

```
# import pandas library
import pandas as pd

# Get the data
column_names = ['user_id', 'item_id', 'rating', 'timestamp']

path = 'https://media.geeksforgeeks.org/wp-content/uploads/file.tsv'

df = pd.read_csv(path, sep='\t', names=column_names)

# Check the head of the data
df.head()
```

↗

	user_id	item_id	rating	timestamp
0	0	50	5	881250949
1	0	172	5	881250949
2	0	133	1	881250949
3	196	242	3	881250949
4	186	302	3	891717742

```
# Check out all the movies and their respective IDs
movie_titles = pd.read_csv('https://media.geeksforgeeks.org/wp-content/uploads/Movie_Id_Title')
movie_titles.head()
```

	item_id	title
0	1	Toy Story (1995)
1	2	GoldenEye (1995)
2	3	Four Rooms (1995)
3	4	Get Shorty (1995)
4	5	Copycat (1995)

```
data = pd.merge(df, movie_titles, on='item_id')
data.head()
```

	user_id	item_id	rating	timestamp	title
0	0	50	5	881250949	Star Wars (1977)

```
# Calculate mean rating of all movies
```

```
data.groupby('title')['rating'].mean().sort_values(ascending=False).head()
```

```
title
They Made Me a Criminal (1939)      5.0
Marlene Dietrich: Shadow and Light (1996)  5.0
Saint of Fort Washington, The (1993)      5.0
Someone Else's America (1995)          5.0
Star Kid (1997)                      5.0
Name: rating, dtype: float64
```

```
# Calculate count rating of all movies
```

```
data.groupby('title')['rating'].count().sort_values(ascending=False).head()
```

```
title
Star Wars (1977)      584
Contact (1997)        509
Fargo (1996)          508
Return of the Jedi (1983)  507
Liar Liar (1997)       485
Name: rating, dtype: int64
```

```
# creating dataframe with 'rating' count values
```

```
ratings = pd.DataFrame(data.groupby('title')['rating'].mean())
```

```
ratings['num of ratings'] = pd.DataFrame(data.groupby('title')['rating'].count())
```

```
ratings.head()
```

	rating	num of ratings
title		
<b>'Til There Was You (1997)</b>	2.333333	9
<b>1-900 (1994)</b>	2.600000	5
<b>101 Dalmatians (1996)</b>	2.908257	109
<b>12 Angry Men (1957)</b>	4.344000	125
<b>187 (1997)</b>	3.024390	41

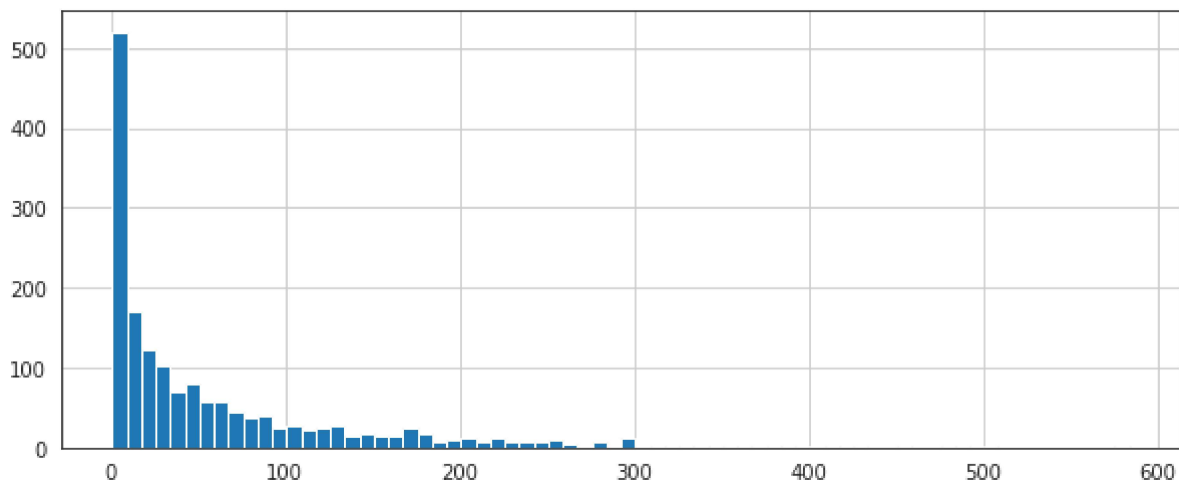
```
import matplotlib.pyplot as plt
import seaborn as sns
```

```
sns.set_style('white')
%matplotlib inline

# plot graph of 'num of ratings column'
plt.figure(figsize=(10, 4))

ratings['num of ratings'].hist(bins = 70)
```

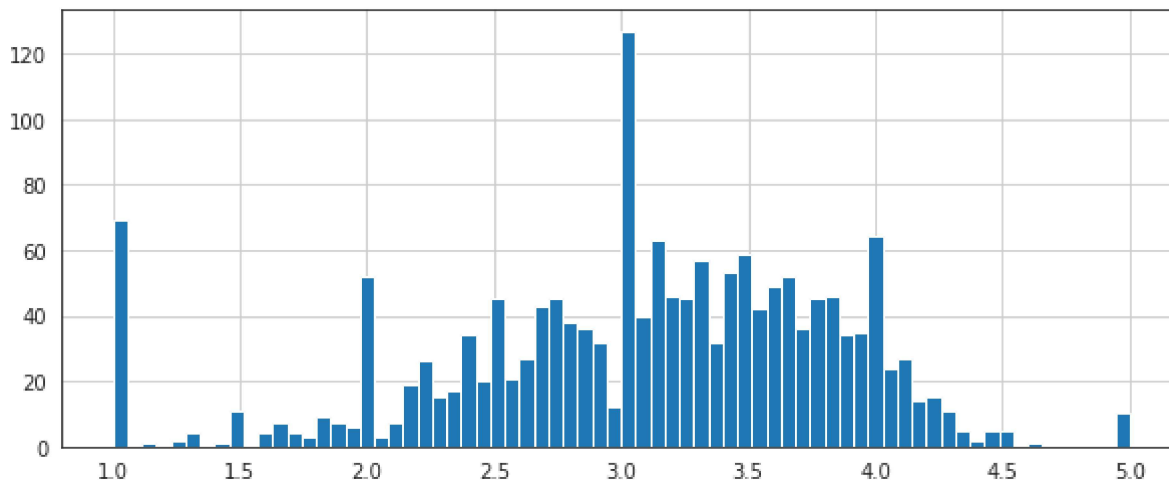
<matplotlib.axes.\_subplots.AxesSubplot at 0x7f8717f53490>



```
# plot graph of 'ratings' column
plt.figure(figsize=(10, 4))

ratings['rating'].hist(bins = 70)
```

<matplotlib.axes.\_subplots.AxesSubplot at 0x7f8717dc7d10>



```
# Sorting values according to
# the 'num of rating column'
moviemat = data.pivot_table(index='user_id',
                             columns='title', values='rating')

moviemat.head()

ratings.sort_values('num of ratings', ascending = False).head(10)
```

	rating	num of ratings
title		
<b>Star Wars (1977)</b>	4.359589	584
<b>Contact (1997)</b>	3.803536	509
<b>Fargo (1996)</b>	4.155512	508
<b>Return of the Jedi (1983)</b>	4.007890	507
<b>Liar Liar (1997)</b>	3.156701	485
<b>English Patient, The (1996)</b>	3.656965	481
<b>Scream (1996)</b>	3.441423	478
<b>Toy Story (1995)</b>	3.878319	452
<b>Air Force One (1997)</b>	3.631090	431
<b>Independence Day (ID4) (1996)</b>	3.438228	429

```
# analysing correlation with similar movies
starwars_user_ratings = moviemat['Star Wars (1977)']
liarliar_user_ratings = moviemat['Liar Liar (1997)']
```

```
starwars_user_ratings.head()
```

```
user_id
0      5.0
1      5.0
2      5.0
3      NaN
4      5.0
Name: Star Wars (1977), dtype: float64
```

```
# analysing correlation with similar movies
similar_to_starwars = moviemat.corrwith(starwars_user_ratings)
similar_to_liarliar = moviemat.corrwith(liarliar_user_ratings)

corr_starwars = pd.DataFrame(similar_to_starwars, columns=['Correlation'])
corr_starwars.dropna(inplace = True)

corr_starwars.head()
```

**Correlation****title**

<b>'Til There Was You (1997)</b>	0.872872
----------------------------------	----------

<b>1.000 (1994)</b>	0.645497
---------------------	----------

```
# Similar movies like starwars
```

```
corr_starwars.sort_values('Correlation', ascending = False).head(10)
```

```
corr_starwars = corr_starwars.join(ratings['num of ratings'])
```

```
corr_starwars.head()
```

```
corr_starwars[corr_starwars['num of ratings']>100].sort_values('Correlation', ascending = Fal
```

**Correlation num of ratings****title**

<b>title</b>	<b>Correlation</b>	<b>num of ratings</b>
<b>Star Wars (1977)</b>	1.000000	584
<b>Empire Strikes Back, The (1980)</b>	0.748353	368
<b>Return of the Jedi (1983)</b>	0.672556	507
<b>Raiders of the Lost Ark (1981)</b>	0.536117	420
<b>Austin Powers: International Man of Mystery (1997)</b>	0.377433	130

```
# Similar movies as of liarliar
```

```
corr_liarliar = pd.DataFrame(similar_to_liarliar, columns =['Correlation'])
```

```
corr_liarliar.dropna(inplace = True)
```

```
corr_liarliar = corr_liarliar.join(ratings['num of ratings'])
```

```
corr_liarliar[corr_liarliar['num of ratings']>100].sort_values('Correlation', ascending = Fal
```

**Correlation num of ratings****title**

<b>title</b>	<b>Correlation</b>	<b>num of ratings</b>
<b>Liar Liar (1997)</b>	1.000000	485
<b>Batman Forever (1995)</b>	0.516968	114
<b>Mask, The (1994)</b>	0.484650	129
<b>Down Periscope (1996)</b>	0.472681	101
<b>Con Air (1997)</b>	0.469828	137

