

```
In [1]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
%matplotlib inline
import seaborn as sns
import warnings
warnings.filterwarnings("ignore")
```

1. Load the data file using pandas.

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

```
Out[2]:
```

	App	Category	Rating	Reviews	Size	Installs	Type	Price	Content Rating	Genres	Last Updated	Current Ver	Android Ver
0	Photo Editor & Candy Camera & Glit & StickerBook	ART_AND_DESIGN	4.1	159	19M	10000+	Free	0	Everyone	Art & Design	January 7, 2018	1.0.0	4.0.3 and up
1	Cooking book:moana	ART_AND_DESIGN	3.9	967	14M	500,000+	Free	0	Everyone	Art & Design/Premid Play	January 15, 2018	2.0.0	4.0.3 and up
2	U Launcher Lite – FREE Live Cool Themes, Hide...	ART_AND_DESIGN	4.7	87510	8.7M	5,000,000+	Free	0	Everyone	Art & Design	August 1, 2018	1.2.4	4.0.3 and up
3	Sketch - Draw & Paint	ART_AND_DESIGN	4.5	215644	25M	100,000+	Free	0	Teen	Art & Design	June 8, 2018	Varies with device	4.2 and up
4	Pixol Draw - Number Art Coloring Book	ART_AND_DESIGN	4.3	967	28M	100,000+	Free	0	Everyone	Art & Design/Creativity	June 20, 2018	1.1	4.4 and up

2. Check for null values in the data. Get the number of null values for each column.

```
In [3]: inp.info()
# As we can find there are some columns that has null values..!!

<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                   10841 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: object(12), float64(1)
memory usage: 1.1+ MB

In [4]: inp.isnull().sum()
# Count of null values in each columns..!!

Out[4]:
App                0
Category           0
Rating             14
Reviews            0
Size               0
Installs           0
Type               0
Price              0
Content Rating     0
Genres             0
Last Updated       0
Current Ver        8
Android Ver        3
dtype: int64

3. Drop records with nulls in any of the columns.

In [5]: inp.dropna(how='any', inplace=True)
inp.shape

Out[5]: (9366, 13)
```

4. Variables seem to have incorrect type and inconsistent formatting. You need to fix them :

```
In [6]: # Size column has sizes in Kb as well as Mb, Multiply the value by 1,000, if size is mentioned in Mb.

def new_size(Size):
    if 'K' in Size:
        x = Size[:-1]
        x = float(x)*1000
        return(x)
    elif 'M' in Size[:-1]:
        x = Size[:-1]
        x = float(x)
        return(x)
    else:
        return None

inp['Size'] = inp['Size'].map(new_size)
```

```
In [18]: inp['Size'].isnull().sum()

Out[18]: 1637

inp.Size.fillna(method = "ffill", inplace = True)

In [12]: inp['Size'].isnull().sum()

Out[12]: 0

In [13]: # 2. Reviews is a numeric field that is loaded as a string field. Convert it to numeric (int/float).
inp.Reviews = inp.Reviews.astype("float")

In [14]: # 3. Installs field is currently stored as string and has values like 1,000,000+.
# Treat 1,000,000+ as 1,000,000
# remove '+', ',' from the field, convert it to integer

def clean(val):
    return float(val.replace(",","").replace("+",""))

inp.Installs = inp.Installs.map(clean)

In [16]: # 4. Price field is a string and has $ symbol. Remove '$' sign, and convert it to numeric.
inp['Price'] = inp.Price.map(lambda x: 0 if x=="$" else float(x[1:]))

In [17]: # Checking the data types after conversion..!!
inp.dtypes
```

```
Out[17]:
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

In [ ]: # 5. Sanity checks:
# Reviews should not be more than installs as only those who installed can review the app.
# Reviews should not be more than Installs as only those who installed can review the app. If there are any such records, drop them.
# For free apps (type = "Free"), the price should not be >0. Drop any such rows.

In [19]: inp.Rating.describe()
# The values are within range..!! Min val is 1.0 and Max val is 5.0

Out[19]:
count    9366.000000
mean     4.130238
min      1.000000
max      5.000000
25%      4.000000
50%      4.380000
75%      4.580000
max      5.000000
Name: Rating, dtype: float64

In [23]: # Reviews should not be more than installs as only those who installed can review the app.
# If there are any such records, drop them.

len(inp[inp.Reviews > inp.Installs])

Out[23]: 7

In [24]: inp[inp[inp.Reviews < inp.Installs].copy()]

In [25]: inp.shape

Out[25]: (9355, 13)

In [26]: # 4. For free apps (type = "Free"), the price should not be >0. Drop any such rows.
len(inp[inp.Type == "Free"] & (inp.Price > 0)])

Out[26]: 0
```

5. Performing univariate analysis :

```
In [29]: # Boxplot for Price
# Are there any outliers? Think about the price of usual apps on Play Store.

sns.set_style(style='whitegrid')
sns.boxplot(inp['Price'])
plt.title('Prices of Apps', fontsize=18)
plt.xlabel('Price', fontsize=14)
# I can see that there are many outliers in Price column

Out[29]: Text(8.5, 6, 'Price')

Prices of Apps

```

```
In [30]: # Boxplot for Reviews
# Are there any apps with very high number of reviews? Do the values seem right?

sns.boxplot(inp['Reviews'])
plt.title('Reviews of Apps', fontsize=18)
plt.xlabel('Reviews', fontsize=14)
# I can see that there are many outliers in Price column

Out[30]: Text(8.5, 6, 'Reviews')

Reviews of Apps

```

```
In [33]: # Histogram for Rating
# How are the ratings distributed? Is it more toward higher ratings?

sns.histplot(inp['Rating'], color='green')
plt.title('Ratings of Apps', fontsize=18)
plt.xlabel('Rating', fontsize=14)
plt.ylabel('Frequency', fontsize=14)
# Rating is distributed as left skewed which we can say as Negatively skewed and it is towards higher ratings.
# It is looking a very good distribution as left skewed..!!

Out[33]: Text(0, 0.5, 'Frequency')

Ratings of Apps

```

```
In [34]: # Histogram for Size
sns.histplot(inp['Size'], palette='blue')
plt.title('Size of Apps', fontsize=18)
plt.xlabel('Size', fontsize=14)
plt.ylabel('Frequency', fontsize=14)
# I can see that the distribution is a right skewed and it starts from higher Size and continue to lower Size.

Out[34]: Text(0, 0.5, 'Frequency')

Sizes of Apps

```

6. Outlier treatment:

Price: From the box plot, it seems like there are some apps with very high price. A price of \$200 for an application on the Play Store is very high and suspicious!

```
In [35]: # Check out the records with very high price
# Is 200 indeed a high price?
# Drop those as most apps are not supposed to be junk apps

len(inp[inp.Price > 200])

Out[35]: 15

In [39]: # There are 15 apps with price greater than 200...!!

inp[inp.Price > 200]

Out[39]:
```

	App	Category	Rating	Reviews	Size	Installs	Type	Price	Content Rating	Genres	Last Updated	Current Ver	Android Ver
4197	most expensive app (N)	FAMILY	4.3	6.0	1500.0	100.0	Paid	399.99	Everyone	Entertainment	July 16, 2018	1.0	7.0 and up
4362	fm rich - Trump Edition	LIFESTYLE	3.8	718.0	26000.0	10000.0	Paid	399.99	Everyone	Lifestyle	March 11, 2018	1.0.0	4.4 and up
4367	fm rich - Trump Edition	LIFESTYLE	3.6	275.0	7300.0	10000.0	Paid	400.00	Everyone	Lifestyle	May 3, 2018	1.0.1	4.1 and up
5361	I am rich	FAMILY	3.8	3547.0	1800.0	10000.0	Paid	399.99	Everyone	Lifestyle	January 12, 2018	2.0	4.0.3 and up
5364	I am rich Plus	FAMILY	4.0	856.0	8700.0	10000.0	Paid	399.99	Everyone	Entertainment	May 19, 2018	3.0	4.4 and up
5365	I am rich VIP	LIFESTYLE	3.8	411.0	2600.0	10000.0	Paid	299.99	Everyone	Lifestyle	July 21, 2018	1.1.1	4.3 and up
5366	I AM Rich Premium	FINANCE	4.1	1887.0	4700.0	5000.0	Paid	399.99	Everyone	Finance	November 12, 2017	1.6	4.0 and up
5368	I am extremely rich	LIFESTYLE	2.9	41.0	2900.0	1000.0	Paid	279.99	Everyone	Lifestyle	July 1, 2018	1.0	4.0 and up
5369	I am rich	FINANCE	3.5	92.0	20000.0	1000.0	Paid	399.99	Everyone	Finance	December 11, 2017	1.0	4.1 and up
5369	I am rich(gemium)	FINANCE	3.5	472.0	960.0	5000.0	Paid	399.99	Everyone	Finance	May 1, 2017	3.4	4.4 and up
5362	I Am Rich Pro	FAMILY	4.4	201.0	2700.0	1000.0	Paid	399.99	Everyone	Entertainment	May 30, 2017	1.54	1.6 and up
5364	I am rich (Most expensive app)	FINANCE	4.1	129.0	2700.0	1000.0	Paid	399.99	Teen	Finance	December 6, 2017	2.7	4.0.3 and up
5366	I am Rich	FAMILY	3.6	217.0	4900.0	10000.0	Paid	399.99	Everyone	Entertainment	June 22, 2018	1.5	4.2 and up
5369	I AM Rich	FINANCE	4.3	180.0	3800.0	5000.0	Paid	399.99	Everyone	Finance	March 22, 2018	1.0	4.2 and up
5373	I AM RICH PRO PLUS	FINANCE	4.0	36.0	41000.0	1000.0	Paid	399.99	Everyone	Finance	June 25, 2018	1.0.2	4.1 and up

```
In [48]: # Dropping the Junk apps...!!

inp[inp.Price>298].copy()

In [41]: inp.shape

Out[41]: (9336, 13)

In [43]: # Reviews: Very few apps have very high number of reviews. These are all star apps that don't help with the analysis and,
# In fact, will skew it..!! Drop records having more than 2 million reviews.

len(inp[inp.Reviews>2000000])

Out[43]: 453

In [44]: inp[inp[inp.Reviews>2000000].copy()]

In [45]: inp.shape

Out[45]: (8883, 13)

In [46]: # Installs: There seems to be some outliers in this field too.
# Apps having very high number of installs should be dropped from the analysis.
# Find out the different percentiles - 10, 25, 50, 70, 90, 95, 99

inp.Installs.quantile([0.1, 0.25, 0.5, 0.75, 0.9, 0.95, 0.99])

Out[46]:
0.10    10000.0
0.25    10000.0
0.50    500000.0
0.75    1000000.0
0.90    1000000.0
0.95    1000000.0
0.99    10000000.0
Name: Installs, dtype: float64

In [47]: len(inp[inp.Installs >= 10000000])

Out[47]: 1627

In [48]: inp[inp[inp.Installs>=10000000].copy()]

In [49]: inp.shape

Out[49]: (8494, 13)
```

7. Bivariate analysis: Let's look at how the available predictors relate to the variable of interest, i.e., our target variable rating. Make scatter plots (for numeric features) and box plots (for character features) to assess the relations between rating and the other features.

```
In [51]: # Make scatter plot/joinplot for Rating vs. Price
# What pattern do you observe? Does rating increase with price?

plt.figure(figsize=(11,8))
sns.set(style='whitegrid')
sns.set(font_scale=1.2)

sns.scatterplot(x=inp['Price'], y=inp['Rating'], hue= inp['Rating'], palette='autumn_r')
plt.title('Rating vs price', fontsize=18)
plt.xlabel('Rating', fontsize=14)
plt.ylabel('Price', fontsize=14)
plt.show()
```

Rating Vs price

From the above scatterplot I can see that there is not a very clean pattern, it does look that the higher priced apps have better rating. Although, there are not a lot of apps which gets high priced, but the pattern is significant.

```
In [54]: # Are heavier apps rated better?

plt.figure(figsize=(11,8))
sns.boxplot(inp.Size, inp.Rating)
sns.setSubplot(inp.Size, inp.Rating, hue=inp.Rating)
<AxesSubplot: xlabel='Size', ylabel='Rating'>
```

Sizes vs Ratings

No clear pattern. There are few low rated apps among the popular ones (maybe poor ones won't get so popularity), after a certain point, the rating does not depend on the popularity.

```
In [56]: # Make boxplot for Rating vs. Content Rating
# Is there any difference in the ratings? Are some types liked better?

plt.figure(figsize=(12,8))
sns.boxplot(x='Content Rating', y = 'Rating', data = inp, palette='hls')
sns.setSubplot(inp.Size, inp.Rating, hue=inp.Rating)
plt.xticks(fontsize=12, rotation=45)
plt.xlabel('Rating', fontsize=14)
plt.ylabel('Content Rating', fontsize=14)
plt.show()
```

Content Rating vs Ratings

While the median rating for most others is similar, the rating for "Adults Only 18+" is the highest.

```
In [57]: # Make boxplot for Ratings vs. Category
# Which genre has the best ratings?

plt.figure(figsize=(18,6))
sns.boxplot(x = 'Category', y = 'Rating', data = inp)
plt.title('Category vs Ratings', fontsize=20)
plt.xticks(fontsize=12, rotation='vertical')
plt.xlabel('Category', fontsize=16)
plt.ylabel('Rating', fontsize=16)

Out[57]: Text(0, 0.5, 'Rating')

Category vs Ratings


Apps around Health & Fitness, Books and Reference, Events seem to have the highest median ratings.



8. Data preprocessing



```
In [ ]: # For the steps below, create a copy of the dataframe to make all the edits. Name it inp1.

In [58]: # Reviews and Install have some values that are still relatively very high. Before building a linear regression model, you need to reduce the skew.
# Apply log transformation (np.log) to Reviews and Installs.

inp1 = inp.copy(deep=True)

In [59]: inp1.head()
```



```
Out[59]:
```



|   | App                                              | Category       | Rating | Reviews | Size    | Installs  | Type | Price | Content Rating | Genres                   | Last Updated     | Current Ver | Android Ver  |
|---|--------------------------------------------------|----------------|--------|---------|---------|-----------|------|-------|----------------|--------------------------|------------------|-------------|--------------|
| 0 | Photo Editor & Candy Camera & Glit & StickerBook | ART_AND_DESIGN | 4.1    | 159     | 19000.0 | 10000.0   | Free | 0.0   | Everyone       | Art & Design             | January 7, 2018  | 1.0.0       | 4.0.3 and up |
| 1 | Cooking book:moana                               | ART_AND_DESIGN | 3.9    | 967     | 14000.0 | 500000.0  | Free | 0.0   | Everyone       | Art & Design/Premid Play | January 15, 2018 | 2.0.0       | 4.0.3 and up |
| 2 | U Launcher Lite – FREE Live Cool Themes, Hide... | ART_AND_DESIGN | 4.7    | 87510.0 | 8700.0  | 5000000.0 | Free | 0.0   | Everyone       | Art & Design             | August 1, 2018   | 1.2.4       | 4.0.3 and up |
| 3 | Sketch - Draw & Paint                            | ART_AND_DESIGN | 4.5    | 215644  | 25000.0 | 100000.0  | Free | 0.0   | Everyone       | Art & Design/Creativity  | June 20, 2018    | 1.1         | 4.4 and up   |
| 4 | Pixol Draw - Number Art Coloring Book            | ART_AND_DESIGN | 4.4    | 167.0   | 5600.0  | 50000.0   | Free | 0.0   | Everyone       | Art & Design             | March 26, 2017   | 1.0         | 2.3 and up   |


```

```
In [61]: inp1.Reviews[inp1.Reviews>1000000]
inp1.Installs[inp1.Installs>10000000]

In [62]: inp1.hist(column=['Reviews', 'Installs'])

Out[62]: array([(('Category', 'Rating', 'Reviews', 'Size', 'Installs', 'Type', 'Price', 'Content Rating', 'Genres'),
<AxesSubplot: title='(center: 'Reviews')>, dtype=object)

Actual Predicted
0 1.7 4.370365
1 3.0 3.761739
2 3.2 3.969185
3 4.0 4.136414
4 4.2 4.103574
5 5.0 3.804942
6 3.9 4.015669
7 4.7 4.489564
8 4.5 4.279194
9 4.2 4.226446

In [79]: fig = plt.figure(figsize=(25))
fig.plot(kind='bar', figsize=(11,8))

Out[79]: <AxesSubplot>
```

Actual Predicted

From the above figure we can observe here that the model has returned pretty good prediction results.

```
In [83]: # Report the R2 on the train set
print('R2_Score =', r2_score(y_test, predict))

R2_Score = 0.34132279149655282

In [81]: print('Prediction Error Percentage is', round((0.50/np.mean(y_test))**100))

Prediction Error Percentage is 12

In [90]: print('Root Mean Squared Error=', np.sqrt(ms(y_test, predict)))

Root Mean Squared Error = 0.587480754931616

In [91]: y_train_preds = linreg.predict(X_train)
r2_score(y_train, y_train_preds)

Out[91]: 0.1639498257788446

In [ ]:
```