

# PYTHON PROJECT (Vinay H)

1.

```
In [1]: import pandas as pd
import numpy as np
import seaborn as sns
```

```
In [2]: data = pd.read_csv('googleplaystore.csv')
```

```
In [3]: data.head()
```

Out[3]:

	App	Category	Rating	Reviews	Size	Installs	Type	Price	Content Rating
0	Photo Editor & Candy Camera & Grid & ScrapBook	ART_AND_DESIGN	4.1	159	19M	10,000+	Free	0	Everyone
1	Coloring book moana	ART_AND_DESIGN	3.9	967	14M	500,000+	Free	0	Everyone
2	U Launcher Lite – FREE Live Cool Themes, Hide ...	ART_AND_DESIGN	4.7	87510	8.7M	5,000,000+	Free	0	Everyone
3	Sketch - Draw & Paint	ART_AND_DESIGN	4.5	215644	25M	50,000,000+	Free	0	Teen
4	Pixel Draw - Number Art Coloring Book	ART_AND_DESIGN	4.3	967	2.8M	100,000+	Free	0	Everyone

```
In [4]: data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 10841 entries, 0 to 10840
Data columns (total 13 columns):
#   Column                Non-Null Count  Dtype
---  -
0   App                    10841 non-null  object
1   Category               10841 non-null  object
2   Rating                 9367 non-null   float64
3   Reviews                10841 non-null  object
4   Size                   10841 non-null  object
5   Installs               10841 non-null  object
6   Type                   10840 non-null  object
7   Price                  10841 non-null  object
8   Content Rating         10840 non-null  object
9   Genres                 10841 non-null  object
10  Last Updated           10841 non-null  object
11  Current Ver            10833 non-null  object
12  Android Ver            10838 non-null  object
dtypes: float64(1), object(12)
memory usage: 1.1+ MB
```

```
In [5]: data.shape
```

```
Out[5]: (10841, 13)
```

## 2.

```
In [6]: data.isnull().any()
```

```
Out[6]: App                    False
Category                   False
Rating                     True
Reviews                    False
Size                       False
Installs                   False
Type                       True
Price                      False
Content Rating             True
Genres                     False
Last Updated               False
Current Ver                True
Android Ver                True
dtype: bool
```

```
In [7]: data.isnull().sum()
```

```
Out[7]: App                0
        Category           0
        Rating            1474
        Reviews            0
        Size               0
        Installs           0
        Type               1
        Price              0
        Content Rating     1
        Genres             0
        Last Updated       0
        Current Ver        8
        Android Ver        3
        dtype: int64
```

### 3.

```
In [8]: data = data.dropna()
```

```
In [9]: data.isnull().any()
```

```
Out[9]: App                False
        Category           False
        Rating             False
        Reviews            False
        Size               False
        Installs           False
        Type               False
        Price              False
        Content Rating     False
        Genres             False
        Last Updated       False
        Current Ver        False
        Android Ver        False
        dtype: bool
```

```
In [10]: data.shape
```

```
Out[10]: (9360, 13)
```


### 4(I).

```
In [11]: data["Size"] = [ float(i.split('M')[0]) if 'M' in i else float(0) for i in dat
```

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

Out[12]:

	App	Category	Rating	Reviews	Size	Installs	Type	Price	Content Rating	
0	Photo Editor & Candy Camera & Grid & ScrapBook	ART_AND_DESIGN	4.1	159	19.0	10,000+	Free	0	Everyone	
1	Coloring book moana	ART_AND_DESIGN	3.9	967	14.0	500,000+	Free	0	Everyone	D
2	U Launcher Lite – FREE Live Cool Themes, Hide ...	ART_AND_DESIGN	4.7	87510	8.7	5,000,000+	Free	0	Everyone	
3	Sketch - Draw & Paint	ART_AND_DESIGN	4.5	215644	25.0	50,000,000+	Free	0	Teen	
4	Pixel Draw - Number Art Coloring Book	ART_AND_DESIGN	4.3	967	2.8	100,000+	Free	0	Everyone	De



```
In [13]: data["Size"] = 1000 * data["Size"]
```

In [14]: data

Out[14]:

	App	Category	Rating	Reviews	Size	Installs	Type	Price
0	Photo Editor & Candy Camera & Grid & ScrapBook	ART_AND_DESIGN	4.1	159	19000.0	10,000+	Free	0
1	Coloring book moana	ART_AND_DESIGN	3.9	967	14000.0	500,000+	Free	0
2	U Launcher Lite – FREE Live Cool Themes, Hide ...	ART_AND_DESIGN	4.7	87510	8700.0	5,000,000+	Free	0
3	Sketch - Draw & Paint	ART_AND_DESIGN	4.5	215644	25000.0	50,000,000+	Free	0
4	Pixel Draw - Number Art Coloring Book	ART_AND_DESIGN	4.3	967	2800.0	100,000+	Free	0
...	...	...	...	...	...	...	...	...
10834	FR Calculator	FAMILY	4.0	7	2600.0	500+	Free	0
10836	Sya9a Maroc - FR	FAMILY	4.5	38	53000.0	5,000+	Free	0
10837	Fr. Mike Schmitz Audio Teachings	FAMILY	5.0	4	3600.0	100+	Free	0
10839	The SCP Foundation DB fr nn5n	BOOKS_AND_REFERENCE	4.5	114	0.0	1,000+	Free	0
10840	iHoroscope - 2018 Daily Horoscope & Astrology	LIFESTYLE	4.5	398307	19000.0	10,000,000+	Free	0

9360 rows × 13 columns



## 4(II).

```
In [15]: data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 9360 entries, 0 to 10840
Data columns (total 13 columns):
#   Column          Non-Null Count  Dtype
---  -
0   App              9360 non-null   object
1   Category         9360 non-null   object
2   Rating           9360 non-null   float64
3   Reviews          9360 non-null   object
4   Size             9360 non-null   float64
5   Installs         9360 non-null   object
6   Type             9360 non-null   object
7   Price            9360 non-null   object
8   Content Rating   9360 non-null   object
9   Genres           9360 non-null   object
10  Last Updated     9360 non-null   object
11  Current Ver      9360 non-null   object
12  Android Ver      9360 non-null   object
dtypes: float64(2), object(11)
memory usage: 1023.8+ KB
```

```
In [16]: data["Reviews"] = data["Reviews"].astype(float)
```

```
In [17]: data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 9360 entries, 0 to 10840
Data columns (total 13 columns):
#   Column          Non-Null Count  Dtype
---  -
0   App              9360 non-null   object
1   Category         9360 non-null   object
2   Rating           9360 non-null   float64
3   Reviews          9360 non-null   float64
4   Size             9360 non-null   float64
5   Installs         9360 non-null   object
6   Type             9360 non-null   object
7   Price            9360 non-null   object
8   Content Rating   9360 non-null   object
9   Genres           9360 non-null   object
10  Last Updated     9360 non-null   object
11  Current Ver      9360 non-null   object
12  Android Ver      9360 non-null   object
dtypes: float64(3), object(10)
memory usage: 1023.8+ KB
```

## 4(III).

```
In [29]: data["Installs"] = data["Installs"].astype(int)
```

```
In [30]: data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 9360 entries, 0 to 10840
Data columns (total 13 columns):
#   Column              Non-Null Count  Dtype
---  -
0   App                  9360 non-null   object
1   Category              9360 non-null   object
2   Rating                9360 non-null   float64
3   Reviews               9360 non-null   float64
4   Size                  9360 non-null   float64
5   Installs              9360 non-null   int32
6   Type                  9360 non-null   object
7   Price                 9360 non-null   float64
8   Content Rating        9360 non-null   object
9   Genres                 9360 non-null   object
10  Last Updated          9360 non-null   object
11  Current Ver           9360 non-null   object
12  Android Ver           9360 non-null   object
dtypes: float64(4), int32(1), object(8)
memory usage: 987.2+ KB
```

## 4(IV).

```
In [32]: data["Price"] = data["Price"].astype(int)
```

```
In [33]: data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 9360 entries, 0 to 10840
Data columns (total 13 columns):
#   Column          Non-Null Count  Dtype
---  -
0   App              9360 non-null   object
1   Category         9360 non-null   object
2   Rating           9360 non-null   float64
3   Reviews          9360 non-null   float64
4   Size             9360 non-null   float64
5   Installs         9360 non-null   int32
6   Type             9360 non-null   object
7   Price            9360 non-null   int32
8   Content Rating   9360 non-null   object
9   Genres           9360 non-null   object
10  Last Updated     9360 non-null   object
11  Current Ver      9360 non-null   object
12  Android Ver      9360 non-null   object
dtypes: float64(3), int32(2), object(8)
memory usage: 950.6+ KB
```

## 4(V-A).

```
In [34]: data.shape
```

```
Out[34]: (9360, 13)
```

```
In [35]: data.drop(data[(data['Reviews'] < 1) & (data['Reviews'] > 5)].index, inplace
```

```
In [36]: data.shape
```

```
Out[36]: (9360, 13)
```

## 4(V-B).

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

```
Out[37]: (9360, 13)
```

```
In [38]: data.drop(data[data['Installs'] < data['Reviews'] ].index, inplace = True)
```

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

```
Out[39]: (9353, 13)
```



## 4(V-C).

```
In [40]: data.shape
```

```
Out[40]: (9353, 13)
```

```
In [41]: data.drop(data[(data['Type'] == 'Free') & (data['Price'] > 0)].index, inplace
```

```
In [42]: data.shape
```

```
Out[42]: (9353, 13)
```

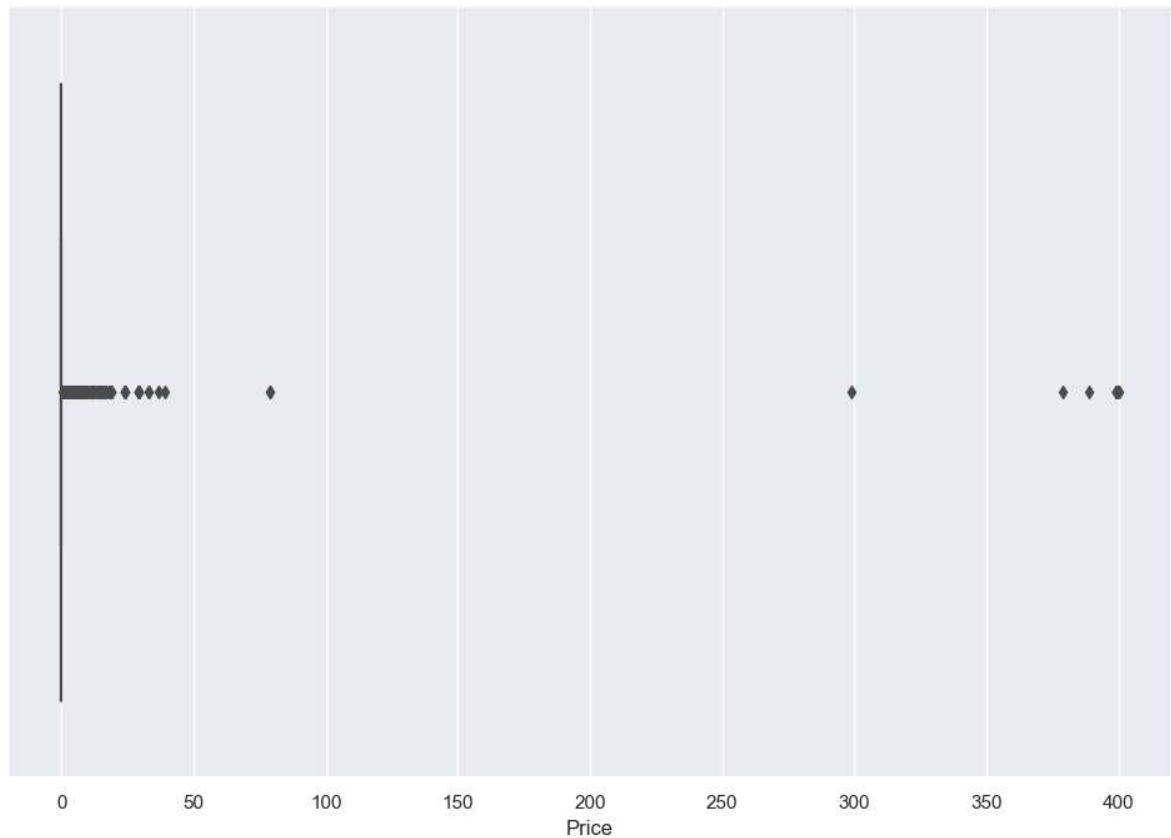
## 5(I).

```
In [43]: sns.set(rc={'figure.figsize':(12,8)})
```

```
In [44]: sns.boxplot(data['Price'])
```

C:\Users\romit\anaconda3\lib\site-packages\seaborn\\_decorators.py:36: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be `data`, and passing other arguments without an explicit keyword will result in an error or misinterpretation.  
warnings.warn(

```
Out[44]: <AxesSubplot:xlabel='Price'>
```

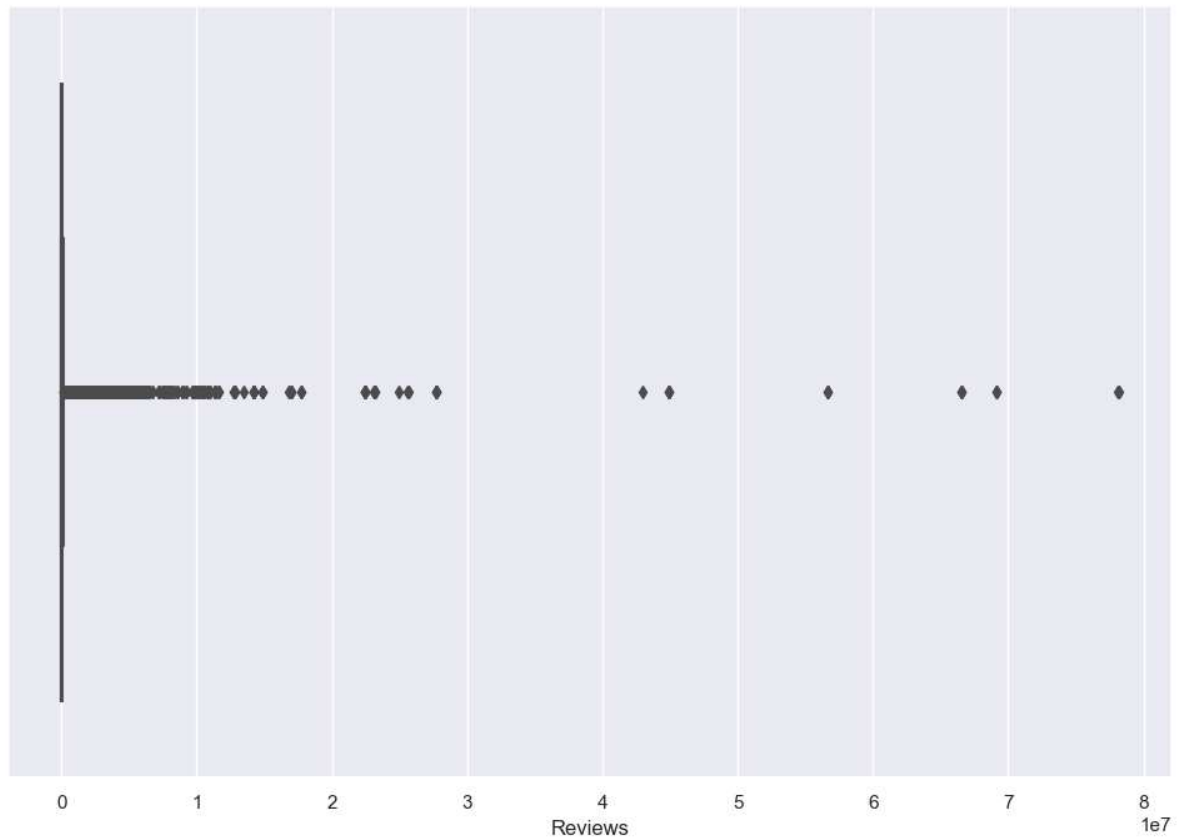


## 5(II).

```
In [45]: sns.boxplot(data['Reviews'])
```

C:\Users\romit\anaconda3\lib\site-packages\seaborn\\_decorators.py:36: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be `data`, and passing other arguments without an explicit keyword will result in an error or misinterpretation.  
warnings.warn(

```
Out[45]: <AxesSubplot:xlabel='Reviews'>
```

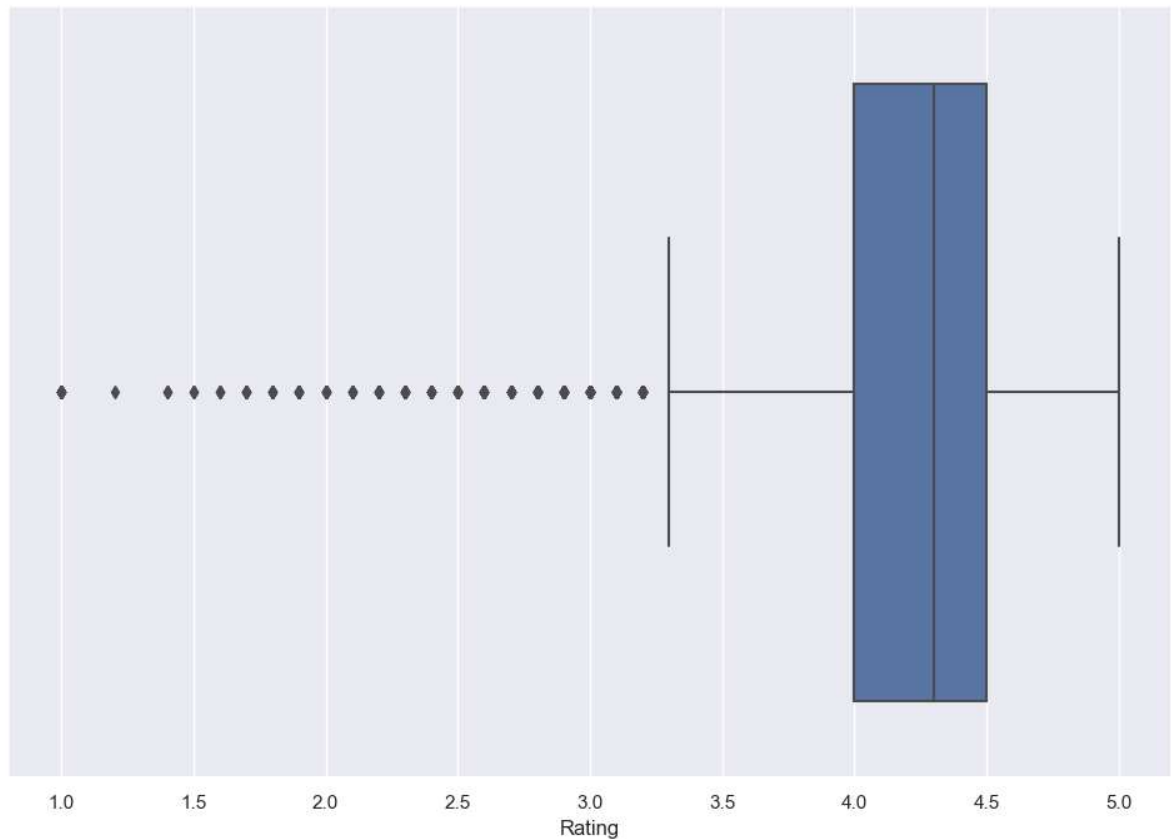


## 5(III).

```
In [46]: sns.boxplot(data['Rating'])
```

C:\Users\romit\anaconda3\lib\site-packages\seaborn\\_decorators.py:36: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be `data`, and passing other arguments without an explicit keyword will result in an error or misinterpretation.  
warnings.warn(

```
Out[46]: <AxesSubplot:xlabel='Rating'>
```

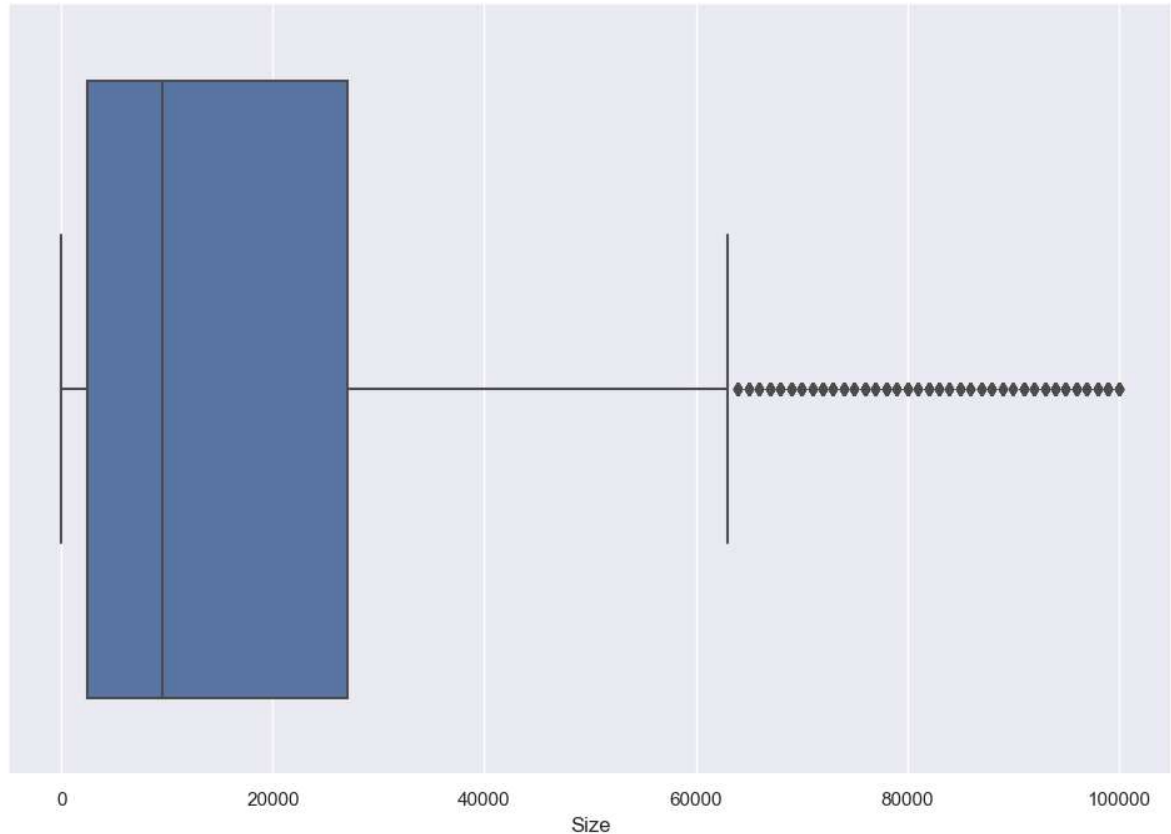


## 5(IV).

```
In [47]: sns.boxplot(data['Size'])
```

C:\Users\romit\anaconda3\lib\site-packages\seaborn\\_decorators.py:36: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be `data`, and passing other arguments without an explicit keyword will result in an error or misinterpretation.  
warnings.warn(

```
Out[47]: <AxesSubplot:xlabel='Size'>
```



## 6(I).

```
In [48]: more = data.apply(lambda x : True  
                           if x['Price'] > 200 else False, axis = 1)
```

```
In [49]: more_count = len(more[more == True].index)
```

```
In [50]: data.shape
```

```
Out[50]: (9353, 13)
```

```
In [51]: data.drop(data[data['Price'] > 200].index, inplace = True)
```

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

```
Out[52]: (9338, 13)
```

## 6(II).

```
In [53]: data.drop(data[data['Reviews'] > 2000000].index, inplace = True)
```

```
In [54]: data.shape
```

```
Out[54]: (8885, 13)
```

## 6(III).

```
In [55]: data.quantile([.1, .25, .5, .70, .90, .95, .99], axis = 0)
```

```
Out[55]:
```

	Rating	Reviews	Size	Installs	Price
<b>0.10</b>	3.5	18.00	0.0	1000.0	0.0
<b>0.25</b>	4.0	159.00	2600.0	10000.0	0.0
<b>0.50</b>	4.3	4290.00	9500.0	500000.0	0.0
<b>0.70</b>	4.5	35930.40	23000.0	1000000.0	0.0
<b>0.90</b>	4.7	296771.00	50000.0	10000000.0	0.0
<b>0.95</b>	4.8	637298.00	68000.0	10000000.0	1.0
<b>0.99</b>	5.0	1462800.88	95000.0	100000000.0	7.0

```
In [56]: # dropping more than 10000000 Installs value  
data.drop(data[data['Installs'] > 10000000].index, inplace = True)
```

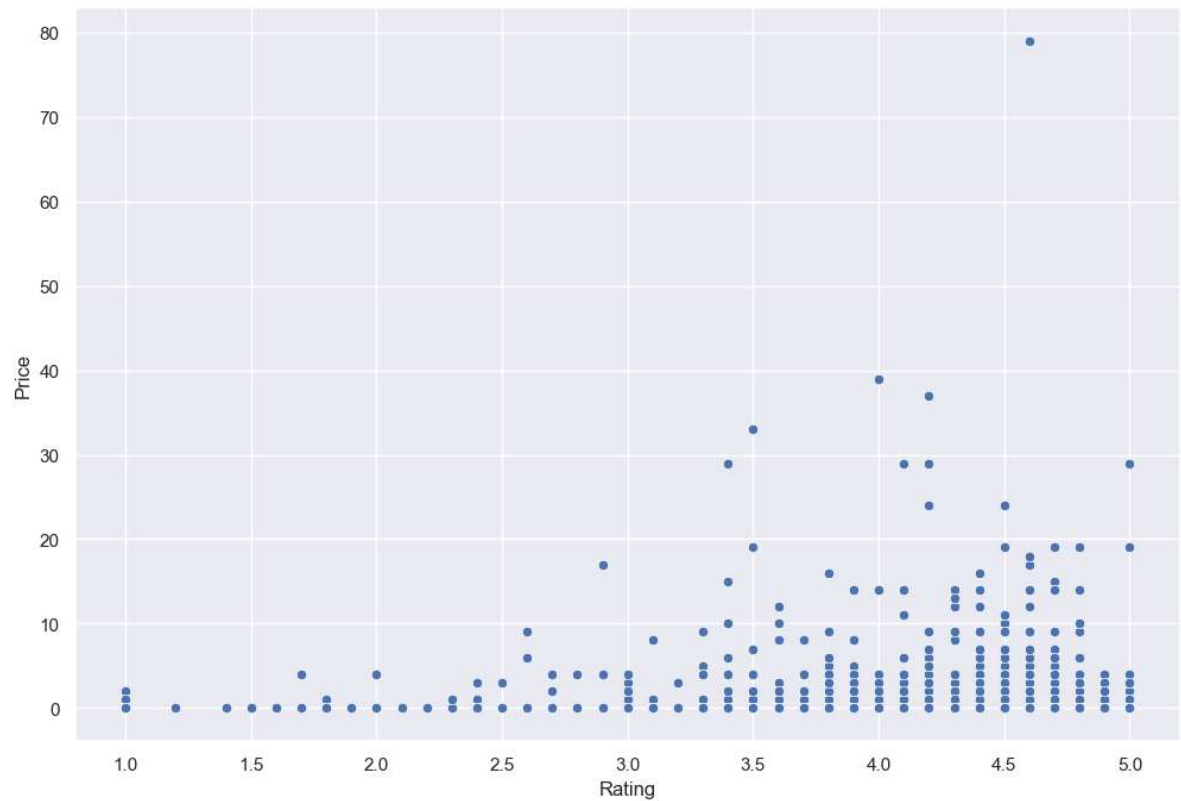
```
In [57]: data.shape
```

```
Out[57]: (8496, 13)
```

7(l).

```
In [58]: sns.scatterplot(x='Rating',y='Price',data=data)
```

```
Out[58]: <AxesSubplot:xlabel='Rating', ylabel='Price'>
```

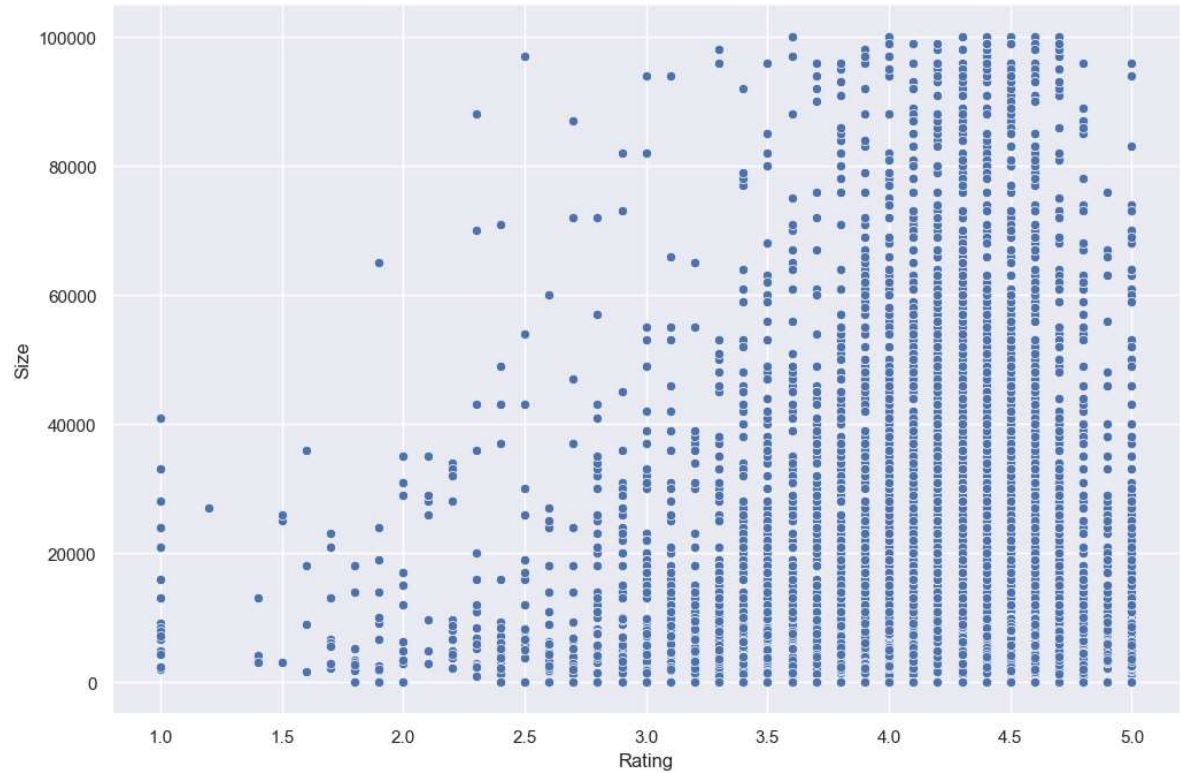


Yes, Paid apps are higher ratings compared to free apps.

## 7(II).

```
In [59]: sns.scatterplot(x='Rating',y='Size',data=data)
```

```
Out[59]: <AxesSubplot:xlabel='Rating', ylabel='Size'>
```



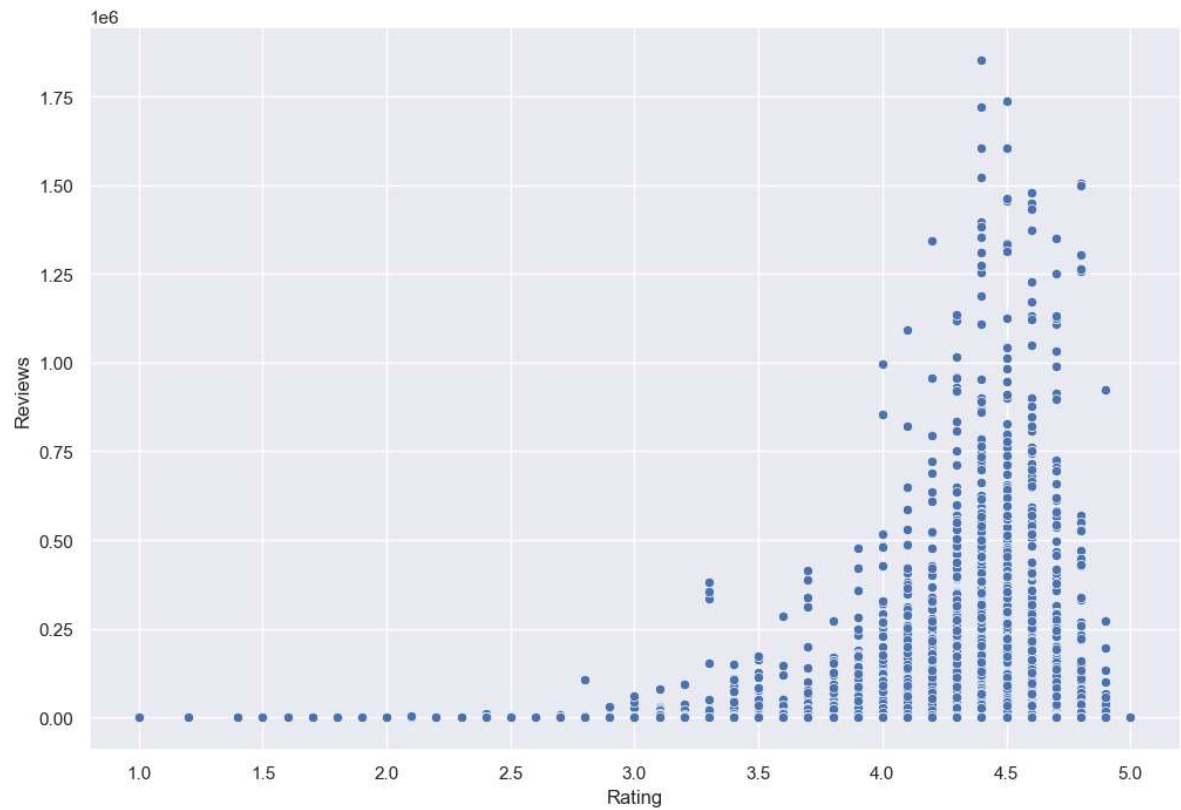
Yes it is clear that heavier apps are rated better.



## 7(III).

```
In [60]: sns.scatterplot(x='Rating',y='Reviews',data=data)
```

```
Out[60]: <AxesSubplot:xlabel='Rating', ylabel='Reviews'>
```

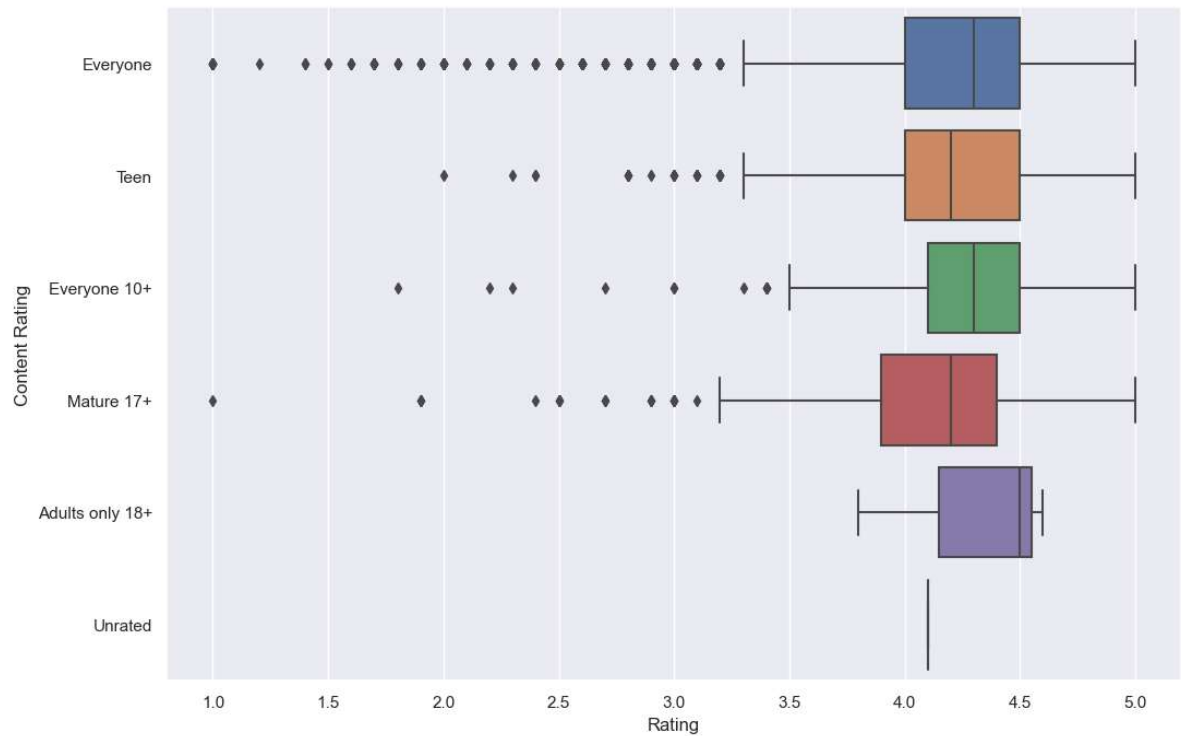


It is clear that more reviews makes app rating better.

## 7(IV).

```
In [61]: sns.boxplot(x="Rating", y="Content Rating", data=data)
```

```
Out[61]: <AxesSubplot:xlabel='Rating', ylabel='Content Rating'>
```

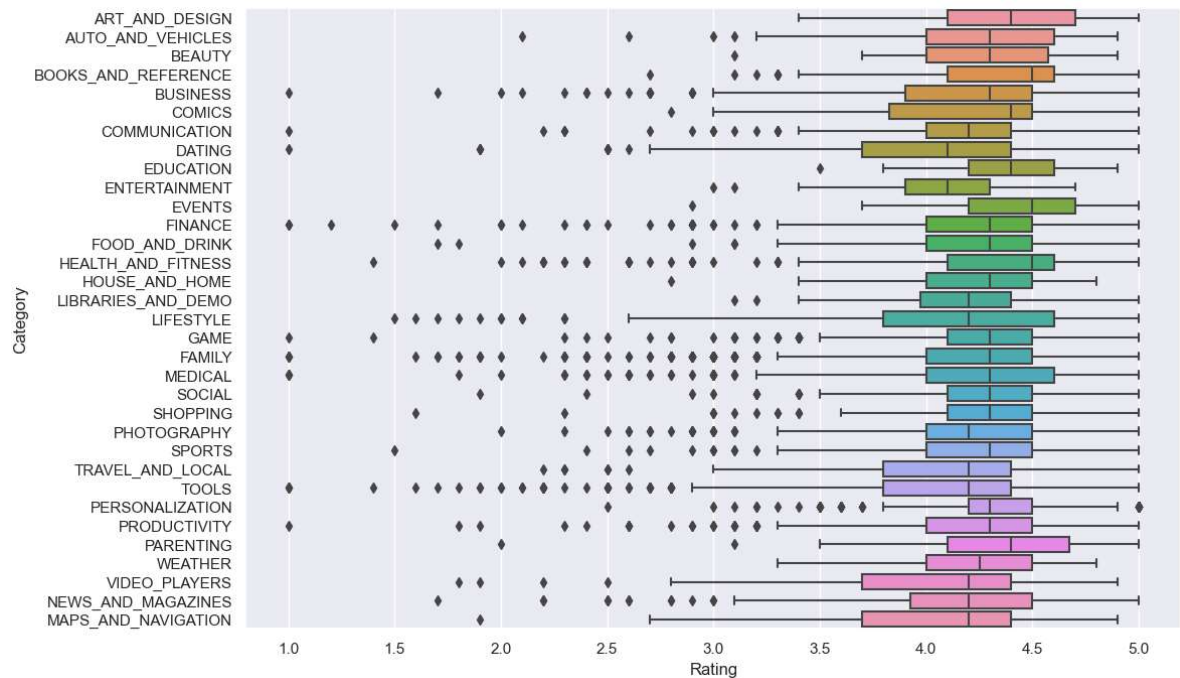


Apps which are for everyone has more bad ratings compare to other sections as it has so much outliers value, while 18+ apps have better ratings.

## 7(V).

```
In [62]: sns.boxplot(x="Rating", y="Category", data=data)
```

```
Out[62]: <AxesSubplot:xlabel='Rating', ylabel='Category'>
```



Events category has best ratings compare to others.


## 8(I).

```
In [63]: inp1 = data
```

```
In [64]: inp1.head()
```

Out[64]:

	App	Category	Rating	Reviews	Size	Installs	Type	Price	Content Rating	
0	Photo Editor & Candy Camera & Grid & ScrapBook	ART_AND_DESIGN	4.1	159.0	19000.0	10000	Free	0	Everyone	
1	Coloring book moana	ART_AND_DESIGN	3.9	967.0	14000.0	500000	Free	0	Everyone	C
2	U Launcher Lite – FREE Live Cool Themes, Hide ...	ART_AND_DESIGN	4.7	87510.0	8700.0	5000000	Free	0	Everyone	
4	Pixel Draw - Number Art Coloring Book	ART_AND_DESIGN	4.3	967.0	2800.0	100000	Free	0	Everyone	De
5	Paper flowers instructions	ART_AND_DESIGN	4.4	167.0	5600.0	50000	Free	0	Everyone	



```
In [65]: inp1.skew()
```

```
C:\Users\romit\AppData\Local\Temp\ipykernel_21424\3545313420.py:1: FutureWarning: Dropping of nuisance columns in DataFrame reductions (with 'numeric_only=None') is deprecated; in a future version this will raise TypeError. Select only valid columns before calling the reduction.
inp1.skew()
```

Out[65]:

Rating	-1.749753
Reviews	4.576494
Size	1.655917
Installs	1.543697
Price	18.074542

dtype: float64

```
In [66]: reviewskew = np.log1p(inp1['Reviews'])
inp1['Reviews'] = reviewskew
```

```
In [67]: reviewskew.skew()
```

Out[67]: -0.20039949659264134

```
In [68]: installsskew = np.log1p(inp1['Installs'])
inp1['Installs']
```

```
Out[68]: 0          10000
1          500000
2          5000000
4          100000
5           50000
...
10834         500
10836         5000
10837          100
10839         1000
10840       10000000
Name: Installs, Length: 8496, dtype: int32
```

```
In [69]: installsskew.skew()
```

```
Out[69]: -0.5097286542754812
```

```
In [70]: inp1.head()
```

```
Out[70]:
```

	App	Category	Rating	Reviews	Size	Installs	Type	Price	Content Rating
0	Photo Editor & Candy Camera & Grid & ScrapBook	ART_AND_DESIGN	4.1	5.075174	19000.0	10000	Free	0	Everyone
1	Coloring book moana	ART_AND_DESIGN	3.9	6.875232	14000.0	500000	Free	0	Everyone
2	U Launcher Lite – FREE Live Cool Themes, Hide ...	ART_AND_DESIGN	4.7	11.379520	8700.0	5000000	Free	0	Everyone
4	Pixel Draw - Number Art Coloring Book	ART_AND_DESIGN	4.3	6.875232	2800.0	100000	Free	0	Everyone
5	Paper flowers instructions	ART_AND_DESIGN	4.4	5.123964	5600.0	50000	Free	0	Everyone

## 8(II).

```
In [71]: inp1.drop(["Last Updated", "Current Ver", "Android Ver", "App", "Type"], axis=1, inplace=True)
```

```
In [72]: inp1.head()
```

Out[72]:

	Category	Rating	Reviews	Size	Installs	Price	Content Rating	Genres
0	ART_AND_DESIGN	4.1	5.075174	19000.0	10000	0	Everyone	Art & Design
1	ART_AND_DESIGN	3.9	6.875232	14000.0	500000	0	Everyone	Art & Design;Pretend Play
2	ART_AND_DESIGN	4.7	11.379520	8700.0	5000000	0	Everyone	Art & Design
4	ART_AND_DESIGN	4.3	6.875232	2800.0	100000	0	Everyone	Art & Design;Creativity
5	ART_AND_DESIGN	4.4	5.123964	5600.0	50000	0	Everyone	Art & Design

```
In [73]: inp1.shape
```

Out[73]: (8496, 8)

## 8(III).

```
In [74]: inp2 = inp1
```

```
In [75]: inp2.head()
```

Out[75]:

	Category	Rating	Reviews	Size	Installs	Price	Content Rating	Genres
0	ART_AND_DESIGN	4.1	5.075174	19000.0	10000	0	Everyone	Art & Design
1	ART_AND_DESIGN	3.9	6.875232	14000.0	500000	0	Everyone	Art & Design;Pretend Play
2	ART_AND_DESIGN	4.7	11.379520	8700.0	5000000	0	Everyone	Art & Design
4	ART_AND_DESIGN	4.3	6.875232	2800.0	100000	0	Everyone	Art & Design;Creativity
5	ART_AND_DESIGN	4.4	5.123964	5600.0	50000	0	Everyone	Art & Design

## Let's apply Dummy EnCoding on Column "Category"

```
In [76]: #get unique values in Column "Category"
inp2.Category.unique()
```

```
Out[76]: array(['ART_AND_DESIGN', 'AUTO_AND_VEHICLES', 'BEAUTY',
                'BOOKS_AND_REFERENCE', 'BUSINESS', 'COMICS', 'COMMUNICATION',
                'DATING', 'EDUCATION', 'ENTERTAINMENT', 'EVENTS', 'FINANCE',
                'FOOD_AND_DRINK', 'HEALTH_AND_FITNESS', 'HOUSE_AND_HOME',
                'LIBRARIES_AND_DEMO', 'LIFESTYLE', 'GAME', 'FAMILY', 'MEDICAL',
                'SOCIAL', 'SHOPPING', 'PHOTOGRAPHY', 'SPORTS', 'TRAVEL_AND_LOCAL',
                'TOOLS', 'PERSONALIZATION', 'PRODUCTIVITY', 'PARENTING', 'WEATHER',
                'VIDEO_PLAYERS', 'NEWS_AND_MAGAZINES', 'MAPS_AND_NAVIGATION'],
              dtype=object)
```

```
In [77]: inp2.Category = pd.Categorical(inp2.Category)

x = inp2[['Category']]
del inp2['Category']

dummies = pd.get_dummies(x, prefix = 'Category')
inp2 = pd.concat([inp2,dummies], axis=1)
inp2.head()
```

```
Out[77]:
```

	Rating	Reviews	Size	Installs	Price	Content Rating	Genres	Category_ART_AND_DE
0	4.1	5.075174	19000.0	10000	0	Everyone	Art & Design	
1	3.9	6.875232	14000.0	500000	0	Everyone	Design;Pretend Play	Art &
2	4.7	11.379520	8700.0	5000000	0	Everyone	Art & Design	
4	4.3	6.875232	2800.0	100000	0	Everyone	Design;Creativity	Art &
5	4.4	5.123964	5600.0	50000	0	Everyone	Art & Design	

5 rows × 40 columns

```
In [78]: inp2.shape
```

```
Out[78]: (8496, 40)
```

## Let's apply Dummy EnCoding on Column "Genres"

```
In [79]: #get unique values in Column "Genres"
inp2["Genres"].unique()
```

```
Out[79]: array(['Art & Design', 'Art & Design;Pretend Play',
               'Art & Design;Creativity', 'Auto & Vehicles', 'Beauty',
               'Books & Reference', 'Business', 'Comics', 'Comics;Creativity',
               'Communication', 'Dating', 'Education', 'Education;Creativity',
               'Education;Education', 'Education;Music & Video',
               'Education;Action & Adventure', 'Education;Pretend Play',
               'Education;Brain Games', 'Entertainment',
               'Entertainment;Brain Games', 'Entertainment;Creativity',
               'Entertainment;Music & Video', 'Events', 'Finance', 'Food & Drink',
               'Health & Fitness', 'House & Home', 'Libraries & Demo',
               'Lifestyle', 'Lifestyle;Pretend Play', 'Card', 'Casual', 'Puzzle',
               'Action', 'Arcade', 'Word', 'Racing', 'Casual;Creativity',
               'Sports', 'Board', 'Simulation', 'Role Playing', 'Adventure',
               'Strategy', 'Simulation;Education', 'Action;Action & Adventure',
               'Trivia', 'Casual;Brain Games', 'Simulation;Action & Adventure',
               'Educational;Creativity', 'Puzzle;Brain Games',
               'Educational;Education', 'Card;Brain Games',
               'Educational;Brain Games', 'Educational;Pretend Play',
               'Casual;Action & Adventure', 'Entertainment;Education',
               'Casual;Education', 'Casual;Pretend Play', 'Music;Music & Video',
               'Racing;Action & Adventure', 'Arcade;Pretend Play',
               'Adventure;Action & Adventure', 'Role Playing;Action & Adventure',
               'Simulation;Pretend Play', 'Puzzle;Creativity',
               'Sports;Action & Adventure', 'Educational;Action & Adventure',
               'Arcade;Action & Adventure', 'Entertainment;Action & Adventure',
               'Puzzle;Action & Adventure', 'Strategy;Action & Adventure',
               'Music & Audio;Music & Video', 'Health & Fitness;Education',
               'Adventure;Education', 'Board;Brain Games',
               'Board;Action & Adventure', 'Board;Pretend Play',
               'Casual;Music & Video', 'Role Playing;Pretend Play',
               'Entertainment;Pretend Play', 'Video Players & Editors;Creativity',
               'Card;Action & Adventure', 'Medical', 'Social', 'Shopping',
               'Photography', 'Travel & Local',
               'Travel & Local;Action & Adventure', 'Tools', 'Tools;Education',
               'Personalization', 'Productivity', 'Parenting',
               'Parenting;Music & Video', 'Parenting;Brain Games',
               'Parenting;Education', 'Weather', 'Video Players & Editors',
               'Video Players & Editors;Music & Video', 'News & Magazines',
               'Maps & Navigation', 'Health & Fitness;Action & Adventure',
               'Music', 'Educational', 'Casino', 'Adventure;Brain Games',
               'Lifestyle;Education', 'Books & Reference;Education',
               'Puzzle;Education', 'Role Playing;Brain Games',
               'Strategy;Education', 'Racing;Pretend Play',
               'Communication;Creativity', 'Strategy;Creativity'], dtype=object)
```

=> Since, There are too many categories under Genres. Hence, we will try to reduce some categories which have very few samples under them and put them under one new common category i.e. "Other".



```
In [80]: lists = []
        for i in inp2.Genres.value_counts().index:
            if inp2.Genres.value_counts()[i]<20:
                lists.append(i)
        inp2.Genres = ['Other' if i in lists else i for i in inp2.Genres]
```

```
In [81]: inp2["Genres"].unique()
```

```
Out[81]: array(['Art & Design', 'Other', 'Auto & Vehicles', 'Beauty',
                'Books & Reference', 'Business', 'Comics', 'Communication',
                'Dating', 'Education', 'Education;Education',
                'Education;Pretend Play', 'Entertainment',
                'Entertainment;Music & Video', 'Events', 'Finance', 'Food & Drink',
                'Health & Fitness', 'House & Home', 'Libraries & Demo',
                'Lifestyle', 'Card', 'Casual', 'Puzzle', 'Action', 'Arcade',
                'Word', 'Racing', 'Sports', 'Board', 'Simulation', 'Role Playing',
                'Adventure', 'Strategy', 'Trivia', 'Educational;Education',
                'Casual;Pretend Play', 'Medical', 'Social', 'Shopping',
                'Photography', 'Travel & Local', 'Tools', 'Personalization',
                'Productivity', 'Parenting', 'Weather', 'Video Players & Editors',
                'News & Magazines', 'Maps & Navigation', 'Educational', 'Casino'],
              dtype=object)
```

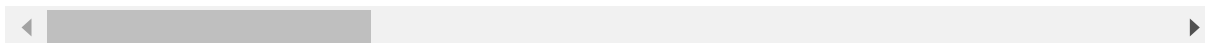
```
In [82]: inp2.Genres = pd.Categorical(inp2['Genres'])
        x = inp2[["Genres"]]
        del inp2['Genres']
        dummies = pd.get_dummies(x, prefix = 'Genres')
        inp2 = pd.concat([inp2,dummies], axis=1)
```

```
In [83]: inp2.head()
```

```
Out[83]:
```

	Rating	Reviews	Size	Installs	Price	Content Rating	Category_ART_AND_DESIGN	Category_
0	4.1	5.075174	19000.0	10000	0	Everyone	1	
1	3.9	6.875232	14000.0	500000	0	Everyone	1	
2	4.7	11.379520	8700.0	5000000	0	Everyone	1	
4	4.3	6.875232	2800.0	100000	0	Everyone	1	
5	4.4	5.123964	5600.0	50000	0	Everyone	1	

5 rows × 91 columns



```
In [84]: inp2.shape
```

```
Out[84]: (8496, 91)
```

## Let's apply Dummy EnCoding on Column "Content Rating"

```
In [85]: #get unique values in Column "Content Rating"
inp2["Content Rating"].unique()
```

```
Out[85]: array(['Everyone', 'Teen', 'Everyone 10+', 'Mature 17+',
               'Adults only 18+', 'Unrated'], dtype=object)
```

```
In [86]: inp2['Content Rating'] = pd.Categorical(inp2['Content Rating'])

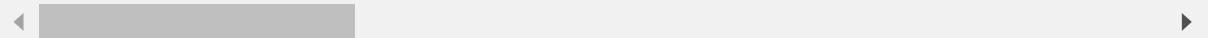
x = inp2[['Content Rating']]
del inp2['Content Rating']

dummies = pd.get_dummies(x, prefix = 'Content Rating')
inp2 = pd.concat([inp2,dummies], axis=1)
inp2.head()
```

```
Out[86]:
```

	Rating	Reviews	Size	Installs	Price	Category_ART_AND_DESIGN	Category_AUTO_AND_DESIGN
0	4.1	5.075174	19000.0	10000	0	1	
1	3.9	6.875232	14000.0	500000	0	1	
2	4.7	11.379520	8700.0	5000000	0	1	
4	4.3	6.875232	2800.0	100000	0	1	
5	4.4	5.123964	5600.0	50000	0	1	

5 rows × 96 columns



```
In [87]: inp2.shape
```

```
Out[87]: (8496, 96)
```

## 9 and 10.

```
In [88]: from sklearn.model_selection import train_test_split as tts
from sklearn.linear_model import LinearRegression as LR
from sklearn.metrics import mean_squared_error as mse
```

```
In [89]: d1 = inp2
X = d1.drop('Rating',axis=1)
y = d1['Rating']

Xtrain, Xtest, ytrain, ytest = tts(X,y, test_size=0.3, random_state=5)
```

## 11.

```
In [90]: reg_all = LR()
reg_all.fit(Xtrain,ytrain)
```

Out[90]: LinearRegression()

```
In [91]: R2_train = round(reg_all.score(Xtrain,ytrain),3)
print("The R2 value of the Training Set is : {}".format(R2_train))
```

The R2 value of the Training Set is : 0.074

## 12.

```
In [92]: R2_test = round(reg_all.score(Xtest,ytest),3)
print("The R2 value of the Testing Set is : {}".format(R2_test))
```

The R2 value of the Testing Set is : 0.063