

1

```
In [46]: import pandas as pd
```

```
In [47]: import numpy as np
```

```
In [48]: import matplotlib.pyplot as plt
```

```
In [49]: import seaborn as sns
```

```
In [50]: import statistics as stc
```

```
In [51]: df=pd.read_csv("googleplaystore.csv")
```

```
In [52]: df.head(5)
```

```
Out[52]:
```

	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 [54]: df.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 [55]: df.shape
```

```
Out[55]: (10841, 13)
```

2

```
In [56]: df.isnull().any()
```

```

Out[56]: 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 [57]: df.isna().sum()
```

```

Out[57]: 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 [58]: df.dropna(inplace=True)
```

```
In [59]: df.isnull().any()
```

```
Out[59]: 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 [60]: df.isna().sum()
```

```
Out[60]: App          0
Category          0
Rating            0
Reviews           0
Size              0
Installs          0
Type              0
Price             0
Content Rating    0
Genres            0
Last Updated      0
Current Ver       0
Android Ver       0
dtype: int64
```

```
In [61]: df.shape
```

```
Out[61]: (9360, 13)
```

4(I)

```
In [62]: df.columns
```

```
Out[62]: Index(['App', 'Category', 'Rating', 'Reviews', 'Size', 'Installs', 'Type',
              'Price', 'Content Rating', 'Genres', 'Last Updated', 'Current Ver',
              'Android Ver'],
              dtype='object')
```

```
In [63]: df["Size"].unique()
```

```
Out[63]: array(['19M', '14M', '8.7M', '25M', '2.8M', '5.6M', '29M', '33M', '3.1M',
               '28M', '12M', '20M', '21M', '37M', '5.5M', '17M', '39M', '31M',
               '4.2M', '23M', '6.0M', '6.1M', '4.6M', '9.2M', '5.2M', '11M',
               '24M', 'Varies with device', '9.4M', '15M', '10M', '1.2M', '26M',
               '8.0M', '7.9M', '56M', '57M', '35M', '54M', '201k', '3.6M', '5.7M',
               '8.6M', '2.4M', '27M', '2.7M', '2.5M', '7.0M', '16M', '3.4M',
               '8.9M', '3.9M', '2.9M', '38M', '32M', '5.4M', '18M', '1.1M',
               '2.2M', '4.5M', '9.8M', '52M', '9.0M', '6.7M', '30M', '2.6M',
               '7.1M', '22M', '6.4M', '3.2M', '8.2M', '4.9M', '9.5M', '5.0M',
               '5.9M', '13M', '73M', '6.8M', '3.5M', '4.0M', '2.3M', '2.1M',
               '42M', '9.1M', '55M', '23k', '7.3M', '6.5M', '1.5M', '7.5M', '51M',
               '41M', '48M', '8.5M', '46M', '8.3M', '4.3M', '4.7M', '3.3M', '40M',
               '7.8M', '8.8M', '6.6M', '5.1M', '61M', '66M', '79k', '8.4M',
               '3.7M', '118k', '44M', '695k', '1.6M', '6.2M', '53M', '1.4M',
               '3.0M', '7.2M', '5.8M', '3.8M', '9.6M', '45M', '63M', '49M', '77M',
               '4.4M', '70M', '9.3M', '8.1M', '36M', '6.9M', '7.4M', '84M', '97M',
               '2.0M', '1.9M', '1.8M', '5.3M', '47M', '556k', '526k', '76M',
               '7.6M', '59M', '9.7M', '78M', '72M', '43M', '7.7M', '6.3M', '334k',
               '93M', '65M', '79M', '100M', '58M', '50M', '68M', '64M', '34M',
               '67M', '60M', '94M', '9.9M', '232k', '99M', '624k', '95M', '8.5k',
               '41k', '292k', '80M', '1.7M', '10.0M', '74M', '62M', '69M', '75M',
               '98M', '85M', '82M', '96M', '87M', '71M', '86M', '91M', '81M',
               '92M', '83M', '88M', '704k', '862k', '899k', '378k', '4.8M',
               '266k', '375k', '1.3M', '975k', '980k', '4.1M', '89M', '696k',
               '544k', '525k', '920k', '779k', '853k', '720k', '713k', '772k',
               '318k', '58k', '241k', '196k', '857k', '51k', '953k', '865k',
               '251k', '930k', '540k', '313k', '746k', '203k', '26k', '314k',
               '239k', '371k', '220k', '730k', '756k', '91k', '293k', '17k',
               '74k', '14k', '317k', '78k', '924k', '818k', '81k', '939k', '169k',
               '45k', '965k', '90M', '545k', '61k', '283k', '655k', '714k', '93k',
               '872k', '121k', '322k', '976k', '206k', '954k', '444k', '717k',
               '210k', '609k', '308k', '306k', '175k', '350k', '383k', '454k',
               '1.0M', '70k', '812k', '442k', '842k', '417k', '412k', '459k',
               '478k', '335k', '782k', '721k', '430k', '429k', '192k', '460k',
               '728k', '496k', '816k', '414k', '506k', '887k', '613k', '778k',
               '683k', '592k', '186k', '840k', '647k', '373k', '437k', '598k',
               '716k', '585k', '982k', '219k', '55k', '323k', '691k', '511k',
               '951k', '963k', '25k', '554k', '351k', '27k', '82k', '208k',
               '551k', '29k', '103k', '116k', '153k', '209k', '499k', '173k',
               '597k', '809k', '122k', '411k', '400k', '801k', '787k', '50k',
               '643k', '986k', '516k', '837k', '780k', '20k', '498k', '600k',
               '656k', '221k', '228k', '176k', '34k', '259k', '164k', '458k',
               '629k', '28k', '288k', '775k', '785k', '636k', '916k', '994k',
               '309k', '485k', '914k', '903k', '608k', '500k', '54k', '562k',
               '847k', '948k', '811k', '270k', '48k', '523k', '784k', '280k',
               '24k', '892k', '154k', '18k', '33k', '860k', '364k', '387k',
               '626k', '161k', '879k', '39k', '170k', '141k', '160k', '144k',
               '143k', '190k', '376k', '193k', '473k', '246k', '73k', '253k',
               '957k', '420k', '72k', '404k', '470k', '226k', '240k', '89k',
```

```
'234k', '257k', '861k', '467k', '676k', '552k', '582k', '619k'],  
dtype=object)
```

```
In [64]: df["Size"] = [ float(i.split('M')[0]) if 'M' in i else float(0) for i in df["Size"]]
```

```
In [65]: df.head(5)
```

Out[65]:	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
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

```
In [66]: df["Size"] = 1000 * df["Size"]
```

```
In [67]: df.head(5)
```

Out[67]:	App	Category	Rating	Reviews	Size	Installs	Type	Price	Cont Rat
0	Photo Editor & Candy Camera & Grid & ScrapBook	ART_AND_DESIGN	4.1	159	19000.0	10,000+	Free	0	Every
1	Coloring book moana	ART_AND_DESIGN	3.9	967	14000.0	500,000+	Free	0	Every

2	U Launcher Lite – FREE Live Cool Themes, Hide ...	ART_AND_DESIGN	4.7	87510	8700.0	5,000,000+	Free	0	Every
3	Sketch - Draw & Paint	ART_AND_DESIGN	4.5	215644	25000.0	50,000,000+	Free	0	T
4	Pixel Draw - Number Art Coloring Book	ART_AND_DESIGN	4.3	967	2800.0	100,000+	Free	0	Every

4(II)

In [68]: `df.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 [69]: `df["Reviews"]=df["Reviews"].astype(float)`

In [70]: `df.dtypes`

```
Out[70]: App              object
Category         object
Rating           float64
Reviews          float64
Size             float64
Installs         object
Type             object
```

```
Price                object
Content Rating       object
Genres               object
Last Updated         object
Current Ver          object
Android Ver          object
dtype: object
```

4(III).

```
In [71]: df["Installs"] = [ float(i.replace('+','').replace(',',' ')) if '+' in i or ',' in i
```

```
In [72]: df.head(5)
```

```
Out[72]:
```

	App	Category	Rating	Reviews	Size	Installs	Type	Price	Conte Rati
0	Photo Editor & Candy Camera & Grid & ScrapBook	ART_AND_DESIGN	4.1	159.0	19000.0	10000.0	Free	0	Everyc
1	Coloring book moana	ART_AND_DESIGN	3.9	967.0	14000.0	500000.0	Free	0	Everyc
2	U Launcher Lite – FREE Live Cool Themes, Hide ...	ART_AND_DESIGN	4.7	87510.0	8700.0	5000000.0	Free	0	Everyc
3	Sketch - Draw & Paint	ART_AND_DESIGN	4.5	215644.0	25000.0	50000000.0	Free	0	Te
4	Pixel Draw - Number Art Coloring Book	ART_AND_DESIGN	4.3	967.0	2800.0	100000.0	Free	0	Everyc

```
In [74]: df["Installs"]=df["Installs"].astype(int)
```

```
In [75]: df.dtypes
```

```
Out[75]: App                object
Category                object
Rating                  float64
Reviews                 float64
```

```

Size          float64
Installs      int32
Type          object
Price         object
Content Rating object
Genres        object
Last Updated  object
Current Ver   object
Android Ver   object
dtype: object

```

4(IV).

```
In [76]: df['Price'] = [ float(i.split('$')[1]) if '$' in i else float(0) for i in df['Price'] ]
```

```
In [77]: df.head()
```

```
Out[77]:
```

	App	Category	Rating	Reviews	Size	Installs	Type	Price	Conter Ratin
0	Photo Editor & Candy Camera & Grid & ScrapBook	ART_AND_DESIGN	4.1	159.0	19000.0	10000	Free	0.0	Everyor
1	Coloring book moana	ART_AND_DESIGN	3.9	967.0	14000.0	500000	Free	0.0	Everyor
2	U Launcher Lite – FREE Live Cool Themes, Hide ...	ART_AND_DESIGN	4.7	87510.0	8700.0	5000000	Free	0.0	Everyor
3	Sketch - Draw & Paint	ART_AND_DESIGN	4.5	215644.0	25000.0	50000000	Free	0.0	Tee
4	Pixel Draw - Number Art Coloring Book	ART_AND_DESIGN	4.3	967.0	2800.0	100000	Free	0.0	Everyor

```
In [78]: df['Price']=df['Price'].astype(int)
```

```
In [79]: df.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 [80]: df.shape
```

```
Out[80]: (9360, 13)
```

```
In [81]: df[(df['Rating']<=5) &(df['Rating']>=1)]
```

```
Out[81]:
```

	App	Category	Rating	Reviews	Size	Installs	Type	Pr
0	Photo Editor & Candy Camera & Grid & ScrapBook	ART_AND_DESIGN	4.1	159.0	19000.0	10000	Free	
1	Coloring book moana	ART_AND_DESIGN	3.9	967.0	14000.0	500000	Free	
2	U Launcher Lite – FREE Live Cool Themes, Hide ...	ART_AND_DESIGN	4.7	87510.0	8700.0	5000000	Free	
3	Sketch - Draw & Paint	ART_AND_DESIGN	4.5	215644.0	25000.0	50000000	Free	
4	Pixel Draw - Number Art	ART_AND_DESIGN	4.3	967.0	2800.0	100000	Free	

	Coloring Book						
...
10834	FR Calculator	FAMILY	4.0	7.0	2600.0	500	Free
10836	Sya9a Maroc - FR	FAMILY	4.5	38.0	53000.0	5000	Free
10837	Fr. Mike Schmitz Audio Teachings	FAMILY	5.0	4.0	3600.0	100	Free
10839	The SCP Foundation DB fr nn5n	BOOKS_AND_REFERENCE	4.5	114.0	0.0	1000	Free
10840	iHoroscope - 2018 Daily Horoscope & Astrology	LIFESTYLE	4.5	398307.0	19000.0	10000000	Free

9360 rows × 13 columns

In [82]: `df.shape`

Out[82]: (9360, 13)

4(V-B).

In [83]: `df.shape`

Out[83]: (9360, 13)

In [84]: `df.drop(df.index[df['Reviews'] > df['Installs']],axis=0,inplace = True)`

In [85]: `df.shape`

Out[85]: (9353, 13)

4(V-C).

In [86]: `df.shape`

Out[86]: (9353, 13)

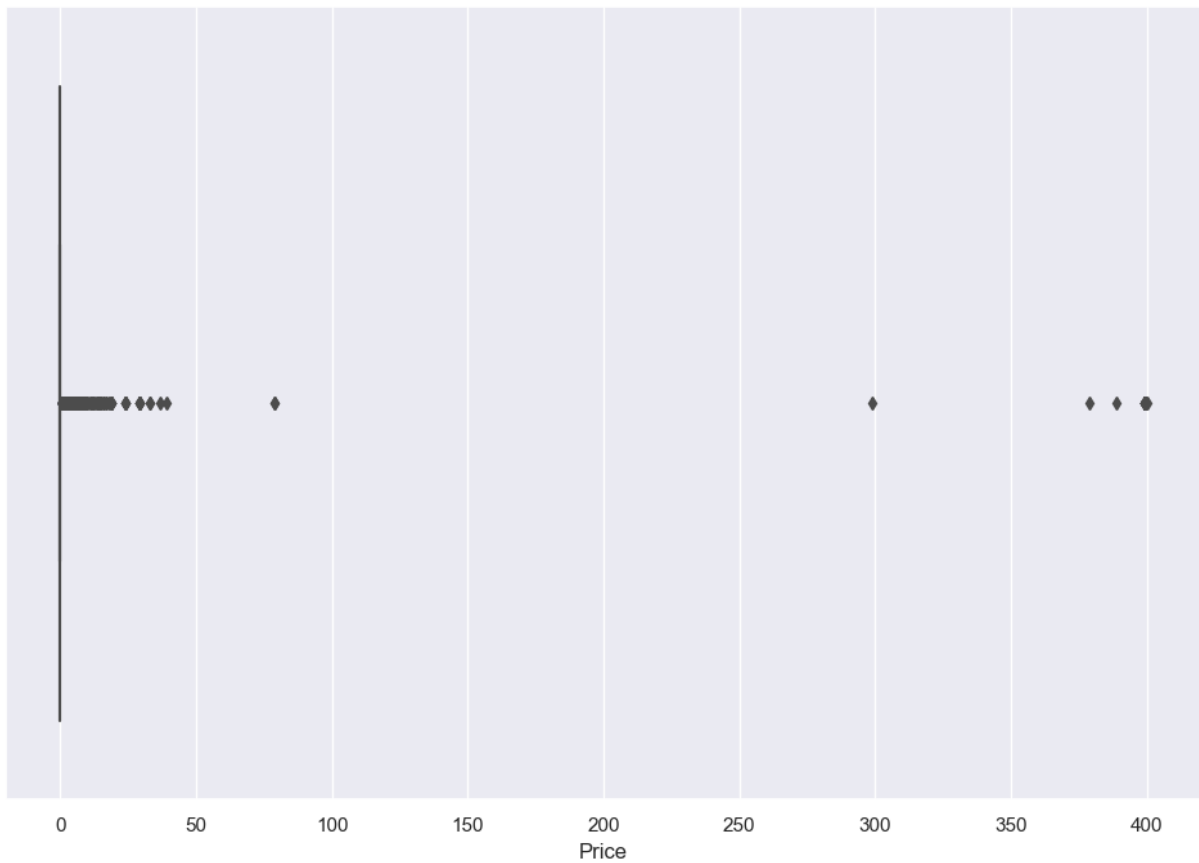
```
In [87]: index_free_and_price_not_0=df.index[((df.Type=='Free')&(df.Price>0))]  
len(index_free_and_price_not_0)  
df.drop(index_free_and_price_not_0,axis=0,inplace=True)
```

```
In [88]: df.shape
```

```
Out[88]: (9353, 13)
```

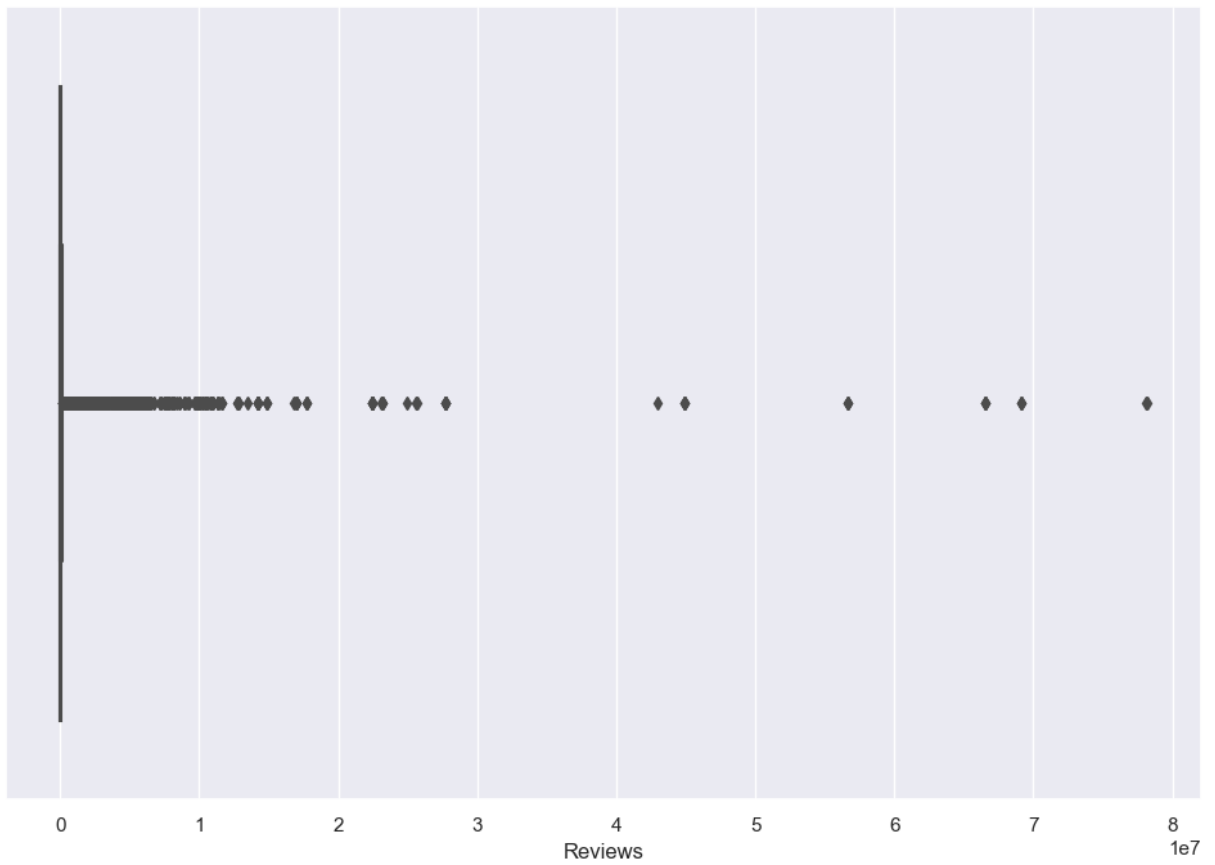
5(I)

```
In [92]: bx_price = sns.boxplot(x='Price',data=df)
```



5(II).

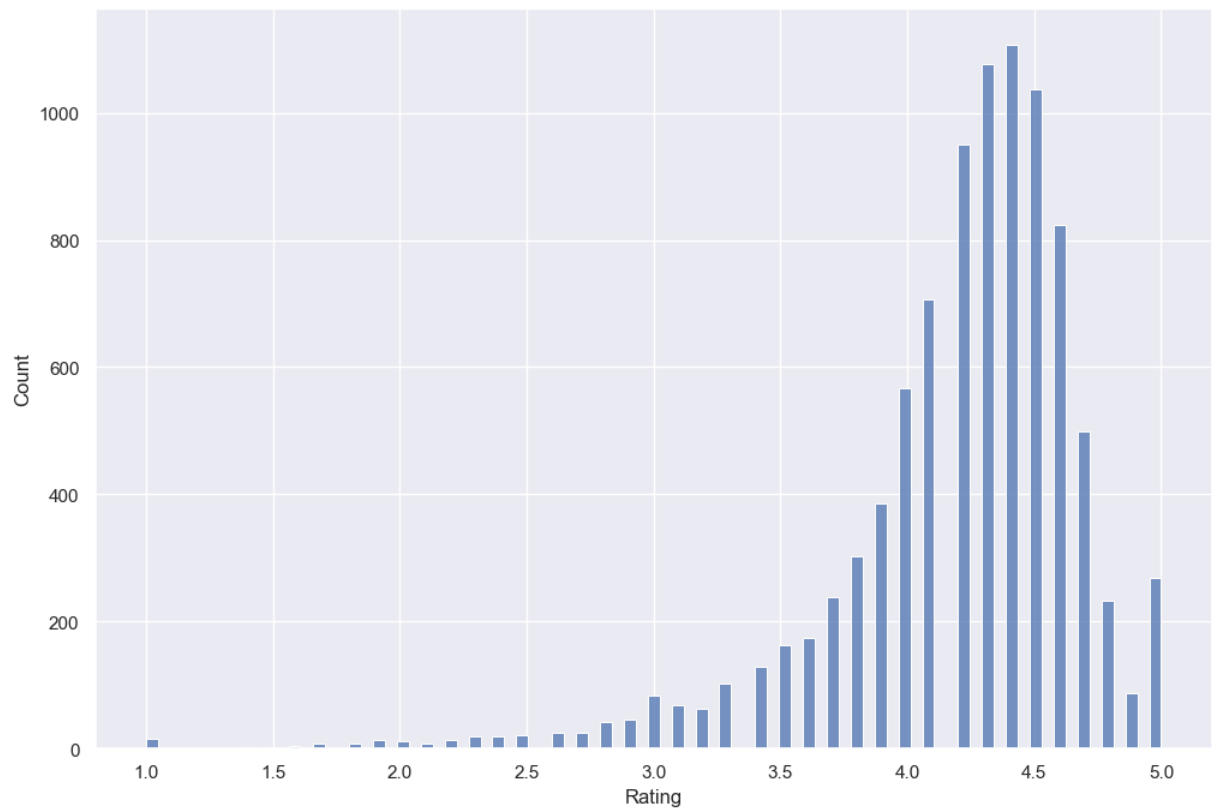
```
In [93]: bx_review = sns.boxplot(x='Reviews',data=df)
```



5(III).

```
In [94]: sns.histplot(df['Rating'])
```

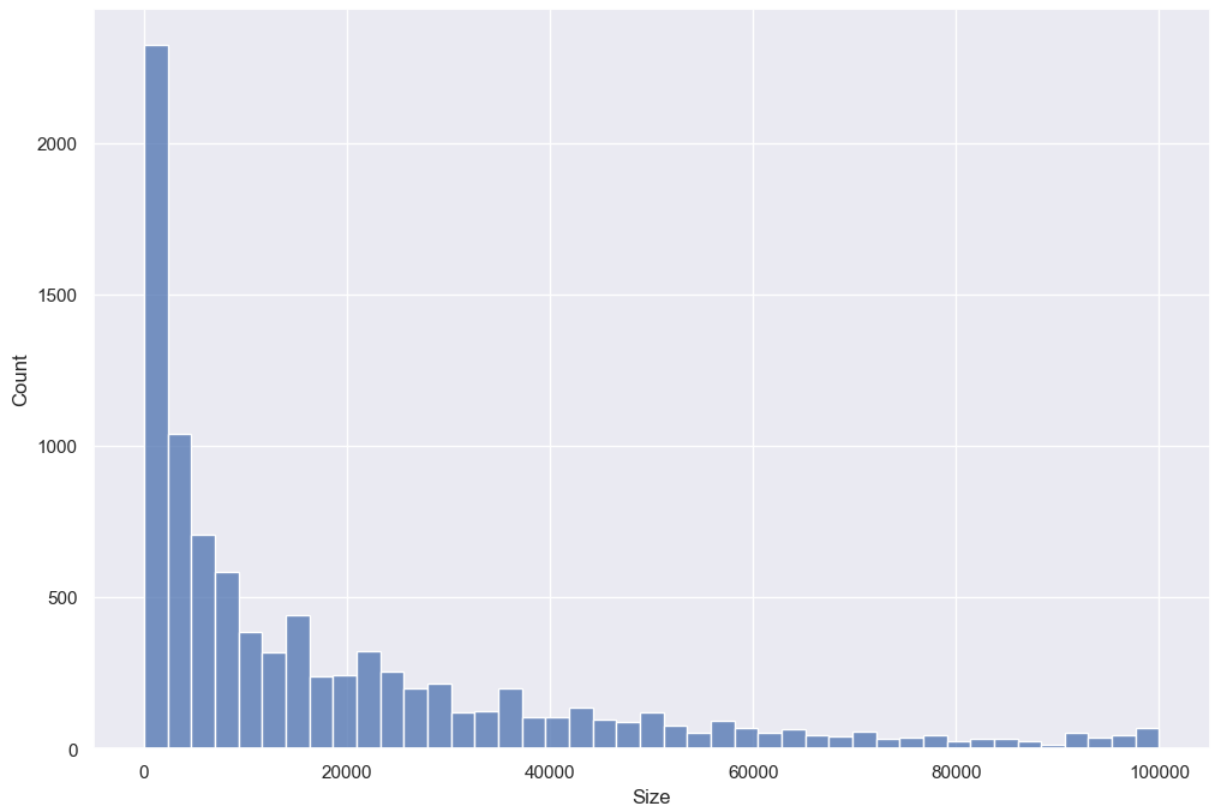
```
Out[94]: <Axes: xlabel='Rating', ylabel='Count'>
```



5(IV).

```
In [95]: sns.histplot(df['Size'])
```

```
Out[95]: <Axes: xlabel='Size', ylabel='Count'>
```



6(I).

```
In [104...] more = df.apply(lambda x : True
                             if x['Price'] > 200 else False, axis = 1)
```

```
In [105...] more_count = len(more[more == True].index)
```

```
In [106...] df.shape
```

```
Out[106]: (9338, 13)
```

```
In [107...] df.drop(df[df['Price'] > 200].index, inplace = True)
```

```
In [108...] df.shape
```

```
Out[108]: (9338, 13)
```

6(II).

```
In [111...] df.drop(df[df['Reviews'] > 2000000].index, inplace = True)
```

```
In [112...] df.shape
```

```
Out[112]: (8885, 13)
```

6(III).

```
In [113... df.quantile([.1, .25, .5, .70, .90, .95, .99], axis = 0)
```

C:\Users\mah3\AppData\Local\Temp\ipykernel_24232\2685684270.py:1: FutureWarning: The default value of numeric_only in DataFrame.quantile is deprecated. In a future version, it will default to False. Select only valid columns or specify the value of numeric_only to silence this warning.

```
df.quantile([.1, .25, .5, .70, .90, .95, .99], axis = 0)
```

```
Out[113]:
```

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

```
In [115... # dropping more than 10000000 Installs value  
df.drop(df[df['Installs'] > 10000000].index, inplace = True)
```

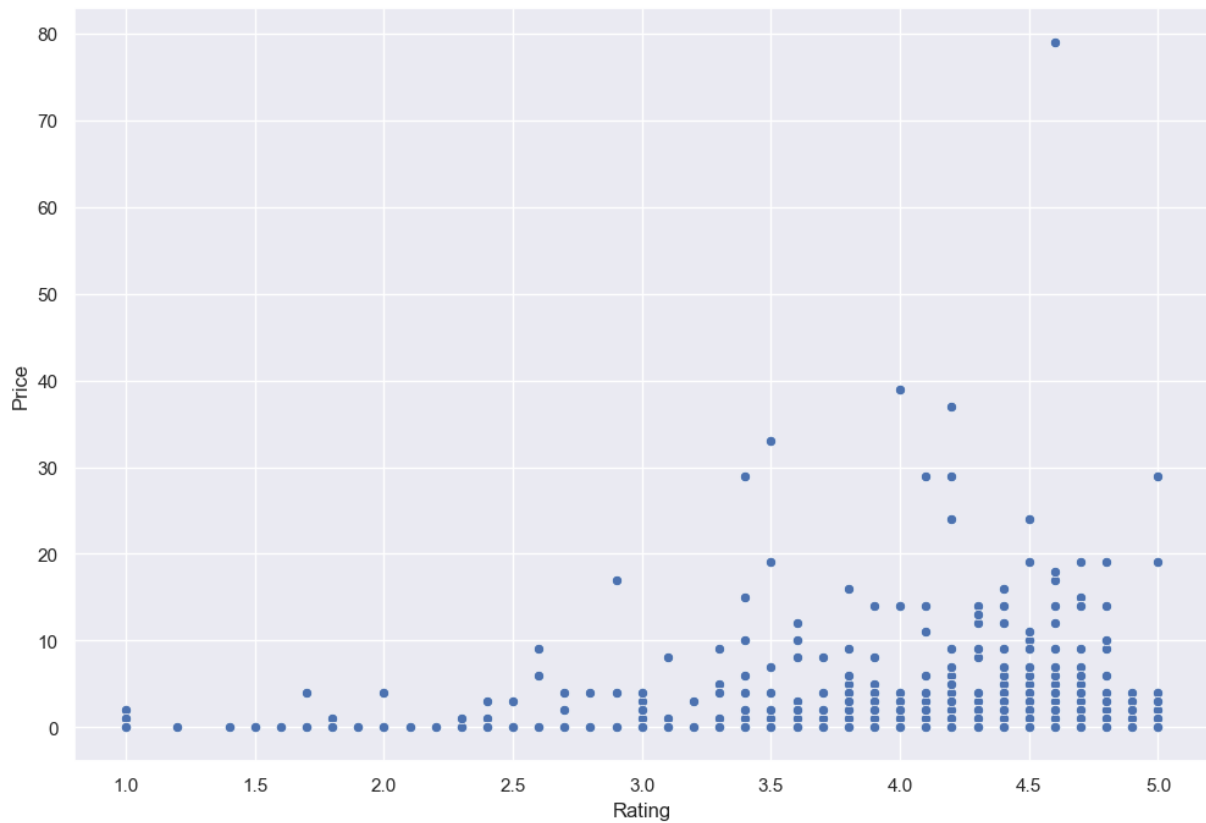
```
In [116... df.shape
```

```
Out[116]: (8496, 13)
```

7(I).

```
In [118... sns.scatterplot(x='Rating',y='Price',data=df)
```

```
Out[118]: <Axes: xlabel='Rating', ylabel='Price'>
```

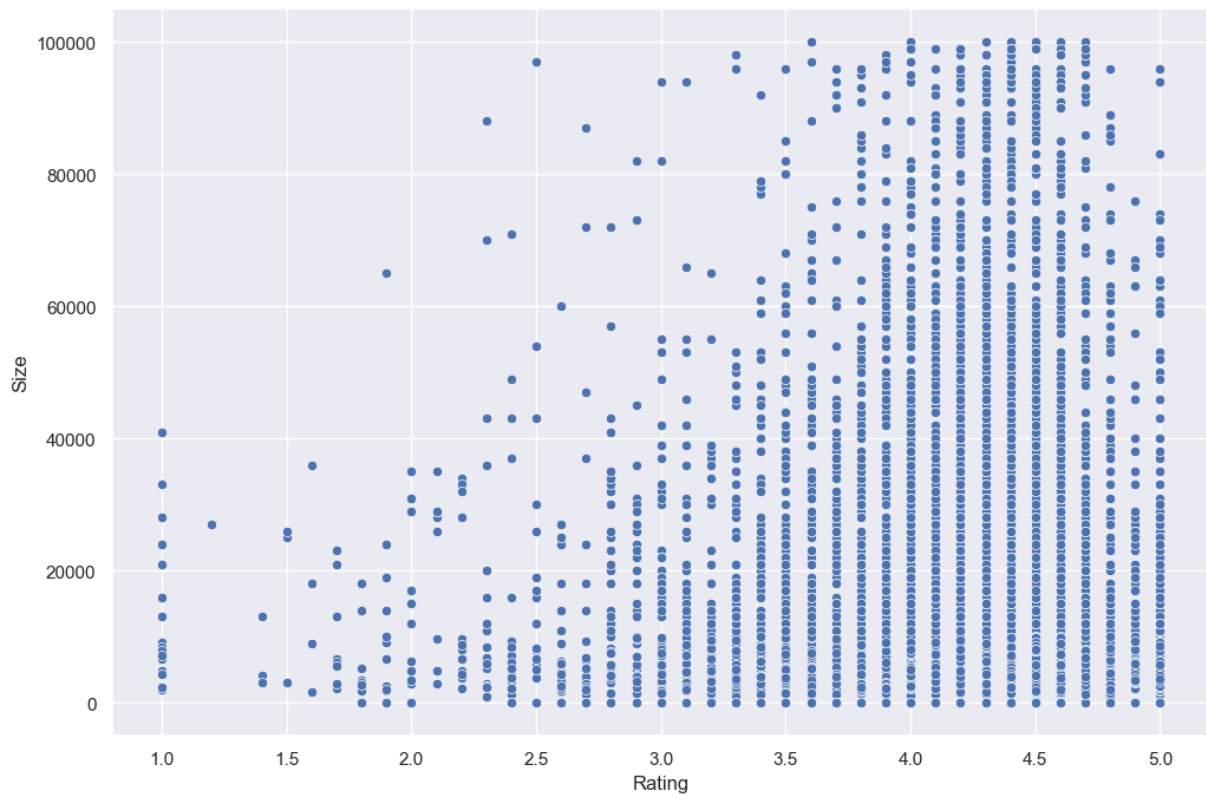


Yes, Paid apps are higher ratings compared to free apps.

7(II).

```
In [120]: sns.scatterplot(x='Rating',y='Size',data=df)
```

```
Out[120]: <Axes: xlabel='Rating', ylabel='Size'>
```

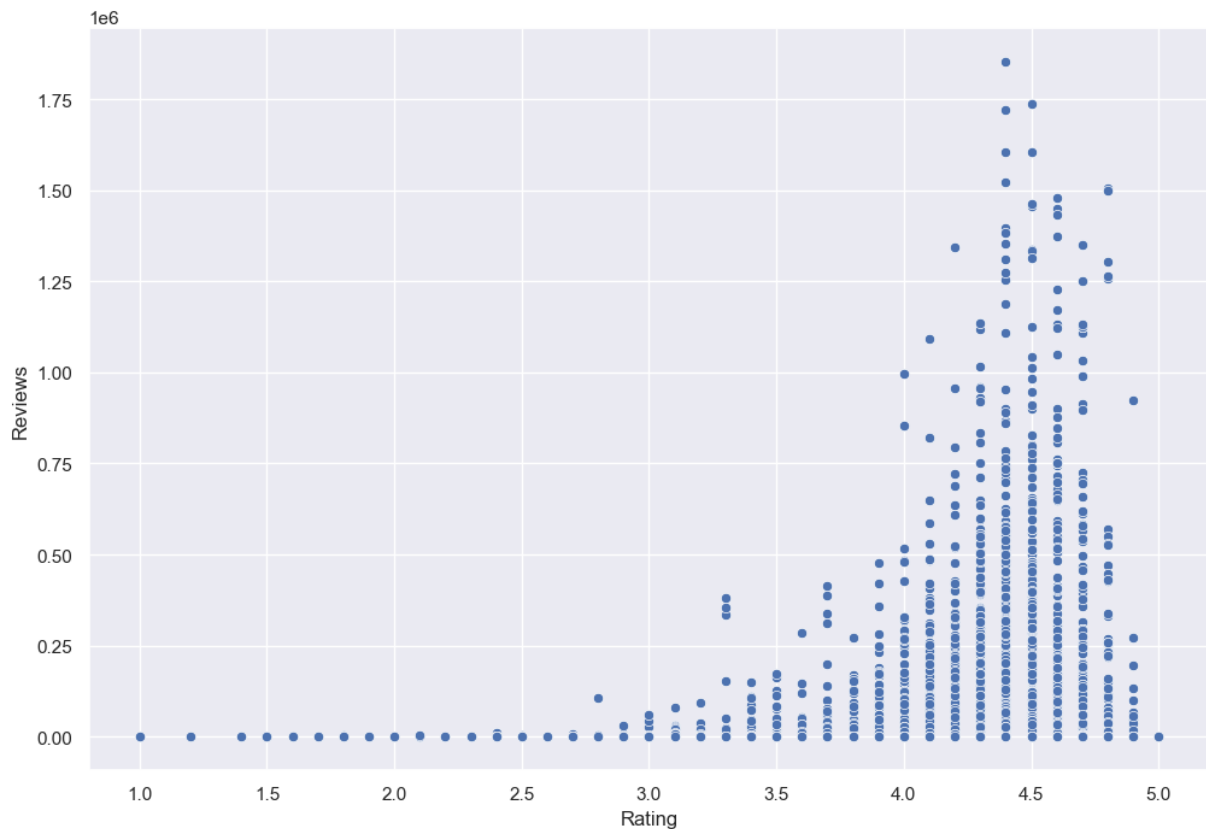



Yes it is clear that heavier apps are rated better.

7(III)

```
In [121]: sns.scatterplot(x='Rating',y='Reviews',data=df)
```

```
Out[121]: <Axes: xlabel='Rating', ylabel='Reviews'>
```

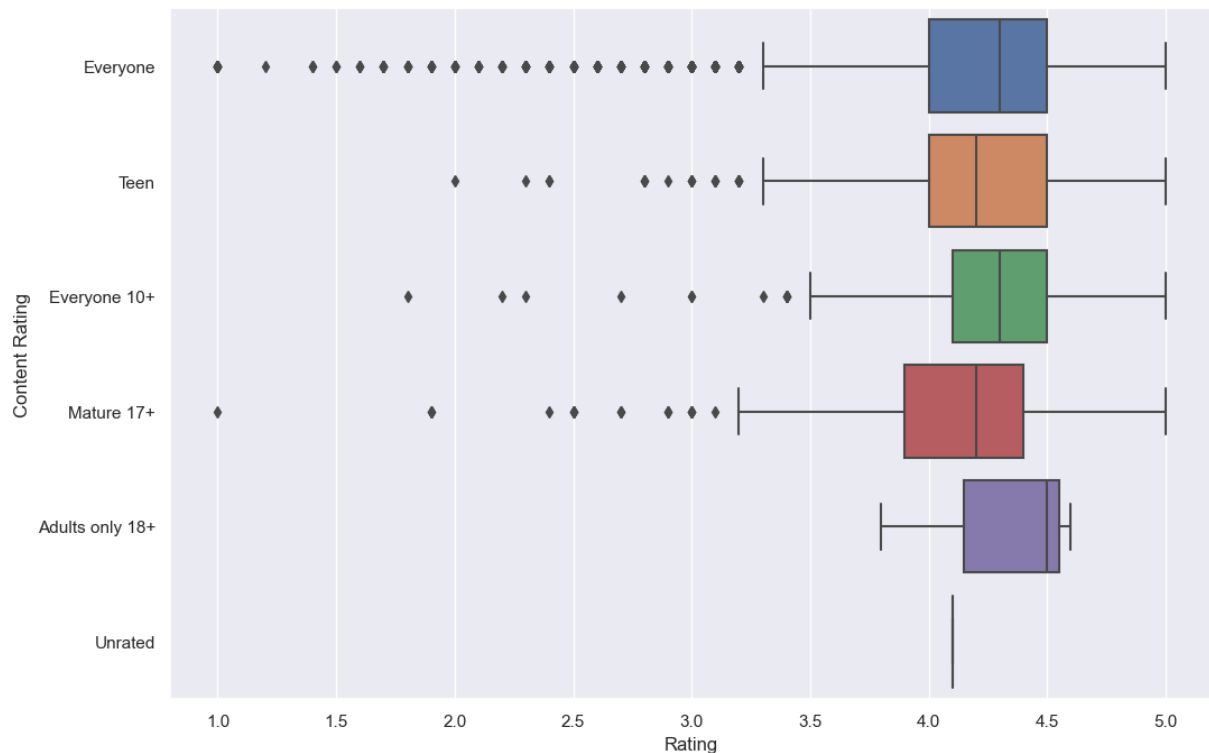


It is clear that more reviews makes app rating better.

7(IV).

```
In [122...] sns.boxplot(x="Rating", y="Content Rating", data=df)
```

```
Out[122]: <Axes: xlabel='Rating', ylabel='Content Rating'>
```

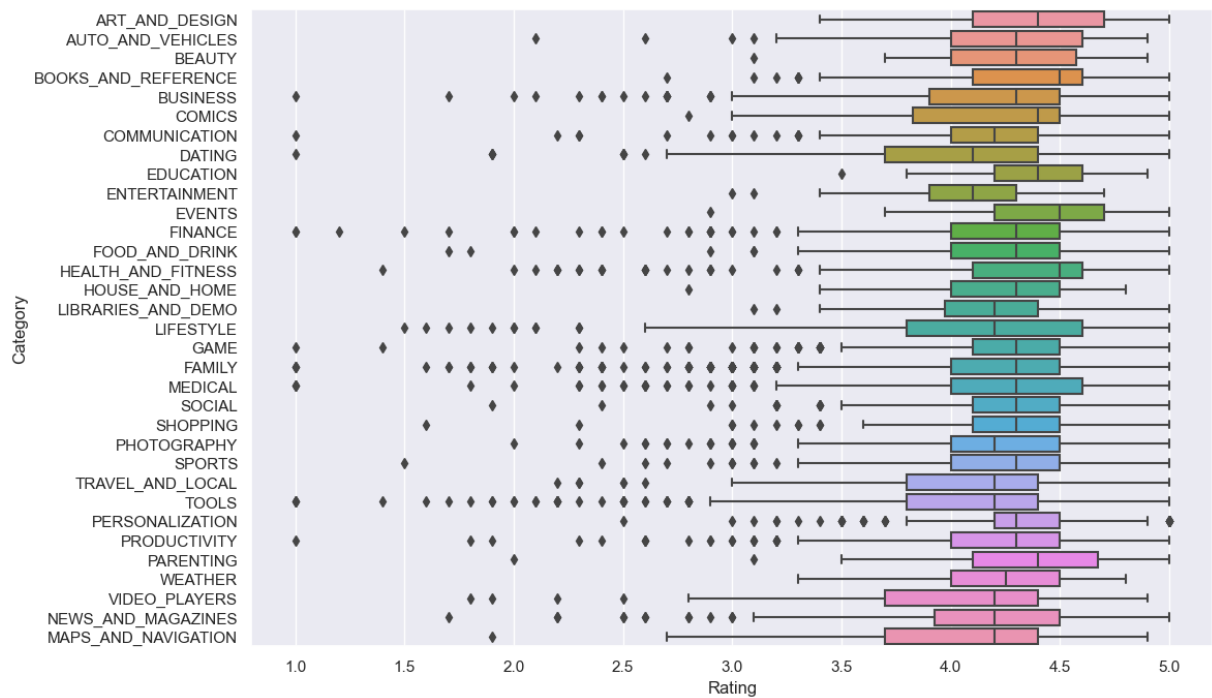


Apps which are for everyone has more bad ratings compare to other sections as it has so much outliers value, while 18+ apps have better ratings.

7(V).

```
In [123...] sns.boxplot(x="Rating", y="Category", data=df)
```

```
Out[123]: <Axes: xlabel='Rating', ylabel='Category'>
```



Events category has best ratings compare to others.

8(I).

In [124... `inp1 = df`

In [125... `inp1.head()`

Out[125]:

	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
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	ART_AND_DESIGN	4.3	967.0	2800.0	100000	Free	0	Everyone

	- Number								
	Art								
	Coloring								
	Book								
	Paper								
5	flowers	ART_AND_DESIGN	4.4	167.0	5600.0	50000	Free	0	Everyone
	instructions								

In [126... `inp1.skew()`

C:\Users\mahi3\AppData\Local\Temp\ipykernel_24232\3545313420.py:1: FutureWarning: The default value of numeric_only in DataFrame.skew is deprecated. In a future version, it will default to False. In addition, specifying 'numeric_only=None' is deprecated. Select only valid columns or specify the value of numeric_only to silence this warning.

`inp1.skew()`

Out[126]: Rating -1.749753
Reviews 4.576494
Size 1.655917
Installs 1.543697
Price 18.074542
dtype: float64

In [127... `reviewskew = np.log1p(inp1['Reviews'])`
`inp1['Reviews'] = reviewskew`

In [128... `reviewskew.skew()`

Out[128]: -0.20039949659264134

In [129... `installsskew = np.log1p(inp1['Installs'])`
`inp1['Installs']`

Out[129]: 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 [130... `installsskew.skew()`

Out[130]: -0.5097286542754812

In [131... `inp1.head()`

Out[131]:

	App	Category	Rating	Reviews	Size	Installs	Type	Price	Conte
--	-----	----------	--------	---------	------	----------	------	-------	-------

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

8(II).

In [132... `inp1.drop(["Last Updated", "Current Ver", "Android Ver", "App", "Type"], axis=1, inplace=`

In [133... `inp1.head()`

Out[133]:

	Category	Rating	Reviews	Size	Installs	Price	Content Rating	Genre
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;Pretenc 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 [134... `inp1.shape`

Out[134]: (8496, 8)

8(III).

```
In [135... inp2=inp1
```

```
In [136... inp2.head()
```

```
Out[136]:
```

	Category	Rating	Reviews	Size	Installs	Price	Content Rating	Genre
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 [137... inp2.Category.unique()
```

```
Out[137]: 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 [138... 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[138]:
```

	Rating	Reviews	Size	Installs	Price	Content Rating	Genres	Category_ART_AND_DESIGN
0	4.1	5.075174	19000.0	10000	0	Everyone	Art & Design	
1	3.9	6.875232	14000.0	500000	0	Everyone	Art &	

								Design;Pretend Play
2	4.7	11.379520	8700.0	5000000	0	Everyone	Art & Design	
4	4.3	6.875232	2800.0	100000	0	Everyone	Art & Design;Creativity	
5	4.4	5.123964	5600.0	50000	0	Everyone	Art & Design	

5 rows × 40 columns

In [139... `inp2.shape`

Out[139]: (8496, 40)

Let's apply Dummy EnCoding on Column "Genres"

In [140... `inp2["Genres"].unique()`

Out[140]: 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 [141... 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 [142... inp2["Genres"].unique()
```

```
Out[142]: 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 [143... 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 [144... inp2.head()
```

```
Out[144]:      Rating  Reviews  Size  Installs  Price  Content  Category_ART_AND_DESIGN  Cate
```

						Rating	
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 [145... inp2.shape
```

Out[145]: (8496, 91)

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

```
In [146... #get unique values in Column "Content Rating"
inp2["Content Rating"].unique()
```

Out[146]: array(['Everyone', 'Teen', 'Everyone 10+', 'Mature 17+',
 'Adults only 18+', 'Unrated'], dtype=object)

```
In [147... 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[147]:

	Rating	Reviews	Size	Installs	Price	Category_ART_AND_DESIGN	Category_AUTC
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 [148... inp2.shape
```

Out[148]: (8496, 96)

9. and 10.

```
In [149... 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 [150... 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 [151... reg_all = LR()
reg_all.fit(Xtrain,ytrain)
```

```
Out[151]: ▾ LinearRegression
LinearRegression()
```

```
In [154... 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 [153... 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