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

```
In [4]:  movies = pd.read_csv(r'C:\Users\Admin\Downloads\archive\movie.csv', sep=',')
```

```
In [5]:  movies
```

Out[5]:

|  | 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 |
| ... | ... | ... | ... |
| 27273 | 131254 | Kein Bund für's Leben (2007) | Comedy |
| 27274 | 131256 | Feuer, Eis & Dosenbier (2002) | Comedy |
| 27275 | 131258 | The Pirates (2014) | Adventure |
| 27276 | 131260 | Rentun Ruusu (2001) | (no genres listed) |
| 27277 | 131262 | Innocence (2014) | Adventure\|Fantasy\|Horror |

27278 rows × 3 columns

```
In [17]:  ratings  = pd.read_csv(r'C:\Users\Admin\Downloads\archive\rating.csv', sep=',', parse_dates=['timestamp'])
```

```
In [18]:  ratings
```

Out[18]:

|  | 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 |
| ... | ... | ... | ... | ... |
| 20000258 | 138493 | 68954 | 4.5 | 2009-11-13 15:42:00 |
| 20000259 | 138493 | 69526 | 4.5 | 2009-12-03 18:31:48 |
| 20000260 | 138493 | 69644 | 3.0 | 2009-12-07 18:10:57 |
| 20000261 | 138493 | 70286 | 5.0 | 2009-11-13 15:42:24 |
| 20000262 | 138493 | 71619 | 2.5 | 2009-10-17 20:25:36 |

20000263 rows × 4 columns

```
In [20]:  tags = pd.read_csv(r'C:\Users\Admin\Downloads\archive\tag.csv', sep=',')
```

In [21]: `tags`

Out[21]:

|  | 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 |
| ... | ... | ... | ... | ... |
| 465559 | 138446 | 55999 | dragged | 2013-01-23 23:29:32 |
| 465560 | 138446 | 55999 | Jason Bateman | 2013-01-23 23:29:38 |
| 465561 | 138446 | 55999 | quirky | 2013-01-23 23:29:38 |
| 465562 | 138446 | 55999 | sad | 2013-01-23 23:29:32 |
| 465563 | 138472 | 923 | rise to power | 2007-11-02 21:12:47 |

465564 rows × 4 columns

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 [23]: 
```python
del ratings['timestamp']
```

In [24]: 
```python
del tags['timestamp']
```

In [27]: 
```python
row_0 = tags.iloc[0]
```

In [28]: 
```python
row_0
```

Out[28]: 
```
userId                18
movieId             4141
tag          Mark Waters
Name: 0, dtype: object
```

In [30]: 
```python
row_0.index
```

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

In [31]: 
```python
row_0['userId']
```

Out[31]: 18

In [32]: 
```python
'rating' in row_0
```

Out[32]: False

In [33]: `row_0.name`

Out[33]: 0

In [34]: `row_0 = row_0.rename('firstRow')`

In [35]: `row_0`

Out[35]:
```
userId                18
movieId             4141
tag         Mark Waters
Name: firstRow, dtype: object
```

In [36]: `tags.head()`

Out[36]:

|   | 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 [37]: `tags.index`

Out[37]: `RangeIndex(start=0, stop=465564, step=1)`

In [38]: `tags.columns`

Out[38]: `Index(['userId', 'movieId', 'tag'], dtype='object')`

In [42]: `tags.iloc[ [0,11,500] ]`

Out[42]:

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

In [43]: `ratings['rating'].describe()`

Out[43]:
```
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 [44]: ratings

Out[44]:

|  | 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 |
| ... | ... | ... | ... |
| 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 |

20000263 rows × 3 columns

In [45]: ratings.describe()

Out[45]:

|  | 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 [46]: ratings['rating'].mean()

Out[46]: 3.5255285642993797

In [47]: ratings.mean()

Out[47]:
```
userId      69045.872583
movieId      9041.567330
rating          3.525529
dtype: float64
```

In [48]: ratings['rating'].min()

Out[48]: 0.5

In [49]: ratings['rating'].max()

Out[49]: 5.0

In [50]: ratings['rating'].std()

Out[50]: 1.051988919275684

```
In [51]: ratings['rating'].mode()
```

```
Out[51]: 0    4.0
         Name: rating, dtype: float64
```

```
In [52]: ratings.corr()
```

Out[52]:

|         | 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 [53]: filter1 = ratings['rating'] > 10
```

```
In [54]: 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 [55]: filter1.any()
```

```
Out[55]: False
```

```
In [56]: filter2 = ratings['rating'] > 0
```

```
In [57]: filter2.all()
```

```
Out[57]: True
```

```
In [58]: movies.shape
```

```
Out[58]: (27278, 3)
```

```
In [62]: movies.isnull().any().any()
```

```
Out[62]: False
```

```
In [63]: ratings.shape
```

```
Out[63]: (20000263, 3)
```

```
In [64]: ratings.isnull().any().any()
```

```
Out[64]: False
```

```
In [65]: tags.shape
```

```
Out[65]: (465564, 3)
```

In [66]: `tags.isnull().any().any()`

Out[66]: True

In [67]: `tags=tags.dropna()`

In [68]: `tags`

Out[68]:

|        | 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 |
| ...    | ...    | ...     | ... |
| 465559 | 138446 | 55999   | dragged |
| 465560 | 138446 | 55999   | Jason Bateman |
| 465561 | 138446 | 55999   | quirky |
| 465562 | 138446 | 55999   | sad |
| 465563 | 138472 | 923     | rise to power |

465548 rows × 3 columns
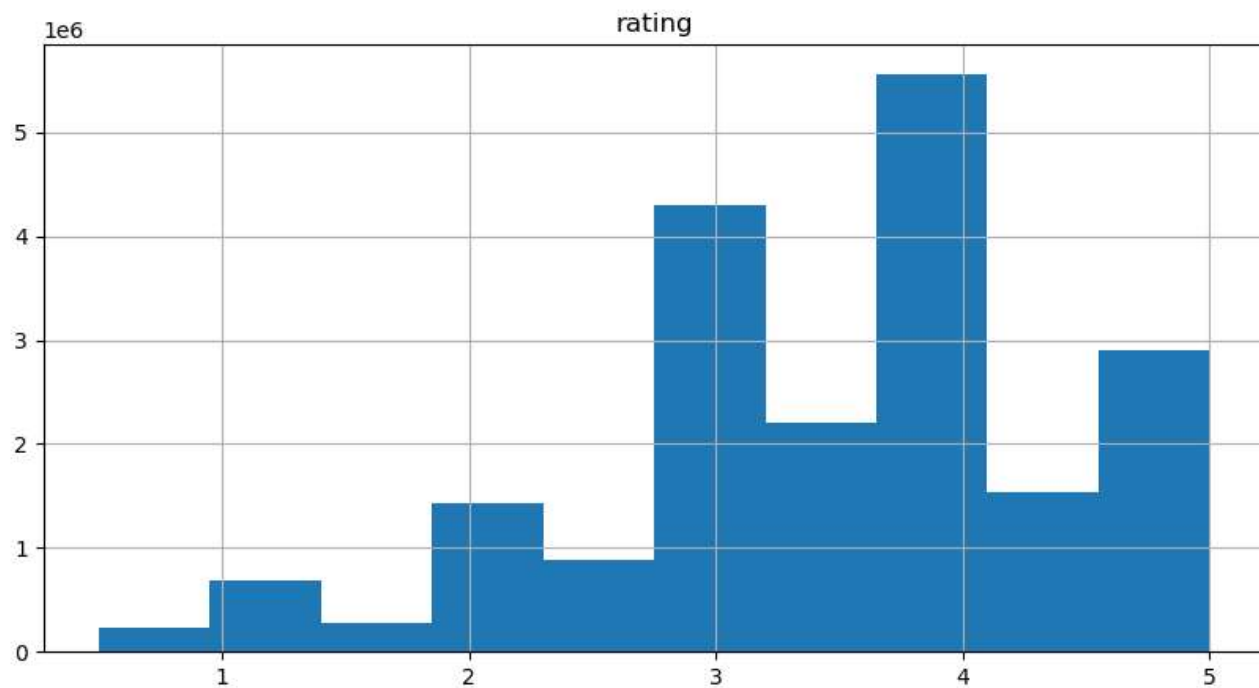
In [69]: `tags.shape`

Out[69]: (465548, 3)

In [70]: `tags.isnull().any().any()`
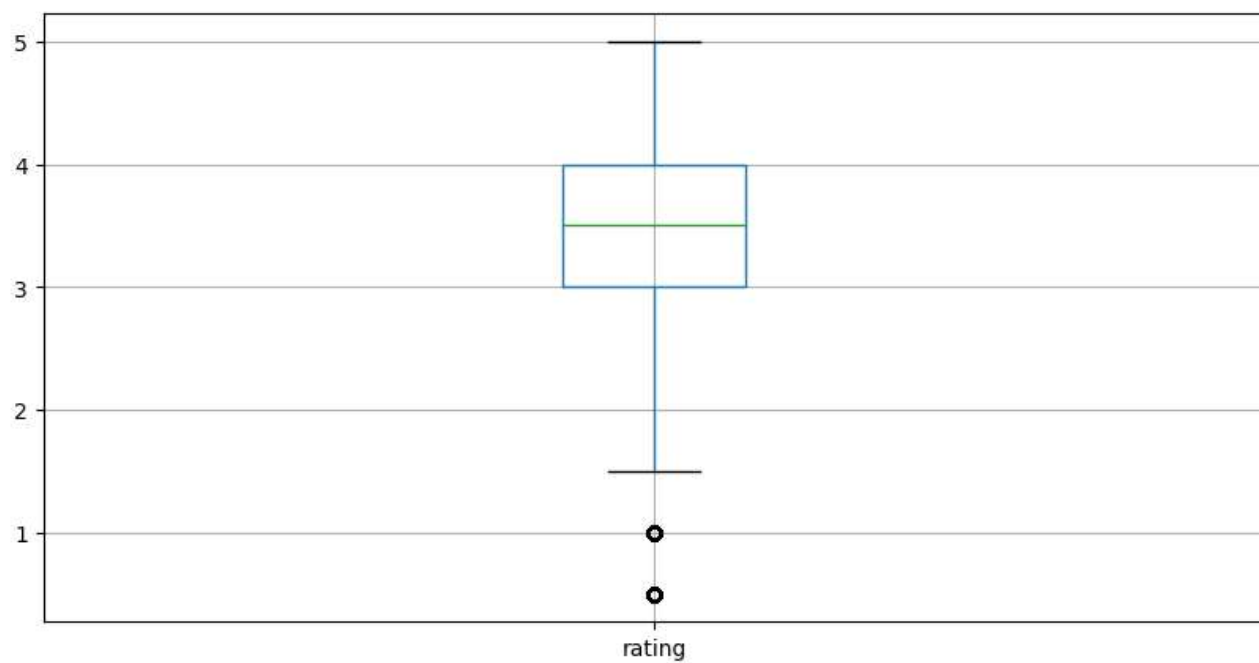
Out[70]: False

In [71]: 
```
%matplotlib inline

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

Out[71]: array([[<Axes: title={'center': 'rating'}>]], dtype=object)



In [72]: 
```
ratings.boxplot(column='rating', figsize=(10,5))
```

Out[72]: <Axes: >

```
In [73]: tags['tag'].head()
```

```
Out[73]: 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()
```

Out[74]:

|   | 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 [75]: ratings[-10:]
```

Out[75]:

|   | 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 [76]: tag_counts = tags['tag'].value_counts()
```

```
In [77]: tag_counts[-10:]
```
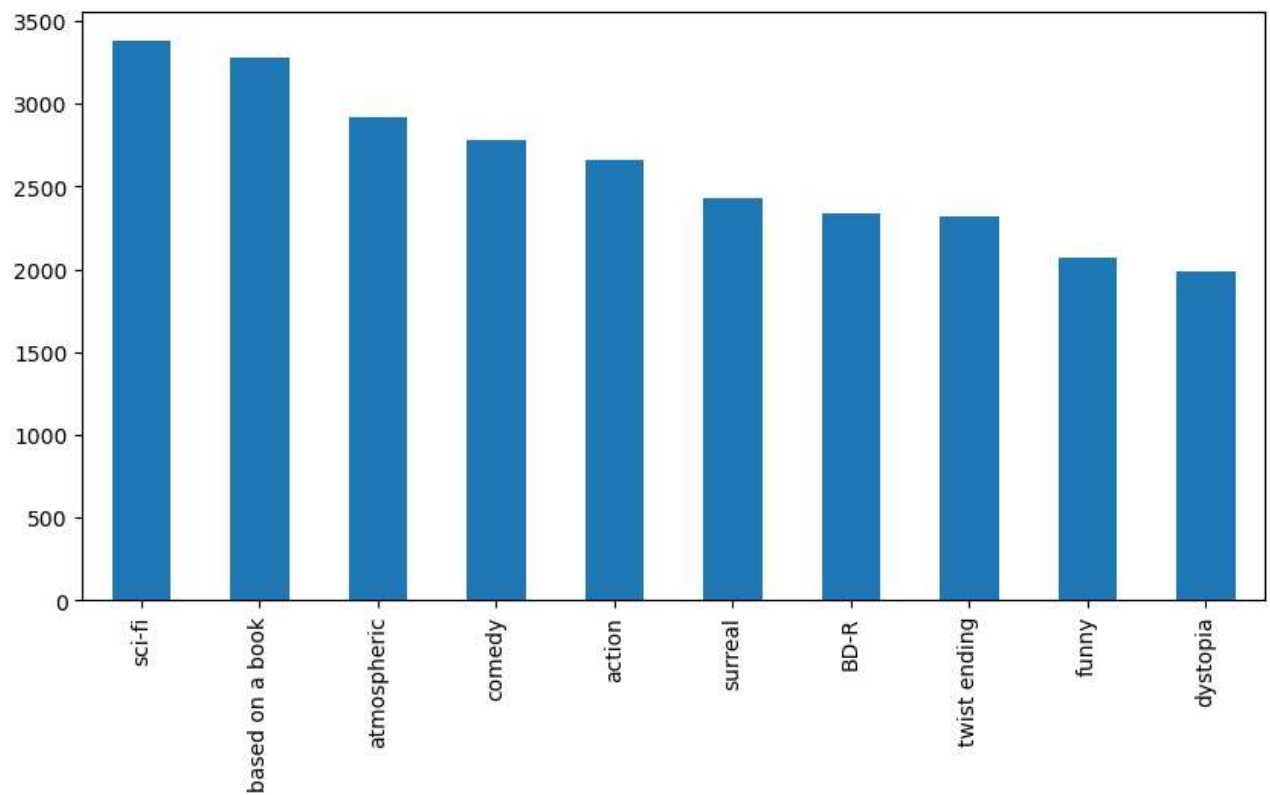
```
Out[77]: missing child                1
         Ron Moore                    1
         Citizen Kane                 1
         mullet                       1
         biker gang                   1
         Paul Adelstein               1
         the wig                      1
         killer fish                  1
         genetically modified monsters    1
         topless scene                1
         Name: tag, dtype: int64
```

In [78]: `tag_counts[:10].plot(kind='bar', figsize=(10,5))`

Out[78]: `<Axes: >`



In [79]: `tag_counts[:10]`

Out[79]:
```
sci-fi             3384
based on a book    3281
atmospheric        2917
comedy             2779
action             2657
surreal            2427
BD-R               2334
twist ending       2323
funny              2072
dystopia           1991
Name: tag, dtype: int64
```

In [ ]: