep = ',')
```

```
In [9]: tags.head()
```

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

In [10]: `ratings = pd.read_csv('C:/Users/abhin/Desktop/Full Stack DS & AI/Works/CSV/rating.csv', sep = ',', parse_dates = ['timestamp'])`

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

Out[12]:

	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 [13]: `del ratings['timestamp']`

In [14]: `ratings.head()`

Out[14]:

	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 [15]: `del tags['timestamp']`

In [16]: `tags.head()`

```
Out[16]:
```

	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 [17]: row_0 = tags.iloc[0]
```

```
In [18]: type(row_0)
```

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

```
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]: 18
```

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

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

```
In [23]: row_0 = row_0.rename('firstRow')
```

```
In [24]: 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 [26]: tags.index
```

```
Out[26]: RangeIndex(start=0, stop=465564, step=1)
```

```
In [27]: tags.columns
```

```
Out[27]: Index(['userId', 'movieId', 'tag'], dtype='object')
```

```
In [28]: tags.iloc[[0,11,500]]
```

```
Out[28]:
```

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

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

```
Out[29]: 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 [30]: ratings.describe()
```

```
Out[30]:
```

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

```
Out[32]: 3.5255285642993797
```

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

```
Out[33]:
```

userId	69045.872583
movieId	9041.567330
rating	3.525529
dtype:	float64

```
In [34]: ratings['rating'].min()
```

```
Out[34]: 0.5
```

```
In [35]: ratings['rating'].max()
```

```
Out[35]: 5.0
```

```
In [36]: ratings['rating'].std()
```

```
Out[36]: 1.051988919275684
```

```
In [37]: ratings['rating'].mode()
```

```
Out[37]:
```

0	4.0
---	-----

Name: rating, dtype: float64

```
In [38]: ratings.corr()
```

```
Out[38]:
```

	userId	movieId	rating
<b>userId</b>	1.000000	-0.000850	0.001175
<b>movieId</b>	-0.000850	1.000000	0.002606
<b>rating</b>	0.001175	0.002606	1.000000

```
In [39]:
```

```
filter1 = ratings['rating'] > 10
```

```
In [40]: 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 [41]: filter1.any()
```

```
Out[41]: False
```

```
In [42]: filter2 = ratings['rating'] > 0
```

```
In [43]: filter2.all()
```

```
Out[43]: True
```

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

```
Out[44]: (27278, 3)
```

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

```
Out[45]: False
```

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

```
Out[46]: (20000263, 3)
```

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

```
Out[47]: False
```

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

```
Out[48]: (465564, 3)
```

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

```
Out[49]: True
```

```
In [50]: tags = tags.dropna()
```

```
In [51]: tags.isnull().any().any()
```

```
Out[51]: False
```

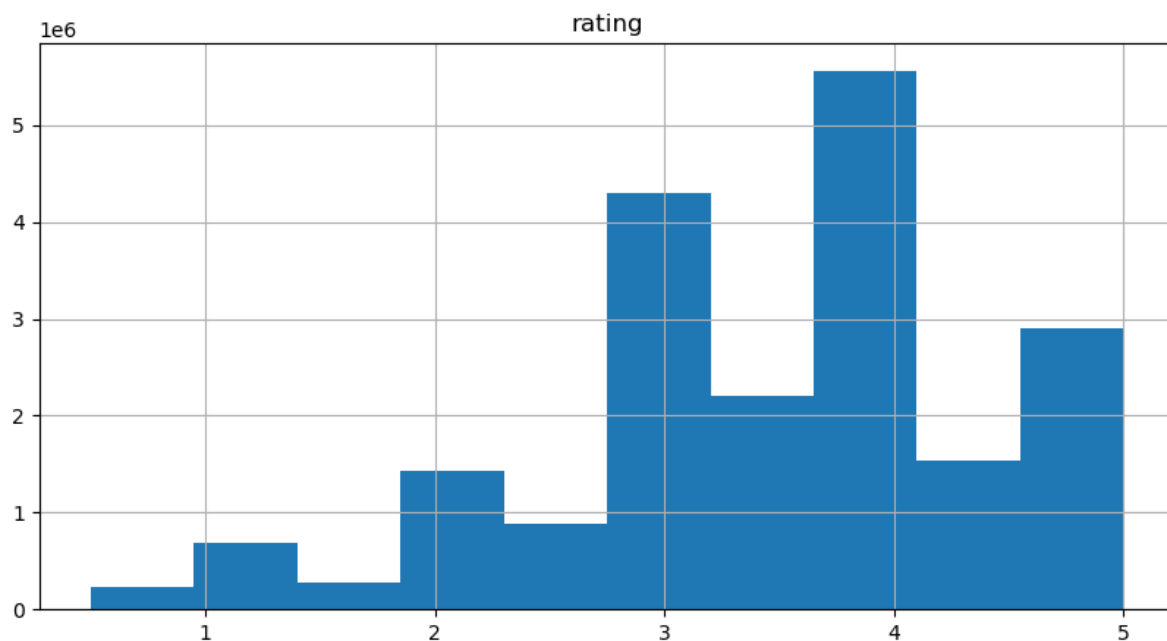
```
In [52]: tags.shape
```

```
Out[52]: (465548, 3)
```

```
In [57]: %matplotlib inline
```

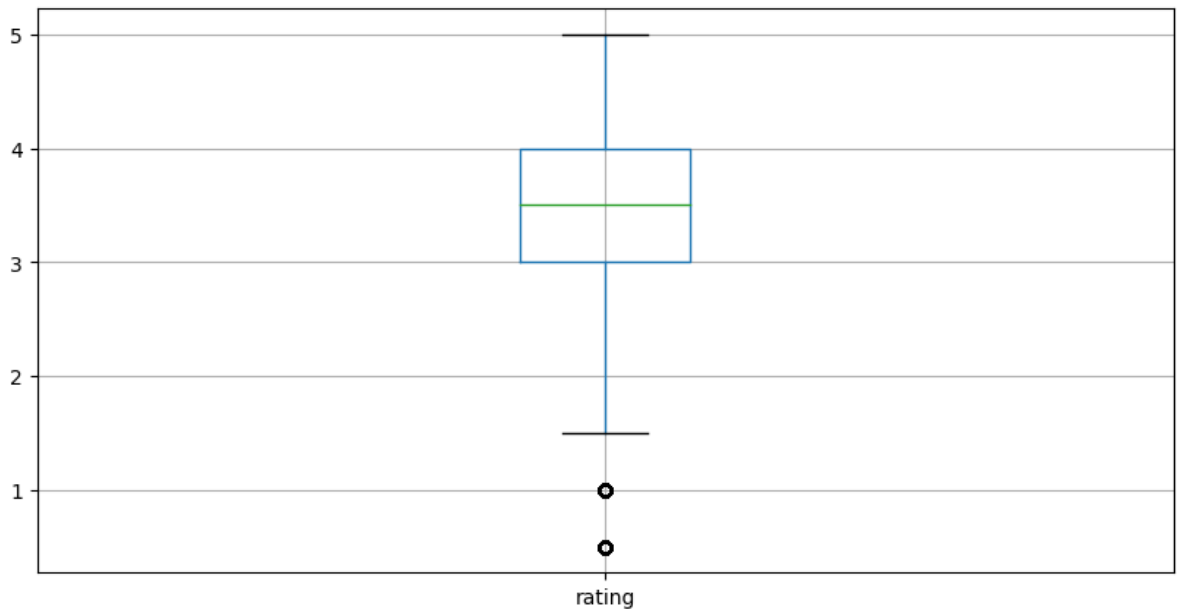
```
In [58]: ratings.hist(column = 'rating', figsize = (10,5))
```

```
Out[58]: array([[<AxesSubplot:title={'center':'rating'}>]], dtype=object)
```



```
In [59]: ratings.boxplot(column = 'rating', figsize = (10,5))
```

```
Out[59]: <AxesSubplot:>
```



```
In [60]: tags['tag'].head()
```

```
Out[60]: 0    Mark Waters
1      dark hero
2      dark hero
3    noir thriller
4      dark hero
Name: tag, dtype: object
```

```
In [63]: movies[['title', 'genres']].head()
```

```
Out[63]:
```

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



Out[64]:

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

In [69]: `tag_counts = tags['tag'].value_counts()`

In [70]: `tag_counts[-10:]`

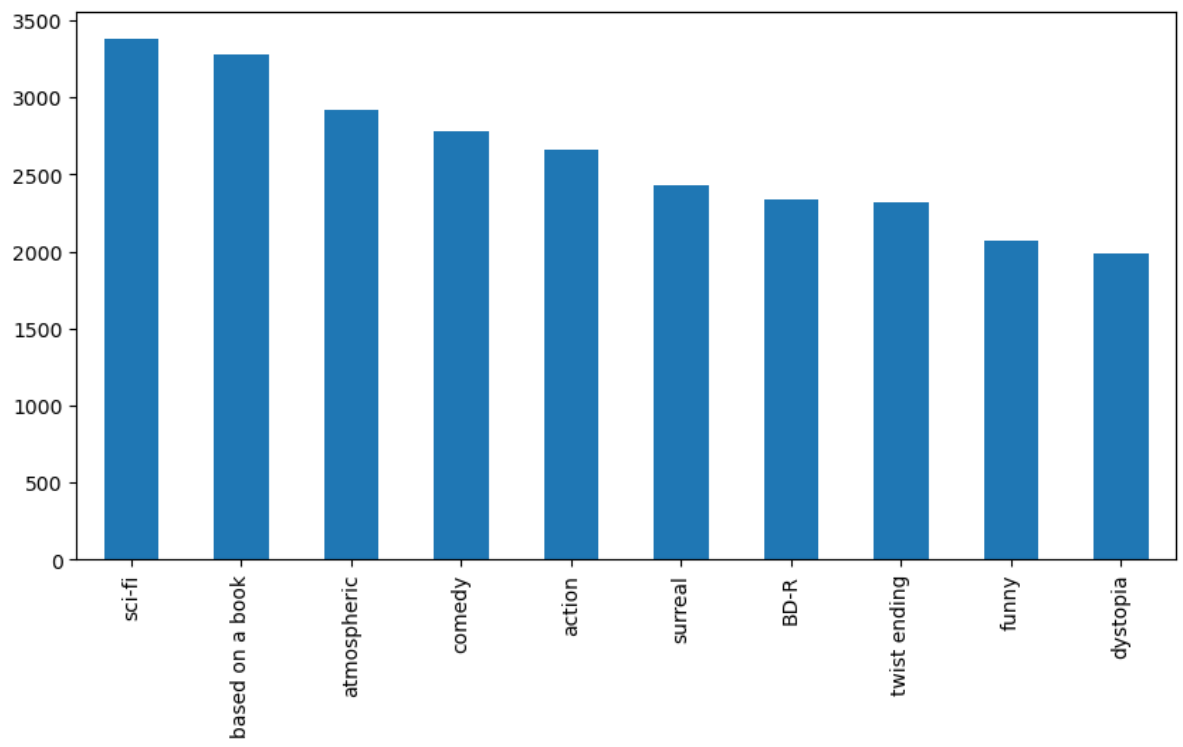
Out[70]:

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: tag, dtype: int64

In [72]: `tag_counts[:10].plot(kind = 'bar', figsize = (10,5))`

Out[72]: `<AxesSubplot:>`



In [ ]: