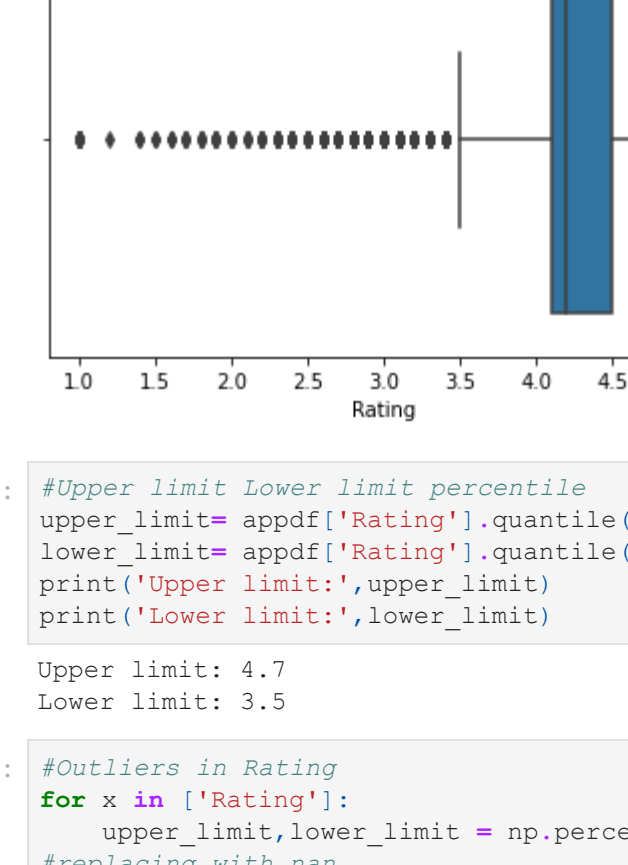


C:\Users\yujoydutta\Anaconda3\lib\site-packages\seaborn\decorators.py:36: FutureWarning: Pass the following va
riable as a keyword arg: x. From version 0.12, the only valid positional argument will be 'data', and passing o
ther arguments without an explicit keyword will result in an error or misinterpretation.

Out[32]:



In [33]:

```
#Upper limit Lower limit percentile  
upper_limit=appdf['Rating'].quantile(0.91)  
lower_limit= appdf['Rating'].quantile(0.09)  
print('Upper limit:',upper_limit)  
print('Lower limit:',lower_limit)
```

Upper limit: 4.7
Lower limit: 3.5

In [34]:

```
#Outliers in Rating  
for x in ['Rating']:  
    upper_limit,lower_limit = np.percentile(appdf.loc[:,x],[91,9])  
    #Replacing with nan  
    appdf.loc[appdf[x] < lower_limit,x] = np.nan  
    appdf.loc[appdf[x] > upper_limit,x] = np.nan
```

In [35]:

```
#Replacing outliers using mean  
appdf['Rating'].replace(np.nan, avgRating, inplace=True)
```

In [36]:

```
#checking value count  
appdf['Rating'].value_counts()
```

Out[36]:

Rating	count
4.13231	2664
4.40000	835
4.50000	822
4.30000	818
4.20000	740
4.60000	646
4.15000	594
4.00000	481
4.70000	421
3.90000	333
3.80000	272
3.70000	198
3.60000	158
3.50000	153

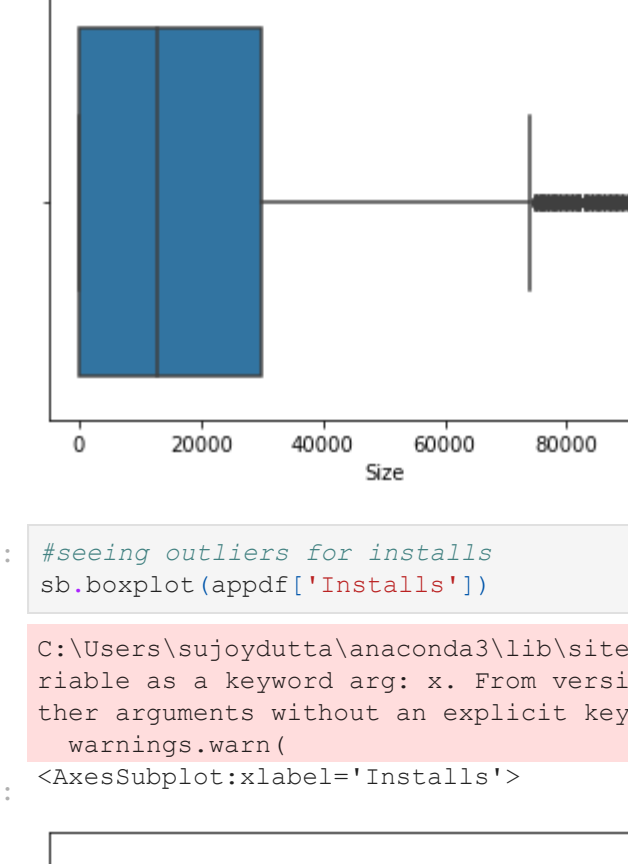
Name: Rating, dtype: int64

In [37]:

```
#seeing the outliers for Size  
sb.boxplot(appdf['Size'])
```

C:\Users\yujoydutta\Anaconda3\lib\site-packages\seaborn\decorators.py:36: FutureWarning: Pass the following va
riable as a keyword arg: x. From version 0.12, the only valid positional argument will be 'data', and passing o
ther arguments without an explicit keyword will result in an error or misinterpretation.

Out[37]:



In [38]:

```
#Upper limit Lower limit percentile  
upper_limit= appdf['Size'].quantile(0.91)  
lower_limit= appdf['Size'].quantile(0.09)  
print('Upper limit:',upper_limit)  
print('Lower limit:',lower_limit)
```

Upper limit: 58000.0
Lower limit: 2.8

In [39]:

```
#Replacing outliers using mean  
appdf['Size'].replace(np.nan, avgSize, inplace=True)
```

In [40]:

```
#checking value count  
appdf['Size'].value_counts()
```

Out[40]:

Size	count
1000.0	198
10000.0	196
14000.0	194
13000.0	191
15000.0	184

...

Size	count
425.0	1
200.0	1
480.0	1
720.0	1
619.0	1

Name: Size, Length: 458, dtype: int64

In [41]:

```
#seeing impact  
sb.boxplot(appdf['Size'])
```

C:\Users\yujoydutta\Anaconda3\lib\site-packages\seaborn\decorators.py:36: FutureWarning: Pass the following va
riable as a keyword arg: x. From version 0.12, the only valid positional argument will be 'data', and passing o
ther arguments without an explicit keyword will result in an error or misinterpretation.

Out[41]:

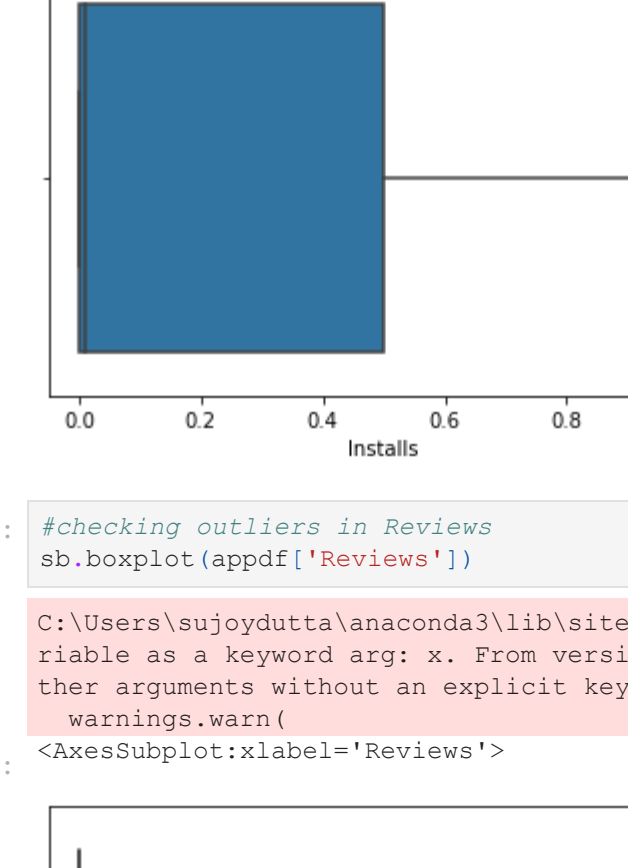


In [42]:

```
#seeing outliers for installs  
sb.boxplot(appdf['Installs'])
```

C:\Users\yujoydutta\Anaconda3\lib\site-packages\seaborn\decorators.py:36: FutureWarning: Pass the following va
riable as a keyword arg: x. From version 0.12, the only valid positional argument will be 'data', and passing o
ther arguments without an explicit keyword will result in an error or misinterpretation.

Out[42]:



In [43]:

```
#Upper limit Lower limit percentile  
upper_limit= appdf['Installs'].quantile(0.91)  
lower_limit= appdf['Installs'].quantile(0.09)  
print('Upper limit:',upper_limit)  
print('Lower limit:',lower_limit)
```

Upper limit: 10000000.0
Lower limit: 100.0

In [44]:

```
#Outliers in Install replacement with nan  
for x in ['Installs']:  
    upper_limit,lower_limit = np.percentile(appdf.loc[:,x],[91,9])  
    #Replacing with nan  
    appdf.loc[appdf[x] < lower_limit,x] = np.nan  
    appdf.loc[appdf[x] > upper_limit,x] = np.nan
```

In [45]:

```
#Replacing outliers using mean  
appdf['Installs'].replace(np.nan, avgInst, inplace=True)
```

In [46]:

```
#checking value count  
appdf['Installs'].value_counts()
```

Out[46]:

Installs	count
1.000000e+06	1303
7.122631e+06	1111
1.000000e+05	1054
1.000000e+04	1011
1.000000e+03	873
1.000000e+07	825
1.000000e+02	702
5.000000e+06	535
5.000000e+05	492
5.000000e+03	461
5.000000e+04	446
5.000000e+02	322

Name: Installs, dtype: int64

In [47]:

```
#seeing impact  
sb.boxplot(appdf['Installs'])
```

C:\Users\yujoydutta\Anaconda3\lib\site-packages\seaborn\decorators.py:36: FutureWarning: Pass the following va
riable as a keyword arg: x. From version 0.12, the only valid positional argument will be 'data', and passing o
ther arguments without an explicit keyword will result in an error or misinterpretation.

Out[47]:



In [48]:

```
#checking outliers in Reviews  
sb.boxplot(appdf['Reviews'])
```

C:\Users\yujoydutta\Anaconda3\lib\site-packages\seaborn\decorators.py:36: FutureWarning: Pass the following va
riable as a keyword arg: x. From version 0.12, the only valid positional argument will be 'data', and passing o
ther arguments without an explicit keyword will result in an error or misinterpretation.

Out[48]:



In [49]:

```
#Upper limit Lower limit percentile  
upper_limit= appdf['Reviews'].quantile(0.91)  
lower_limit= appdf['Reviews'].quantile(0.09)  
print('Upper limit:',upper_limit)  
print('Lower limit:',lower_limit)
```

Upper limit: 242684.44000000032
Lower limit: 1.0

In [50]:

```
#Outliers in Rating  
for x in ['Reviews']:  
    upper_limit,lower_limit = np.percentile(appdf.loc[:,x],[91,9])  
    #Replacing with nan  
    appdf.loc[appdf[x] < lower_limit,x] = np.nan  
    appdf.loc[appdf[x] > upper_limit,x] = np.nan
```

In [51]:

```
#Replacing outliers using mean  
appdf['Reviews'].replace(np.nan, avgRating, inplace=True)
```

In [52]:

```
#checking value count  
appdf['Reviews'].value_counts()
```

Out[52]:

Reviews	count
4.17321	1380
1.00000	128
2.00000	209
3.00000	173
4.00000	131

...

Reviews	count
4724.00000	1
43960.00000	1
229.00000	1
6626.00000	1
1195.00000	1

Name: Reviews, Length: 3909, dtype: int64

In [53]:

```
#seeing impact  
sb.boxplot(appdf['Reviews'])
```

C:\Users\yujoydutta\Anaconda3\lib\site-packages\seaborn\decorators.py:36: FutureWarning: Pass the following va
riable as a keyword arg: x. From version 0.12, the only valid positional argument will be 'data', and passing o
ther arguments without an explicit keyword will result in an error or misinterpretation.

Out[53]:

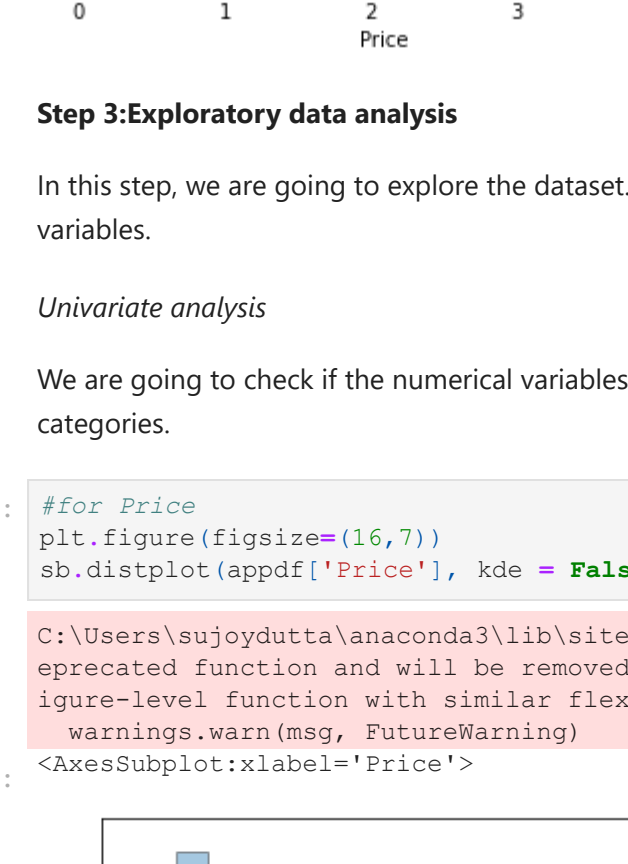


In [54]:

```
#checking outliers in Price  
sb.boxplot(appdf['Price'])
```

C:\Users\yujoydutta\Anaconda3\lib\site-packages\seaborn\decorators.py:36: FutureWarning: Pass the following va
riable as a keyword arg: x. From version 0.12, the only valid positional argument will be 'data', and passing o
ther arguments without an explicit keyword will result in an error or misinterpretation.

Out[54]:



In [55]:

```
#Upper limit Lower limit percentile  
upper_limit= appdf['Price'].quantile(0.91)  
lower_limit= appdf['Price'].quantile(0.09)  
print('Upper limit:',upper_limit)  
print('Lower limit:',lower_limit)
```

Upper limit: 0.0
Lower limit: 0.0

In [56]:

```
#Outliers in Price  
for x in ['Price']:  
    upper_limit,lower_limit = np.percentile(appdf.loc[:,x],[91,9])  
    #Replacing with nan  
    appdf.loc[appdf[x] < lower_limit,x] = np.nan  
    appdf.loc[appdf[x] > upper_limit,x] = np.nan
```

In [57]:

```
#Replacing outliers using mean  
appdf['Price'].replace(np.nan, avgRating, inplace=True)
```

In [58]:

```
#checking value count  
appdf['Price'].value_counts()
```

Out[58]:

Price	count
0.00000	8416
4.17321	719

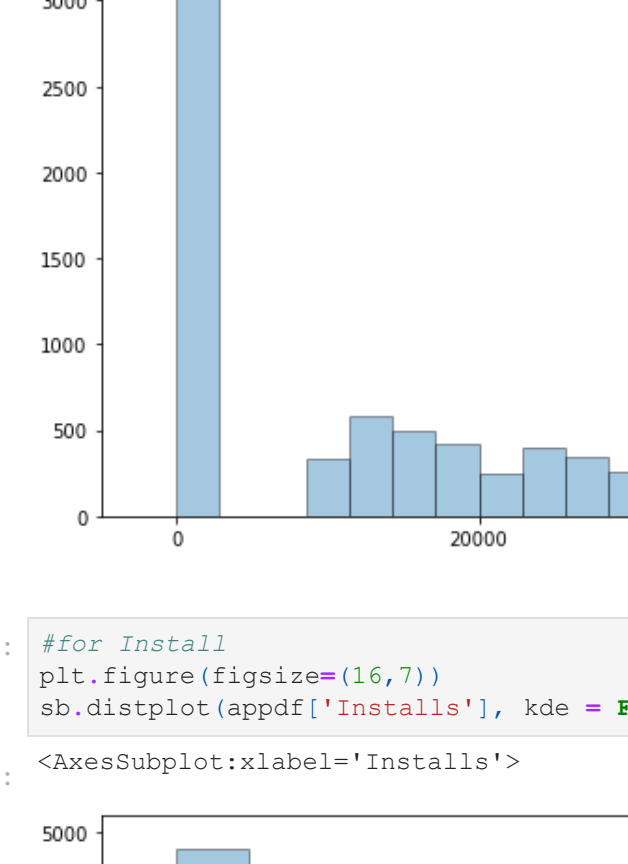
Name: Price, dtype: int64

In [59]:

```
#seeing impact  
sb.boxplot(appdf['Price'])
```

C:\Users\yujoydutta\Anaconda3\lib\site-packages\seaborn\decorators.py:36: FutureWarning: Pass the following va
riable as a keyword arg: x. From version 0.12, the only valid positional argument will be 'data', and passing o
ther arguments without an explicit keyword will result in an error or misinterpretation.

Out[59]:



Step 3:Exploratory data analysis

In this step, we are going to explore the dataset. Perform hypothesis tests,bivariate/univariate analysis and check for correlation between variables.

Univariate analysis

We are going to check if the numerical variables are skewed or not. We are also going to make bins and categorize data into quantitive categories.

In [60]:

```
#for Price  
plt.figure(figsize=(15,7))  
sb.boxplot(appdf['Price'], kde = False,hist_kws=dict(edgecolor="k"))
```

C:\Users\yujoydutta\Anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: 'distplot' is a d
eprated function and will be removed in a future version. Please adapt your code to use either 'displot' (a f
igure-level function with similar flexibility) or 'histplot' (an axes-level function for histograms).

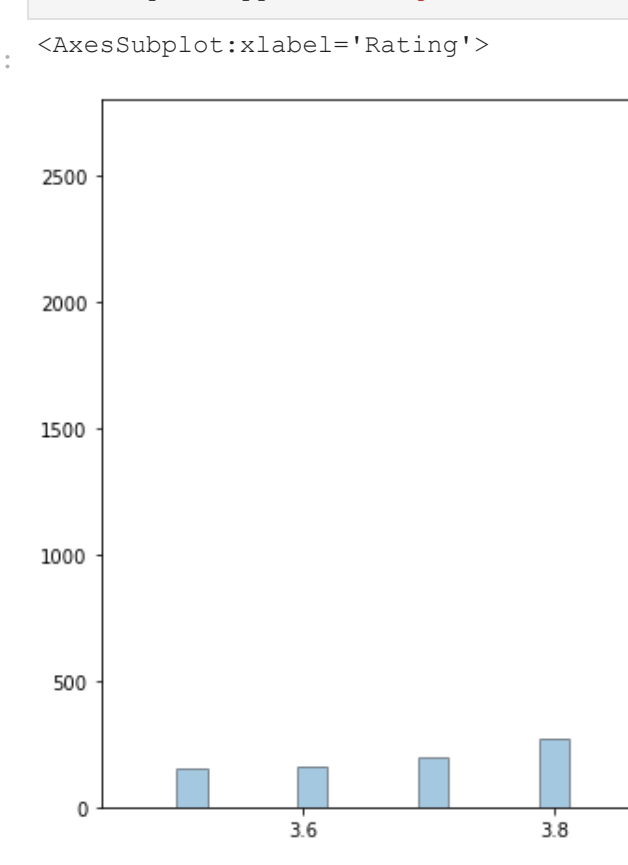
Out[60]:



In [61]:

```
#for Size  
plt.figure(figsize=(15,7))  
sb.distplot(appdf['Size'], kde = False,hist_kws=dict(edgecolor="k"))
```

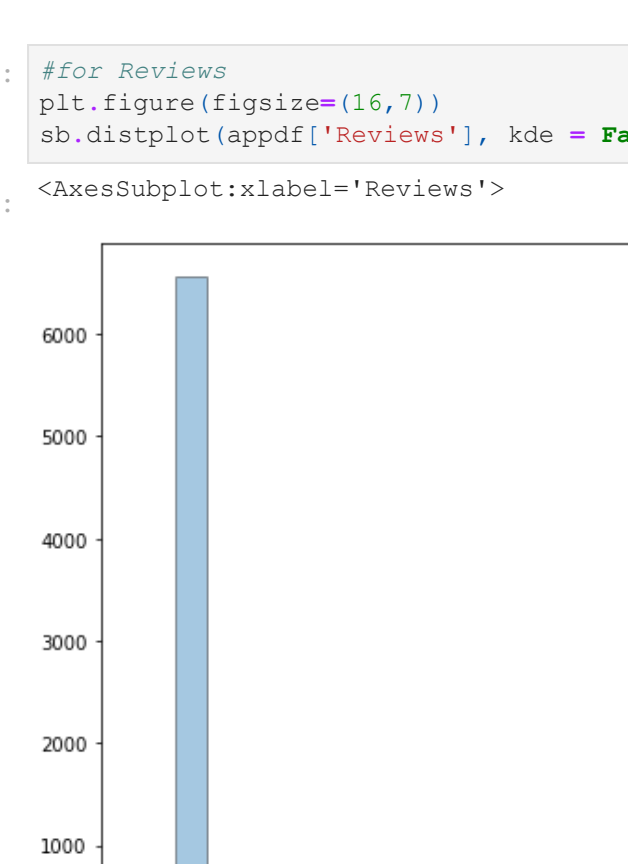
Out[61]:



In [62]:

```
#for Install  
plt.figure(figsize=(15,7))  
sb.distplot(appdf['Installs'], kde = False,hist_kws=dict(edgecolor="k"))
```

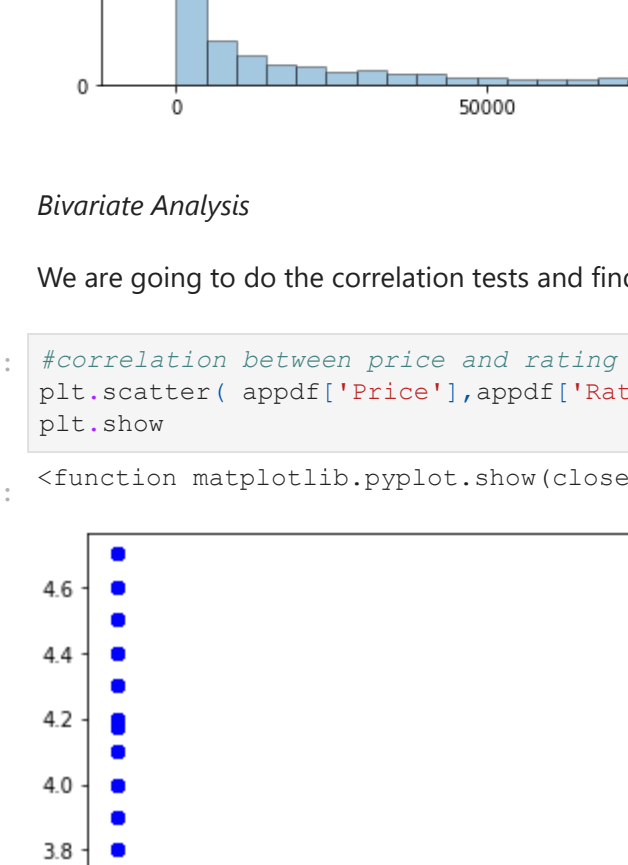
Out[62]:



In [63]:

```
#for Ratings  
plt.figure(figsize=(15,7))  
sb.distplot(appdf['Rating'], kde = False,hist_kws=dict(edgecolor="k"))
```

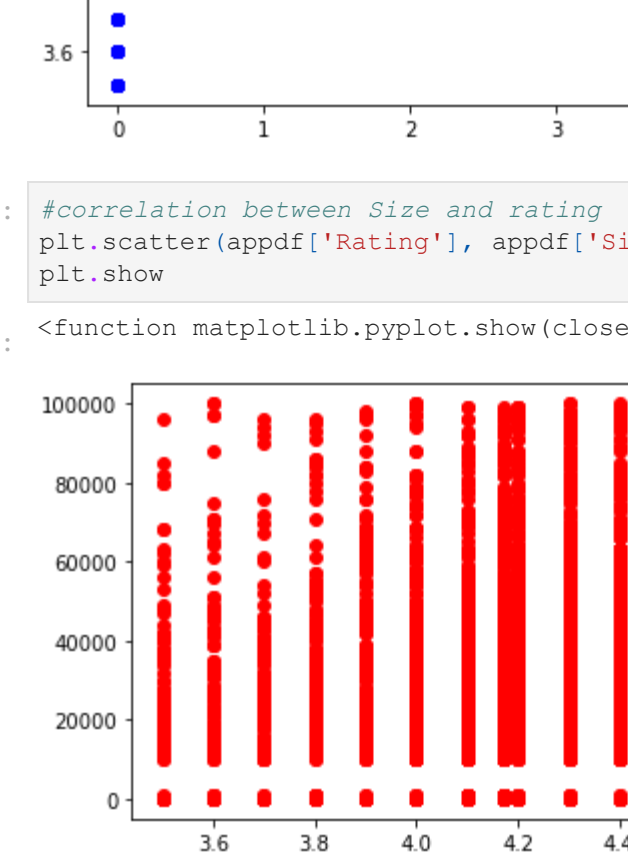
Out[63]:



In [64]:

```
#for Reviews  
plt.figure(figsize=(15,7))  
sb.distplot(appdf['Reviews'], kde = False,hist_kws=dict(edgecolor="k"))
```

Out[64]:



Bivariate Analysis

We are going to do the correlation tests and find out which variables are going to impact our dependent variable the most.

In [65]:

```
#correlation between price and rating  
plt.scatter(appdf['Price'],appdf['Rating'], c ="blue")  
plt.show
```

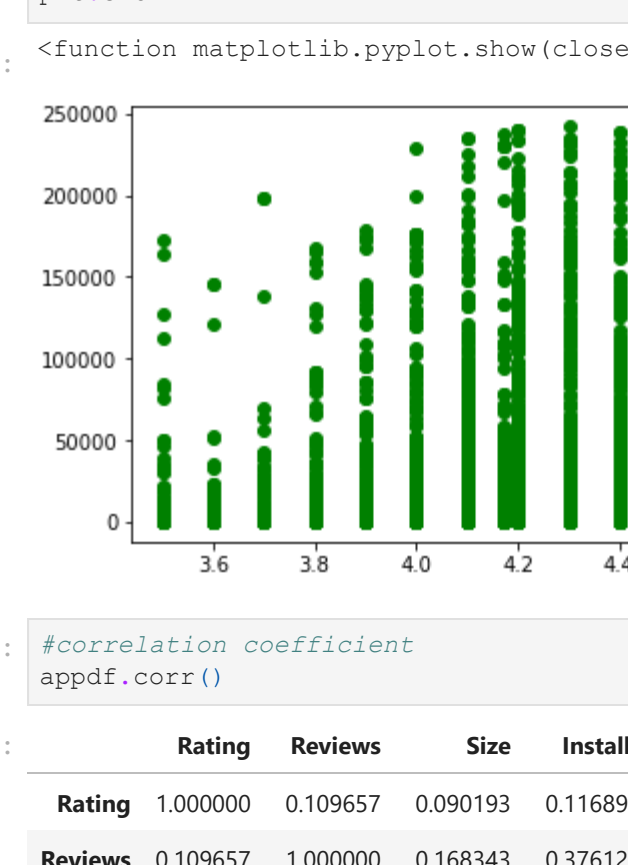
Out[65]:



In [66]:

```
#correlation between size and rating  
plt.scatter(appdf['Rating'], appdf['Size'], c ="red")  
plt.show
```

Out[66]:



In [67]:

```
#correlation between size and rating  
plt.scatter(appdf['Rating'], appdf['Installs'], c ="pink")  
plt.show
```

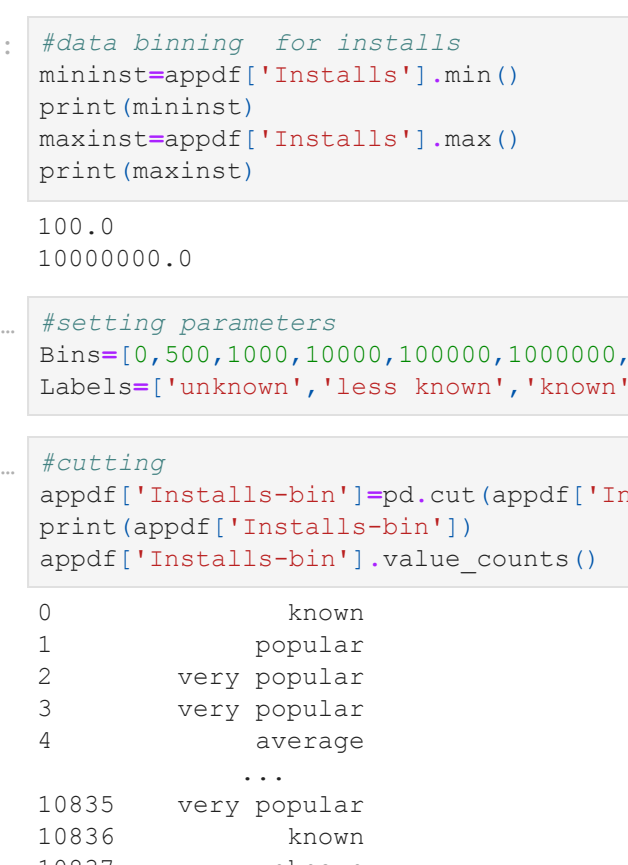
Out[67]:



In [68]:

```
#correlation between reviews and rating  
plt.scatter(appdf['Rating'], appdf['Reviews'], c ="Green")  
plt.show
```

Out[68]:



In [69]:

```
#correlation coefficient  
appdf.corr()
```

Out[69]:

	Rating	Reviews	Size	Installs	Price
Rating	1.000000	0.109657	0.090193	0.116897	0.039352
Reviews	0.109657	1.000000	0.168343	0.376120	-0.084486
Size	0.090193	0.168343	1.000000	0.229774	-0.024926
Installs	0.116897	0.376120	0.229774	1.000000	-0.069327
Price	0.039352	-0.084486	-0.024926	-0.069327	1.000000

In [70]:

```
#summary statistics  
appdf.describe().round(3)
```

Out[70]:

	Rating	Reviews	Size	Installs	Price
count	9135.0000	9135.0000	9135.0000	9135.0000	9135.0000
mean	4.225	15871.610	19579.917	2.2487232e+06	0.328
std	0.256	39003.789	24045.451	3.431325e+06	1.124
min	3.500	1.000	1.000	1.000000e+03	0.000
25%	4.173	6.000	5.600	5.000000e+03	0.000
50%	4.173	204000	13000.000	1.000000e+05	0.000
75%	4.400	7716.500	30000.000	5.000000e+06	0.000
max	4.700	242096.000	100000.000	1.000000e+07	4.173

In [71]:

```
#data binning for installs  
mininst=appdf['Installs'].min()  
print(mininst)  
maxinst=appdf['Installs'].max()  
print(maxinst)
```

mininst: 100.0
maxinst: 1000000.0

In [71]:

```
#setting parameters  
bins=(0,500,1000,10000,100000,1000000,10000000)  
labels=['unknown','less known','known','average','popular','very popular']
```

In [72]:

```
#cutting  
appdf['Installs-bin']=mpd.cut(appdf['Installs'],bins=bins,labels=labels)  
print(appdf['Installs-bin'])  
appdf['Installs-bin'].value_counts()
```

0 unknown
1 popular
2 very popular
3 very popular
4 average

10835 very popular
10836 known
10837 unknown
10838 less known
10840 very popular

appdf['Installs-bin'].length: 9135, dtype: category
Categories (6, object): ['unknown' < 'less known' < 'known' < 'average' < 'popular' < 'very popular']

very popular 2471
popular 1500
known 1472
unknown 1024
less known 873
Name: Installs-bin, dtype: int64

In [72]:

```
#plotting the histogram  
plt.bar(labels, appdf['Installs-bin'].value_counts())  
plt.xlabel("Number of apps in category")  
plt.ylabel("Number of installs")  
plt.title("Install price bins")
```

Out[72]:



In [73]:

```
#binning for sizes  
mininst=appdf['Size'].min()  
print(mininst)  
maxinst=appdf['Size'].max()  
print(maxinst)
```

mininst: 1.0
maxinst: 100000.0

In [73]:

```
#setting parameters  
bins=(0,5000,10000,50000)  
labels=['light','average','heavy']
```

In [74]:

```
#cutting  
appdf['Size-bin']=mpd.cut(appdf['Size'],bins=bins,labels=labels)  
appdf['Size-bin'].value_counts()
```

light 3920
average 3006
heavy 1794
Name: Size-bin, dtype: int64

In [75]:

```
#plotting the histogram  
plt.bar(labels, appdf['Size-bin'].value_counts())  
plt.xlabel("Number of apps in category")  
plt.ylabel("Size in kilobytes")  
plt.title("Size bins")
```

Out[75]:



In [76]:

```
#checking distribution of Rating  
stat, p = shapiro(appdf['Rating'])  
print('stat=%.3f, p=%.3f' % (stat, p))  
if p > 0.05:  
    print('Probably Gaussian')  
else:  
    print('Probably not Gaussian')
```

stat=0.942, p=0.000
Probably not Gaussian

C:\Users\yujoydutta\Anaconda3\lib\site-packages\scipy\stats\morestats.py:1760: FutureWarning: p-value may not be accurate for N > 5000.
warnings.warn("p-value may not be accurate for N > 5000.")

In [77]:

```
#checking distribution of Size  
stat, p = shapiro(appdf['Size'])  
print('stat=%.3f, p=%.3f' % (stat, p))  
if p > 0.05:  
    print('Probably Gaussian')  
else:  
    print('Probably not Gaussian')
```

stat=0.802, p=0.000
Probably not Gaussian

In [78]:

```
#checking distribution of Reviews  
stat, p = shapiro(appdf['Reviews'])  
print('stat=%.3f, p=%.3f' % (stat, p))  
if p > 0.05:  
    print('Probably Gaussian')  
else:  
    print('Probably not Gaussian')
```

stat=0.469, p=0.000
Probably not Gaussian

In [79]:

```
#checking distribution of Installs  
stat, p = shapiro(appdf['Installs'])  
print('stat=%.3f, p=%.3f' % (stat, p))  
if p > 0.05:  
    print('Probably Gaussian')  
else:  
    print('Probably not Gaussian')
```

stat=0.665, p=0.000
Probably not Gaussian

In [80]:

```
#checking distribution of Price  
stat, p = shapiro(appdf['Price'])  
print('stat=%.3f, p=%.3f' % (stat, p))  
if p > 0.05:  
    print('Probably Gaussian')  
else:  
    print('Probably not Gaussian')
```

stat=0.465, p=0.000
Probably not Gaussian


```
stat=0.297, p=0.000
Probably not causal.
```

```
In [18]: #checking if price is dependent on app popularity
inscat=appdf[['install', 'price']]
mod = ols("price ~ inscat", data = appdf).fit()
anov_table = sm.stats.anova_lm(mod)
anov_table
```

	df	sum_sq	mean_sq	F	PR(>F)
inscat	5.0	3213.10924	642.62185	52.309415	1.101626e-53
Residual	9129.0	11214.588560	1228501	NaN	NaN

```
In [85]: #checking if price is dependent on app size
sizcat=appdf[['size', 'price']]
mod = ols("price ~ sizcat", data = appdf).fit()
anov_table = sm.stats.anova_lm(mod)
anov_table
```

	df	sum_sq	mean_sq	F	PR(>F)
sizcat	7.0	21357180	10678590	8.383646	0.00003
Residual	8717.0	11103.196534	1273471	NaN	NaN

```
In [86]: #checking if Rating is dependent on app popularity
inscat=appdf[['install', 'rating']]
mod = ols("rating ~ inscat", data = appdf).fit()
anov_table = sm.stats.anova_lm(mod)
anov_table
```

	df	sum_sq	mean_sq	F	PR(>F)
inscat	5.0	9.036342	1.807268	27.897998	3.805860e-28
Residual	9129.0	591.388409	0.064781	NaN	NaN

```
In [87]: appdf.head(10)
```

	App	Category	Rating	Reviews	Size	Installs	Type	Price	Content Rating	Genres	Installs-bin	Size-bin
--	-----	----------	--------	---------	------	----------	------	-------	----------------	--------	--------------	----------

```
In [88]: #checking if Rating differs to apps being free or paid
Type=appdf['Type'].unique()
inscat=appdf[['Type', 'rating']]
mod = ols("Rating ~ Type", data = appdf).fit()
anov_table = sm.stats.anova_lm(mod)
anov_table
```

	df	sum_sq	mean_sq	F	PR(>F)
Type	1.0	0.929794	0.929794	14.164939	0.000169
Residual	9133.0	599.494957	0.065641	NaN	NaN

```
In [89]: #checking if Rating differs according to content type
Cont=appdf[['Content Rating', 'rating']]
mod = ols("Rating ~ Cont", data = appdf).fit()
anov_table = sm.stats.anova_lm(mod)
anov_table
```

	df	sum_sq	mean_sq	F	PR(>F)
Cont	5.0	2.219889	0.443978	6.775393	0.000003
Residual	9129.0	596.204862	0.065528	NaN	NaN

```
In [90]: #checking if Rating is different in different Category
Cat=appdf[['Category', 'rating']]
mod = ols("Rating ~ Cat", data = appdf).fit()
anov_table = sm.stats.anova_lm(mod)
anov_table
```

	df	sum_sq	mean_sq	F	PR(>F)
Cat	32.0	14.122668	0.441333	6.851445	2.276268e-29
Residual	9102.0	586.302083	0.064415	NaN	NaN

```
#seeing average by content rating
appdf.groupby('Content Rating').mean().round(2)
```

	Rating	Reviews	Size	Installs	Price
Content Rating					
Adults only 18+	4.55	37011.00	20502.45	750000.00	0.00
Everyone	4.22	14309.78	17215.27	2102136.54	0.35
Everyone 10+	4.29	29292.61	38949.73	3465137.88	0.38
Mature 17+	4.18	20186.65	23833.02	2164815.76	0.17
Teen	4.23	21291.71	29108.00	2968605.79	0.22
Unrated	4.14	5935.30	4.45	25250.00	0.00

```
In [122]: #which content rating is more rated
x=appdf['Content Rating']
y=appdf['Rating'].mean()
#visualization stage
plt.bar(x, y)
```



```
In [93]: #which type is more popular
partofday=appdf['Type'].unique()
incidents=appdf['Type'].value_counts()
```



```
In [94]: #seeing average by Category
appdf.groupby('Category').mean().round(2)
```

	Rating	Reviews	Size	Installs	Price
Category					
ART AND DESIGN	4.32	17537.86	9470.65	1025538.09	0.20
AUTO AND VEHICLES	4.17	9605.64	17923.26	1438317.04	0.11
BEAUTY	4.26	3975.75	10939.03	442968.74	0.00
BOOKS AND REFERENCE	4.26	9320.97	10354.90	1795130.09	0.57
BUSINESS	4.20	6065.91	11962.18	2139010.27	0.14
COMICS	4.25	12082.24	11999.41	48484.92	0.00
COMMUNICATION	4.17	20788.97	9725.35	287721.29	0.36
DATING	4.14	13204.53	14181.78	1307596.17	0.10
EDUCATION	4.36	43212.06	19299.53	2511729.73	0.15
ENTERTAINMENT	4.16	35470.03	20823.25	3977122.75	0.05
EVENTS	4.25	2222.35	11587.19	93680.45	0.07
FAMILY	4.22	14329.14	26666.52	1906236.09	0.40
FINANCE	4.23	9351.85	15957.58	1078802.26	0.23
FOOD AND DRINK	4.19	18182.41	21021.16	2417921.78	0.04
GAME	4.28	26108.27	43945.59	4033074.90	0.32
HEALTH AND FITNESS	4.27	19349.90	20449.56	2231866.51	0.19
HOUSE AND HOME	4.21	11893.46	12895.09	1418827.83	0.00
LIBRARIES AND DEMO	4.20	8382.88	9622.60	841689.40	0.05
LIFESTYLE	4.18	8412.70	12323.31	15580137.52	0.21
MAPS AND NAVIGATION	4.18	15337.37	14100.56	1553385.05	0.15
MEDICAL	4.22	3288.00	17908.88	1226053.58	0.94
NEWS AND MAGAZINES	4.22	14139.46	10346.70	1925119.03	0.04
PARENTING	4.26	4196.61	20113.51	575807.98	0.08
PERSONALIZATION	4.27	12374.41	8702.93	205039.33	0.92
PHOTOGRAPHY	4.24	29029.71	14911.86	3656296.51	0.31
PRODUCTIVITY	4.22	15239.68	10569.64	2446136.08	0.33
SHOPPING	4.25	27315.52	14271.36	3466175.71	0.04
SOCIAL	4.23	15515.90	15192.95	271567.41	0.06
SPORTS	4.23	22954.83	22730.01	2809365.12	0.30
TOOLS	4.18	11814.32	6174.50	187353.70	0.38
TRAVEL AND LOCAL	4.19	19211.85	23263.59	2228605.78	0.22
VIDEO PLAYERS	4.18	19839.93	33779.25	249877.06	0.03
WEATHER	4.22	22181.95	10591.69	1941541.98	0.44

```
In [95]: #which content rating is more popular
x=appdf['Category']
y=appdf['Rating'].mean()
#visualization stage
plt.bar(x, y)
```



```
In [96]: #seeing average by Category
appdf.groupby('Genres').mean().round(2)
```

	Rating	Reviews	Size	Installs	Price
Genres					
Action	4.28	22778.07	50164.34	4201242.30	0.31
Action&Action & Adventure	4.30	34379.34	53006.00	4876281.25	0.78
Adventure	4.19	27886.89	36160.36	2000915.84	0.73
Adventure&Action & Adventure	4.42	7866.89	49693.37	5404129.10	0.64
Adventure&Brain Games	4.60	7148.00	8.10	100000.00	0.17
...
Video Players & Editors	4.18	18764.68	13708.33	2479532.85	0.03
Video Players & Editors&Creativity	4.10	159620.50	23006.00	5000000.00	0.00
Video Players & Editors&Music & Video	4.00	39779.00	19000.00	5000000.00	0.00
Weather	4.22	22181.95	10591.69	1941541.98	0.44
Word	4.29	78167.07	34481.24	4736710.47	0.00

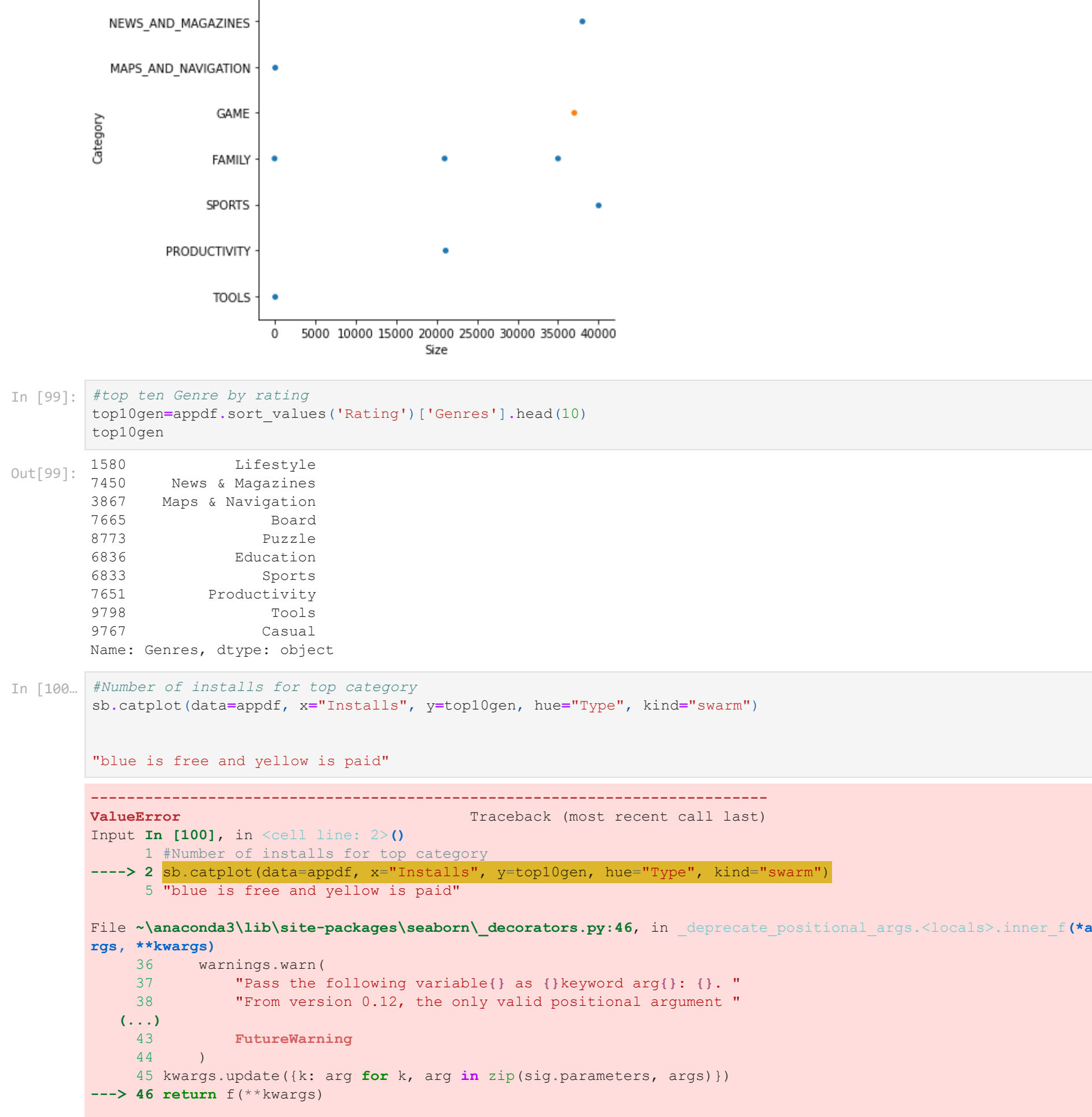
116 rows x 5 columns

```
In [97]: #top ten category by rating
top10gen=appdf.sort_values('Rating')['Category'].head(10)
top10cat
```

1580	LIFESTYLE
7450	NEWS AND MAGAZINES
3867	MAPS & NAVIGATION
7665	GAME
8773	FAMILY
6836	FAMILY
6833	SPORTS
7651	PRODUCTIVITY
9798	TOOLS
9767	FAMILY
Name: Category, dtype: object	

```
In [98]: #which category apps are the heaviest
sb.catplot(data=appdf, x="Size", y=top10cat, hue="Type", kind="swarm")
plt.cparams('figure.figsize') = [50, 10]
```

"blue is free and yellow is paid"



```
In [99]: #top ten Genre by rating
top10gen=appdf.sort_values('Rating')['Genres'].head(10)
top10gen
```

1580	Lifestyle
7450	News & Magazines
3867	Maps & Navigation
7665	Board
8773	Puzzle
6836	Education
6833	Sports
7651	Productivity
9798	Tools
9767	Casual
Name: Genres, dtype: object	

```
In [100]: #Number of installs for top category
sb.catplot(data=appdf, x="Installs", y=top10gen, hue="Type", kind="swarm")
```

"blue is free and yellow is paid"



```
In [101]: #Which Genre is more popular
genotypes=appdf['Genres'].unique()
number=appdf['Genres'].value_counts()
```

#color for each label, 'l', 'q', 'b'



