

```
# Amazon Top 50 Bestselling Books 2009-2022
```

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
%matplotlib inline
```

```
In [... data =pd.read_csv(r"C:\Users\Gulsum\Downloads\bestsellers_with_categories_2022_03_27.csv"
df = pd.DataFrame(data)
```

```
In [8]: df.head()
```

```
Out[... 
```

	Name	Author	User Rating	Reviews	Price	Year	Genre
0	Act Like a Lady, Think Like a Man: What Men Re...	Steve Harvey	4.6	5013	17	2009	Non Fiction
1	Arguing with Idiots: How to Stop Small Minds a...	Glenn Beck	4.6	798	5	2009	Non Fiction
2	Breaking Dawn (The Twilight Saga, Book 4)	Stephenie Meyer	4.6	9769	13	2009	Fiction
3	Crazy Love: Overwhelmed by a Relentless God	Francis Chan	4.7	1542	14	2009	Non Fiction
4	Dead And Gone: A Sookie Stackhouse Novel (Sook...	Charlaine Harris	4.6	1541	4	2009	Fiction

```
In [10]: #Shows descriptive statistics data
df.describe()
```

```
Out[10]: 
```

	User Rating	Reviews	Price	Year
count	700.000000	700.000000	700.000000	700.000000
mean	4.639857	19255.195714	12.700000	2015.500000
std	0.218586	23613.443875	9.915162	4.034011
min	3.300000	37.000000	0.000000	2009.000000
25%	4.500000	4987.250000	7.000000	2012.000000
50%	4.700000	10284.000000	11.000000	2015.500000
75%	4.800000	23358.000000	15.000000	2019.000000
max	4.900000	208917.000000	105.000000	2022.000000

```
In [56]: #cheking for null values
for i in df.columns:
    print(i,"\t-\t", df[i].isna().mean()*100)
```

```
Name      -      0.0
Author    -      0.0
User Rating -      0.0
```

```
Reviews      -      0.0
Price -      0.0
Year -      0.0
Genre -      0.0
```

```
In [11]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 700 entries, 0 to 699
Data columns (total 7 columns):
#   Column          Non-Null Count  Dtype
---  -
0   Name            700 non-null   object
1   Author          700 non-null   object
2   User Rating     700 non-null   float64
3   Reviews         700 non-null   int64
4   Price           700 non-null   int64
5   Year            700 non-null   int64
6   Genre           700 non-null   object
dtypes: float64(1), int64(3), object(3)
memory usage: 38.4+ KB
```

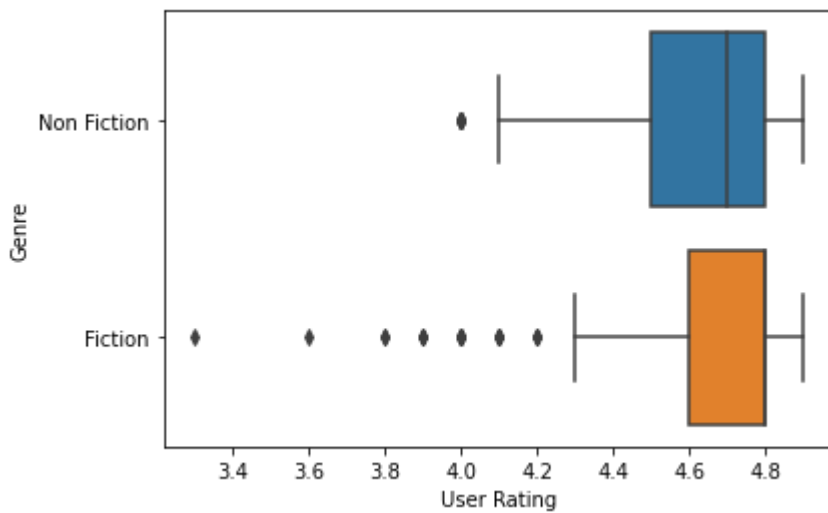
```
In [3... #sort the values of the 10 books in a ascending order from top to bottom by User Ratings
top10=df.sort_values('User Rating',ascending=False)[:10]
```

```
In [32]: top10
```

```
Out[3...
Name Author User Rating Reviews Price Year Genre
605 Brown Bear, Brown Bear, What Do You See? Bill Martin Jr. 4.9 38969 5 2021 Fiction
607 Call Us What We Carry: Poems Amanda Gorman 4.9 2873 14 2021 Fiction
457 Dog Man: Brawl of the Wild: From the Dav Pilkey 4.9 7235 4 2018 Fiction
Creator o...
456 Dog Man and Cat Kid: From the Creator of Dav Pilkey 4.9 5062 6 2018 Fiction
Capta...
223 Oh, the Places You'll Go! Dr. Seuss 4.9 21834 8 2013 Fiction
586 The Deep End (Diary of a Wimpy Kid Book 15) Jeff Kinney 4.9 38674 7 2020 Fiction
227 Rush Revere and the Brave Pilgrims: Time- Rush Limbaugh 4.9 7150 12 2013 Fiction
Trave...
443 The Wonderful Things You Will Be Emily Winfield Martin 4.9 8842 10 2017 Fiction
592 The Very Hungry Caterpillar Eric Carle 4.9 47260 5 2020 Fiction
441 The Very Hungry Caterpillar Eric Carle 4.9 19546 5 2017 Fiction
```

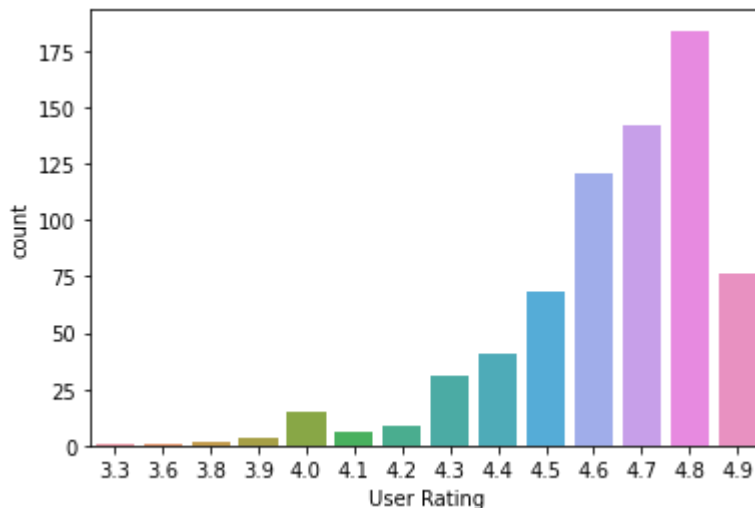
```
In [36]: #Used seaborn to graph the genre and the user ratings.
#Acorrding to the graph fiction is more popular than non-fiction on Amazon
sns.boxplot(x='User Rating', y='Genre',data=df)
```

```
Out[36]:<AxesSubplot:xlabel='User Rating', ylabel='Genre'>
```



In [26]: *#shows a visual reappresentation of user rating*
`sns.countplot(x = df['User Rating'])`

Out[26]:<AxesSubplot:xlabel='User Rating', ylabel='count'>



In [71]: *#Used pandas to extract data from the column user ratings that is equal to 4.9*
#I used the groupby function to group author and user rating coulumn.
#Shows the top authors with the highest ratings.
`bestsellers = df[df['User Rating']==5.0]`
`bestsellers = bestsellers.groupby('Author')['User Rating']`

In [72]: `bestsellers`

Out[72]:<pandas.core.groupby.generic.SeriesGroupBy object at 0x000001F8A1917C10>

In [...]: *# I made the new year set to years from 2009-2022.*
#I used the mean() function to give the average of the other numeric columns.
#Reset index to reset the index after making modifications to the column
`pyear = df.groupby('Year').mean().reset_index()`
`pyear['Year'] = [2009, 2010, 2011, 2012, 2013, 2014, 2015, 2016, 2017, 2018, 2019, 2020,`
`pyear`

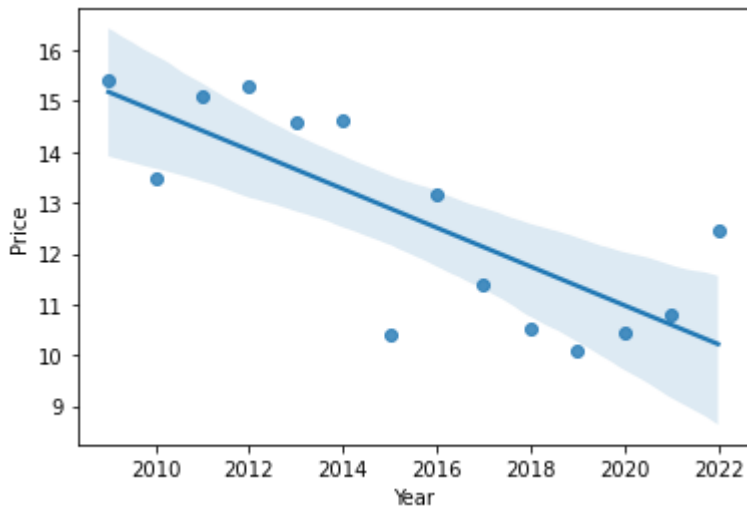
Out[15]:

	Year	User Rating	Reviews	Price
0	2009	4.584	4710.12	15.40

	Year	User Rating	Reviews	Price
1	2010	4.558	5479.62	13.48
2	2011	4.558	8100.82	15.10
3	2012	4.532	13090.92	15.30
4	2013	4.554	13098.14	14.60
5	2014	4.622	15859.94	14.64
6	2015	4.648	14233.38	10.42
7	2016	4.678	14196.00	13.18
8	2017	4.660	12888.40	11.38
9	2018	4.668	13930.42	10.52
10	2019	4.740	15898.34	10.08
11	2020	4.726	52349.94	10.46
12	2021	4.738	44859.48	10.78
13	2022	4.692	40877.22	12.46

In [16]: *#plots the linear regression model of the data from x and y.*
#From the data the amazon price of books declined as the years went on.
`sns.regplot(x="Year", y="Price", data=pyear)`

Out[16]:<AxesSubplot:xlabel='Year', ylabel='Price'>



In []: