

## \*\*Regression\*\*

### Agenda

In this lesson, we will cover the following concepts with the help of a business use case:

- Use Case: Regression
- Regression Algorithms
- Types of Model Evaluation Metrics
- Gradient Descent
- Types of Gradient Descents
- Use Case: Stochastic Gradient Descent (SGD)

### Use Case: Regression

Note: At first, with the help of a use case, we are going to perform all the basic steps to reach the model training and prediction part.

#### \*\*Problem Statement:\*\*

Google Play Store team is about to launch a new feature wherein certain apps that are promising are boosted in visibility. The boost will manifest in multiple ways including higher priority in recommendations sections ("Similar apps", "You might also like", "New and updated games"). These will also get a boost in search results visibility. This feature will help bring more attention to newer apps that have the potential.

#### \*\*Analysis to identify the:\*\*

The problem is to identify the apps that are going to be good for Google to promote. App ratings, which are provided by the customers, are always great indicators of the goodness of the app. The problem reduces to: predict which apps will have high ratings.

#### \*\*Dataset\*\*

Google Play Store data ([googleplaystore.csv](#))

Link: <https://www.dropbox.com/sh/06ohrau3ucfqbml/AAceYXuml56543knDNQFj8ma?dl=0>

#### \*\*Data Dictionary:\*\*

	Variables	Description
	App	Application name
	Category	Category to which the app belongs
	Rating	Overall user rating of the app
	Reviews	Number of user reviews for the app
	Size	Size of the app
	Installs	Number of user downloads/installs for the app
	Type	Paid or Free
	Price	Price of the app
	Content_Rating	Age group the app is targeted at - Children / Mature 21+ / Adult
	Genres	An app can belong to multiple genres (apart from its main category) For example a musical family game will belong to Music, Game, Family genres
	Last_Updated	Date when the app was last updated on Play Store
	Current_Ver	Current version of the app available on Play Store
	Android_Ver	Minimum required Android version

#### \*\*Solution:\*\*

##### \*\*Import Libraries\*\*

```
#Importing libraries
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt, seaborn as sns
import matplotlib
```

##### \*\*Import and Check Dataset\*\*

```
inp0 = pd.read_csv("googleplaystore.csv")
```

```
# Check first five rows
inp0.head(2)
```

	App	Category	Rating	Reviews	Size	Installs	Type	Price	Content_Rating	Genres	Last_Updated	Current_Ver	Android_Ver
0	Photo Editor & Candy Camera & Grid & Scrapbook	ART_AND_DESIGN	4.1	159	19M	10,000+	Free	0	Everyone	Art&Design	January 7, 2018	1.0.0	4.1
1	Coloring book for kids: Dora the Explorer	ART_AND_DESIGN	3.9	967	14M	500,000+	Free	0	Everyone	Art&Design,Prete	January 15, 2018	2.0.0	4.1

##### Observations

The data will be displayed on the screen.

```
#Check number of columns and rows, and data types
inp0.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 10841 entries, 0 to 10840
Data columns (total 13 columns):
#   Non-Null Count  Dtype
---  --
0   App             object
1   Category        object
2   Rating          float64
3   Reviews         object
4   Size            object
5   Installs        object
6   Type            object
7   Price          float64
8   Content_Rating  object
9   Genres          object
10  Last_Updated    object
11  Current_Ver     object
12  Android_Ver     object
dtypes: float64(1), object(12)
memory usage: 1.1+ MB
```

##### \*\*Check Data Types\*\*

```
#checking datatypes
inp0.dtypes
```

```
App             object
Category        object
Rating          float64
Reviews         object
Size            object
Installs        object
Type            object
Price          float64
Content_Rating  object
Genres          object
Last_Updated    object
Current_Ver     object
Android_Ver     object
dtype: object
```

##### \*\*Finding and Treating Null Values\*\*

```
#Finding count of null values
inp0.isnull().sum(axis=0)
```

```
App             0
Category        0
Rating          1474
Reviews         0
Size            0
Installs        0
Type            0
Price           1
Content_Rating  0
Genres          0
Last_Updated    8
Current_Ver     3
Android_Ver     0
dtype: int64
```

```
#Dropping the records with null ratings
#This is done because ratings is our target variable
inp0.dropna(how="any", inplace=True)
```

```
inp0.isnull().sum(axis=0)
```

```
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
```

##### Handling the Variables

###### 1. Clean the price column

```
#Cleaning the price column
inp0.Price.value_counts()[1:5]
```

```
0      8715
$2.99   114
$0.99   106
$4.99    70
$1.99    59
Name: Price, dtype: int64
```

##### Observations

Some have dollars, some have 0

- We need to conditionally handle this.
- First, let's modify the column to take 0 if value is 0, else take the first letter onwards.

```
#Modifying the column
inp0['Price0'] = inp0.Price.map(lambda x: 0 if x=="0" else float(x[1:]))
```

##### The other columns with numeric data are:

1. Reviews
2. Installs
3. Size

###### 2. Convert reviews to numeric

```
#Converting reviews to numeric
inp0.Reviews = inp0.Reviews.astype("int32")
```

```
inp0.Reviews.describe()
```

```
count      9.360000e+03
mean       5.147675e+03
std        3.145023e+06
min        1.000000e+00
25%        1.867500e+02
50%        5.955000e+03
75%        8.162750e+04
max        7.012831e+07
Name: Reviews, dtype: float64
```

###### 3. Handle the installs column

```
#Handling the installs column
inp0.Installs.value_counts()
```

```
1,000,000+      1576
10,000,000+     1252
100,000+        1150
10,000+         1009
5,000,000+       752
1,000+          112
500,000+         537
50,000+          466
5,000+           421
100,000,000+     409
100+             309
50,000,000+      289
500+             201
500,000,000+     72
10+              69
1,000,000,000+   58
5+               56
1+               9
10+              3
Name: Installs, dtype: int64
```

We'll need to remove the commas and the plus signs.

###### Defining function for the same

```
def clean_installs(val):
    return int(val.replace(",","").replace("+",""))
```

```
inp0.Installs = inp0.Installs.map(clean_installs)
```

```
inp0.Installs.describe()
```

```
count      9.360000e+03
mean       1.790875e+07
std        9.126637e+07
min        1.000000e+00
25%        1.000000e+04
50%        5.000000e+05
75%        5.000000e+06
max        1.000000e+09
Name: Installs, dtype: float64
```

###### 4. Handle the app size field

```
#Handling the app size field
def change_size(size):
    if 'M' in size:
        x = size[:-1]
        x = float(x)*1000
        return(x)
    elif 'k' in size[:-1]:
        x = size[:-1]
        x = float(x)
        return(x)
    else:
        return None
```

```
inp0['Size'] = inp0['Size'].map(change_size)
```

```
inp0.Size.describe()
```

```
count      7723.000000
mean       22970.456105
std        23449.628935
min         0.000000
25%        8.500000
50%       14000.000000
75%       33000.000000
max       100000.000000
Name: Size, dtype: float64
```

```
#Filling size which had NA
inp0.Size.fillna(method='ffill', inplace=True)
```

```
#Checking datatypes
inp0.dtypes
```

```
App             object
Category        object
Rating          float64
Reviews         int32
Size            float64
Installs        int64
Type            object
Price          float64
Content_Rating  object
Genres          object
Last_Updated    object
Current_Ver     object
Android_Ver     object
dtype: object
```

##### Sanity checks

1. Average rating should be between 1 and 5, as only these values are allowed on Play Store. Drop any rows that have a value outside this range.

```
#Checking the rating
inp0.Rating.describe()
```

```
count      9360.000000
mean       4.191838
std        0.515263
min        1.000000
25%        4.000000
50%        4.300000
75%        4.500000
max        5.000000
Name: Rating, dtype: float64
```

##### Observations

Min is 1 and max is 5. None of the values have rating outside the range.

1. Reviews should not be more than 10 million as only those who installed can review the app.

Checking if reviews are more than installs. Counting total rows like this.

```
#Checking and counting the rows
len(inp0[inp0.Reviews > inp0.Installs])
```

```
7
```

```
inp0[inp0.Reviews > inp0.Installs]
```

	App	Category	Rating	Reviews	Size	Installs	Type	Price	Content_Rating	Genres	Last_Updated	Current_Ver	Android_Ver
2454	KBA-EZ Health Guide	MEDICAL	5.0	4	25000.0	1	Free	0.00	Everyone	Medical	August 2, 2018	1.0.72	4.0.3 and up
4663	Alarmy Sleep If U Can! - Pro	LIFESTYLE	4.8	10249	30000.0	10000	Paid	2.49	Everyone	Lifestyle	July 30, 2018	Varies with device	Varies with device
5917	Ra Ga Ba	GAME	5.0	2	20000.0	1	Paid	1.49	Everyone	Arcade	February 8, 2017	1.0.4	2.3 and up
6700	Brick Breaker	GAME	5.0	7	19000.0	5	Free	0.00	Everyone	Arcade	July 23, 2018	1	4.1 and up
7400	Tromai se o resc	GAME	5.0	11	61000.0	10	Free	0.00	Everyone	Arcade	March 11, 2017	0.1	2.3 and up
8591	DN Blog	SOCIAL	5.0	20	42000.0	10	Free	0.00	Teen	Social	July 23, 2018	1	4.0 and up
10697	Mu.F.O.	GAME	5.0	2	16000.0	1	Paid	0.99	Everyone	Arcade	March 3, 2017	1	2.3 and up

```
inp0 = inp0[inp0.Reviews <= inp0.Installs].copy()
```

```
inp0.shape
```

```
(9353, 13)
```

1. For free apps (Type == "Free"), the price should not be > 0. Drop any such rows.

```
len(inp0[(inp0.Type == "Free") & (inp0.Price>0)])
```

```
0
```

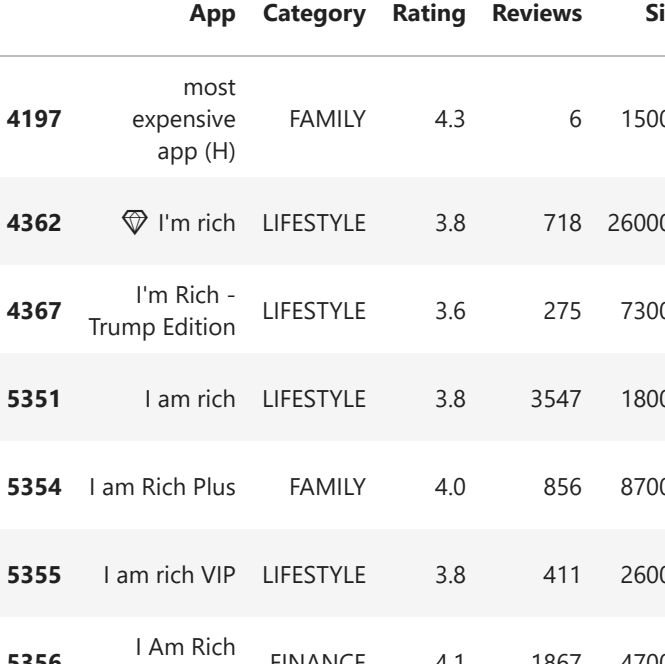
##### EDA

###### Box Plot: Price

Are there any outliers? Think about the price of usual apps on the Play Store.

```
sns.boxplot(inp0.Price)
plt.show()
```

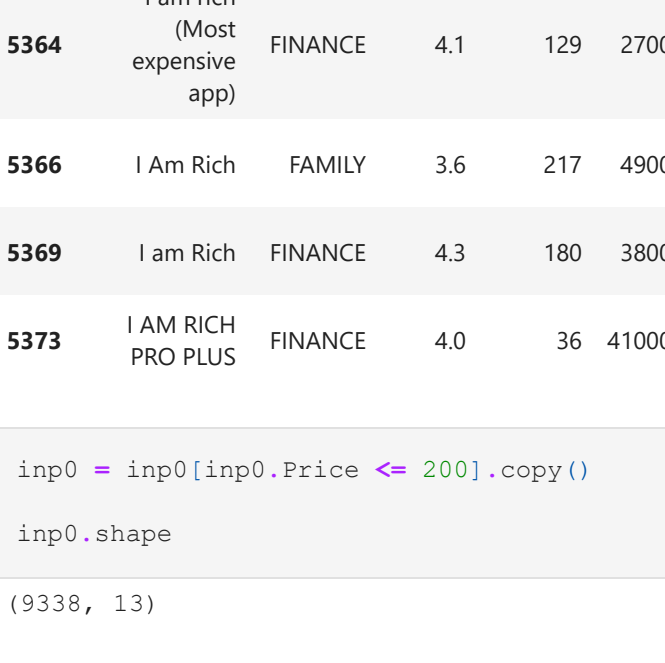
```
C:\Users\alpiika.gupta\Anaconda3\lib\site-packages\seaborn\decorators.py:36: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be "data", and passing other arguments without an explicit keyword will result in an error or misinterpretation.
warnings.warn(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.'))
```



###### Box Plot: Reviews

Are there any apps with very high number of reviews? Do the values seem right?

```
sns.boxplot(inp0.Reviews)
plt.show()
```

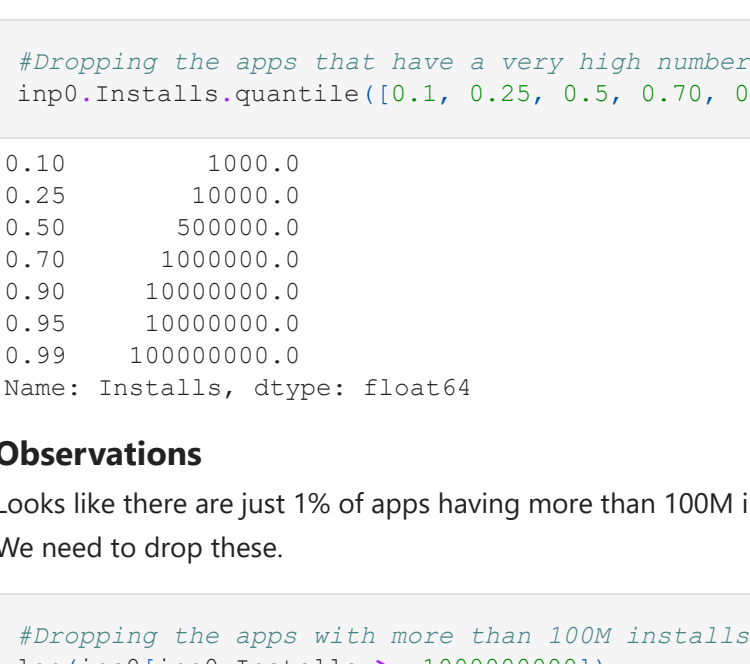


##### Checking Distribution and Skewness:

How are the ratings distributed? Is it more toward higher ratings?

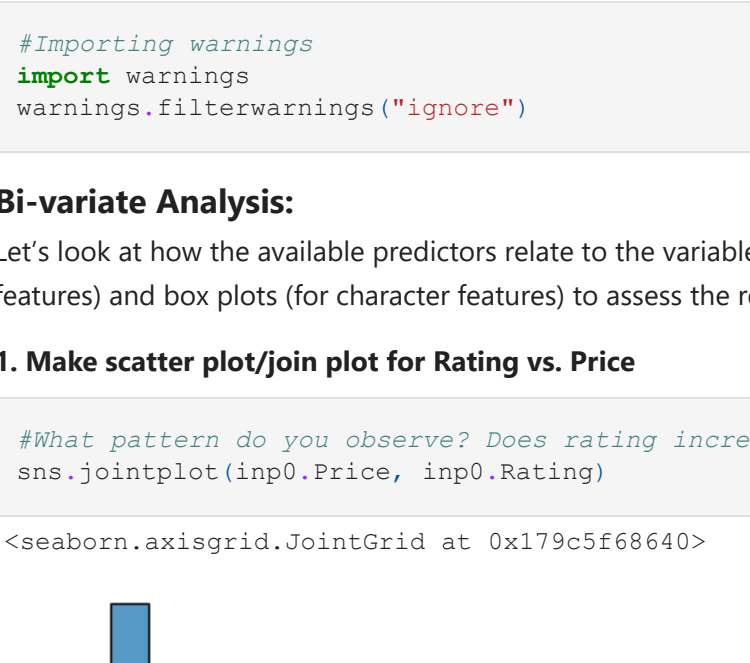
##### Distribution of Ratings

```
#distributing the ratings
inp0.Rating.plot.hist()
#Show plot
plt.show()
```



##### Histogram: Size

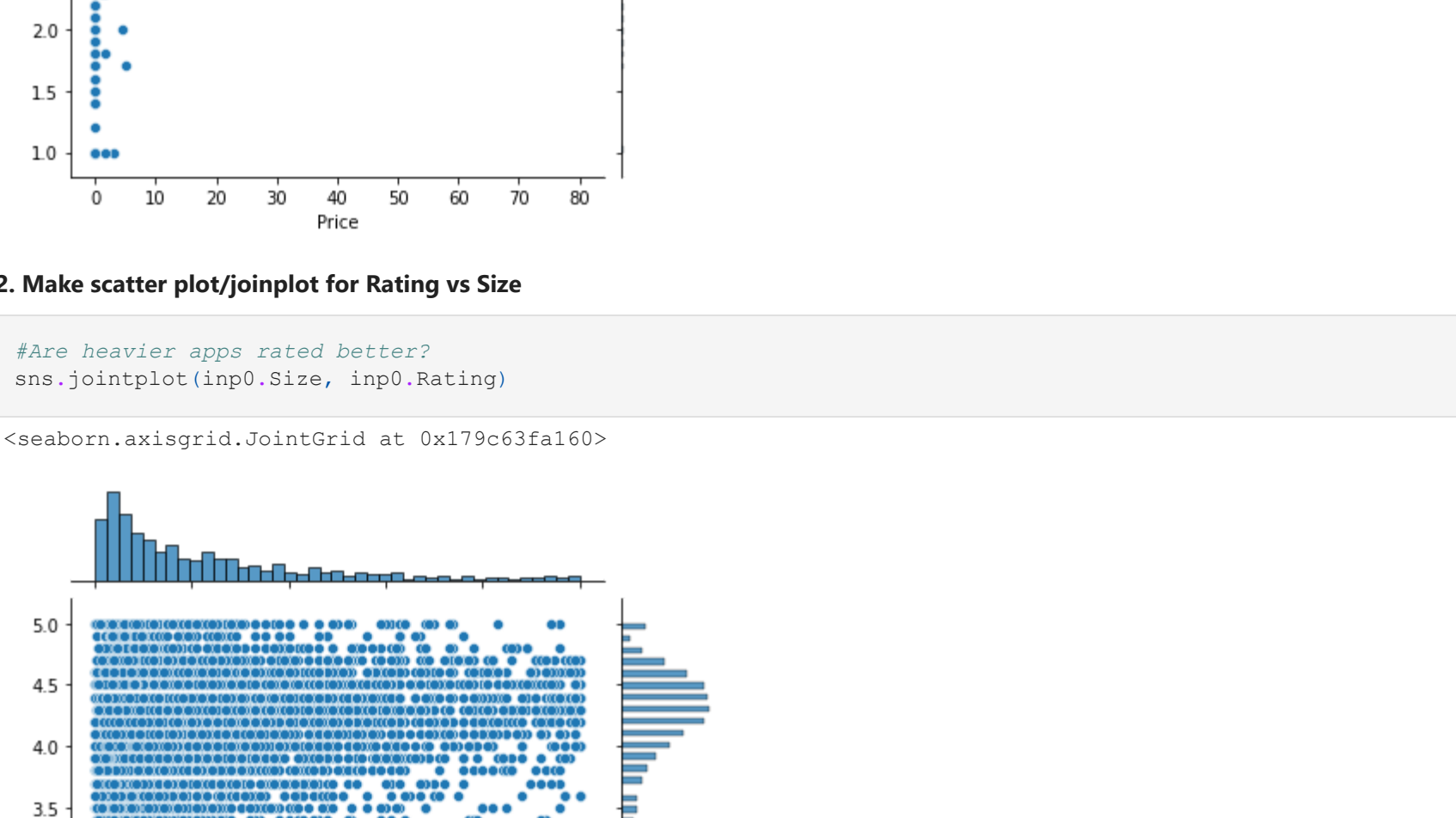
```
inp0['Size'].plot.hist()
#Show plot
plt.show()
```



##### Observations

A histogram is plotted with ratings on the x-axis and frequency on the y-axis, and the ratings are distributed.

```
#Pair plot
sns.pairplot(data=inp0)
```



##### Outlier Treatment:

###### 1. Price:

From the box plot, it seems like there are some apps with very high prices. A price of \$200 for an application on the Play Store is very high and suspicious. Check the records that have very high price: Is 200 a high price?

```
#Checking the records
len(inp0[inp0.Price > 200])
```

```
15
```

```
inp0[inp0.Price > 200]
```

	App	Category	Rating	Reviews	Size	Installs	Type	Price	Content_Rating	Genres	Last_Updated	Current_Ver	Android_Ver
4197	most expensive app (H)	FAMILY	4.3	6	1500.0	100	Paid	399.99	Everyone	Entertainment	July 16, 2018	1	7.0 and up
4362	I'm rich	LIFESTYLE	3.8	718	26000.0	10000	Paid	399.99	Everyone	Lifestyle	March 11, 2018	1.0.0	4.4 and up
4367	I'm Rich - Trump Edition	LIFESTYLE	3.6	275	7300.0	10000	Paid	400.00	Everyone	Lifestyle	May 3, 2018	1.0.1	4.1 and up
5351	I am Rich	LIFESTYLE	3.8	3547	1800.0	100000	Paid	399.99	Everyone	Lifestyle	January 12, 2018	2	4.0.3 and up
5354	I am Rich Plus	FAMILY	4.0	856	8700.0	10000	Paid	399.99	Everyone	Entertainment	May 19, 2018	3	4.4 and up
5355	I am rich VIP	LIFESTYLE	3.8	411	2600.0	10000	Paid	399.99	Everyone	Lifestyle	July 21, 2017	1.1.1	4.3 and up
5356	I Am Rich Premium	FINANCE	4.1	1867	4700.0	50000	Paid	399.99	Everyone	Finance	November 12, 2017	1.6	4.0 and up
5357	I am extremely Rich	LIFESTYLE	2.9	41	2900.0	1000	Paid	379.99	Everyone	Lifestyle	July 1, 2017	1	4.0 and up
5358	I am Rich!	FINANCE	3.8	93	2200.0	1000	Paid	399.99	Everyone	Finance	December 11, 2017	1	4.1 and up
5359	I am rich (premium)	FINANCE	3.5	472	965.0	5000	Paid	399.99	Everyone	Finance	May 1, 2017	3.4	4.4 and up
5362	I Am Rich Pro	FAMILY	4.4	201	2700.0	5000	Paid	399.99	Everyone	Entertainment	July 30, 2017	1.54	1.6 and up
5364	I am rich (Most expensive app)	FINANCE	4.1	129	2700.0	1000	Paid	399.99	Teen	Finance	December 6, 2017	2	4.0.3 and up
5366	I Am Rich	FAMILY	3.6	217	4900.0	10000	Paid	389.99	Everyone	Entertainment	June 22, 2018	1.5	4.2 and up
5369	I am Rich	FINANCE	4.3	180	3800.0	5000	Paid	399.99	Everyone	Finance	March 22, 2018	1	4.2 and up
5373	I AM RICH PRO PLUS	FINANCE	4.0	36	4100.0	1000	Paid	399.99	Everyone	Finance	June 25, 2018	1.0.2	4.1 and up

```
inp0 = inp0[inp0.Price <= 200].copy()
```

```
inp0.shape
```

```
(9353, 13)
```

##### 2. 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.

```
#Dropping the records with more than 2 million reviews
inp0 = inp0[inp0.Reviews <= 2000000]
```

```
inp0.shape
```

```
(8885, 13)
```

##### 3. Installs:

There seem to be some outliers in this field too. Apps having a very high number of installs should be dropped from the analysis. Find out the different percentiles - 10, 25, 50, 70, 90, 95, 99.

Decide a threshold as the cutoff for outliers and drop records having values more than the threshold.

```
#Dropping the apps that have a very high number of installs
inp0.Installs.quantile([0.1, 0.25, 0.5, 0.75, 0.9, 0.95, 0.99])
```

```
0.10      1000.0
0.25     10000.0
0.50    500000.0
0.70   1000000.0
0.90   1000000.0
0.95   1000000.0
0.99   1000000.0
Name: Installs, dtype: float64
```

##### Observations

Looks like there are just 1% of apps having more than 100M installs. These apps might be genuine, but will definitely skew our analysis. We need to drop these.

```
#Dropping the apps with more than 100M installs
len(inp0[inp0.Installs >= 100000000])
```

```
6
```

```
inp0 = inp0[inp0.Installs < 100000000].copy()
```

```
inp0.shape
```

```
(8879, 13)
```

##### Bi-variate Analysis:

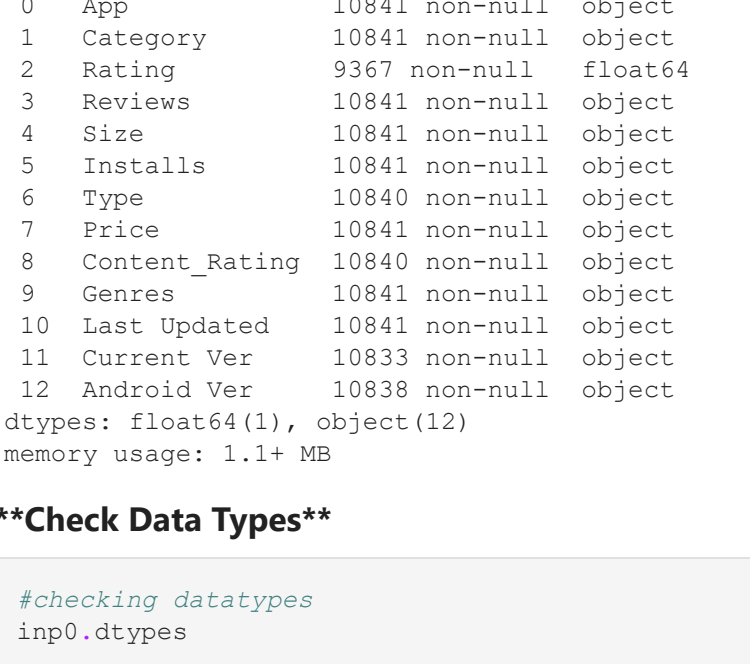
Let's look at how the available predictors relate to the variable of interest, i.e., our target variable rating. Make scatter plots (for numeric features) and box plots (for character features) to assess the relationships between rating and the other features.

###### 1. Make scatter plot/joint plot for Rating vs. Price

```
#What pattern do you observe? Does rating increase with price?
```

```
sns.jointplot(inp0.Price, inp0.Rating)
```

```
<seaborn.axisgrid.JointGrid at 0x179c6f8640>
```

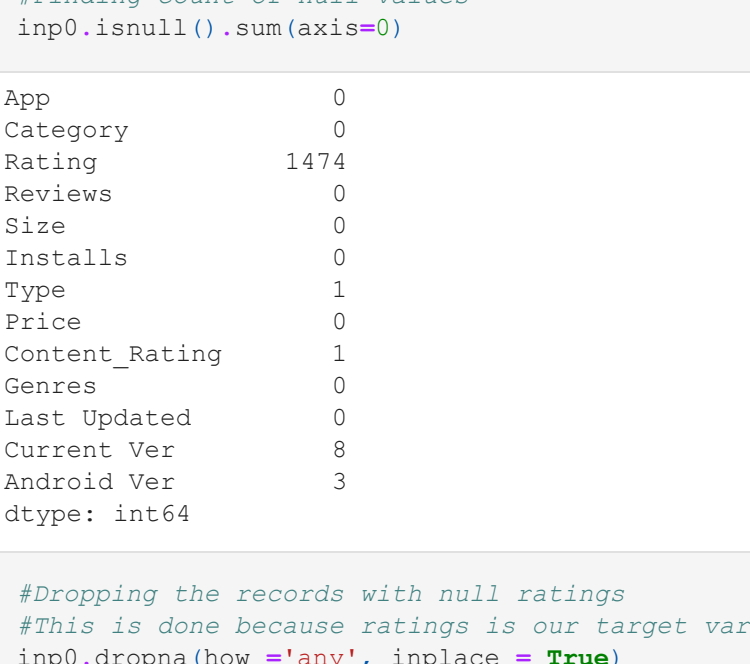


###### 2. Make scatter plot/jointplot for Rating vs Size

```
#Are heavier apps rated better?
```

```
sns.jointplot(inp0.Size, inp0.Rating)
```

```
<seaborn.axisgrid.JointGrid at 0x179c6f8fa60>
```

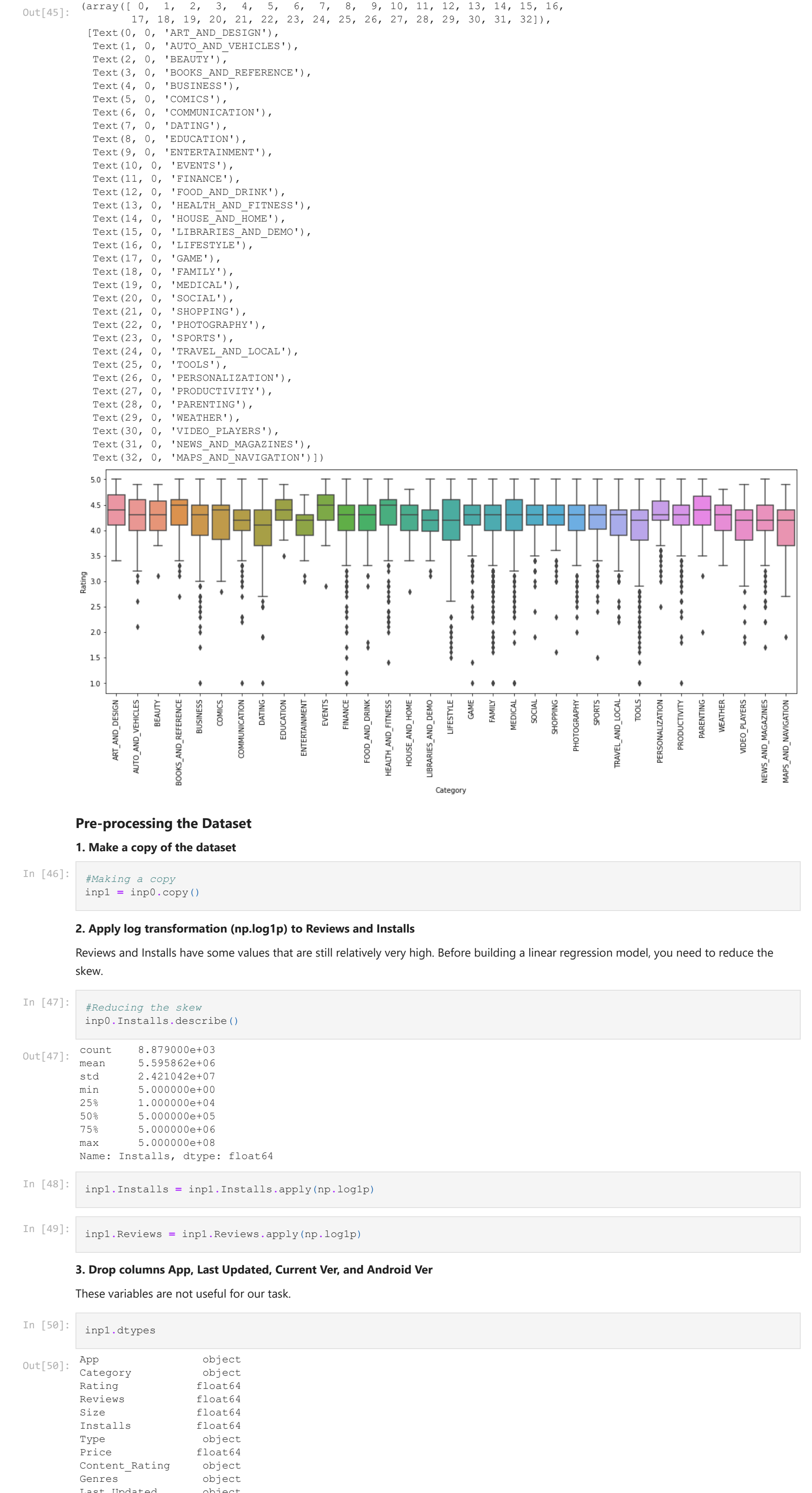


###### 3. Make scatter plot/jointplot for Rating vs Reviews

```
# Does more review mean a better rating always?
```

```
sns.jointplot(inp0.Reviews, inp0.Rating)
```





## Pre-processing the Dataset

### 1. Make a copy of the dataset

```
In [46]: #making a copy
inp1 = inp0.copy()
```

### 2. Apply log transformation (np.log1p) to Reviews and Installs

Reviews and Installs have some values that are still relatively very high. Before building a linear regression model, you need to reduce the skew.

```
In [47]: #Reducing the skew
inp0.Installs.describe()
```

```
Out[47]:
count      8.879000e+03
mean       5.595862e+06
std        2.421042e+07
min        5.000000e+00
25%        1.000000e+04
50%        5.000000e+05
75%        5.000000e+06
max        5.000000e+08
Name: Installs, dtype: float64
```

```
In [48]: inp1.Installs = inp1.Installs.apply(np.log1p)
```

```
In [49]: inp1.Reviews = inp1.Reviews.apply(np.log1p)
```

### 3. Drop columns App, Last Updated, Current Ver, and Android Ver

These variables are not useful for our task.

```
In [50]: inp1.dtypes
```

```
Out[50]:
App                object
Category           object
Rating            float64
Reviews           float64
Size             float64
Installs          float64
Type             float64
Price            float64
Content_Rating    object
Genres           object
Last_Updated     object
Current_Ver      object
Android_Ver      object
dtype: object
```

```
In [51]: #Dropping the variables that are not useful for our task
inp1.drop(["App", "Last_Updated", "Current_Ver", "Android_Ver"], axis=1, inplace=True)
inp1.shape
```

```
Out[51]:
(8879, 9)
```

### 4. Dummy Columns:

Get dummy columns for Category, Genres, and Content Rating. This needs to be done as the models do not understand categorical data, and all data should be numeric. Dummy encoding is one way to convert character fields to numeric fields. Name of the dataframe should be **inp2**.

```
In [52]: inp2 = pd.get_dummies(inp1, drop_first=True)
```

```
In [53]: inp2.columns
```

```
Out[53]:
Index(['Rating', 'Reviews', 'Size', 'Installs', 'Price',
       'Category.AUTO_AND_VEHICLES', 'Category.BEAUTY',
       'Category.BOOKS_AND_REFERENCE', 'Category.COMICS',
       ...,
       'Genres.Tools', 'Genres.Tools_Education', 'Genres.Travel_&Local',
       'Genres.Travel_&Local_Action_&Adventure', 'Genres.Trivia',
       'Genres.Video_Players_&Editors',
       'Genres.Video_Players_&Editors_Creativity',
       'Genres.Video_Players_&Editors_Music_&Video', 'Genres.Weather',
       'Genres.Words'],
      dtype='object', length=157)
```

```
In [54]: inp2.shape
```

```
Out[54]:
(8879, 157)
```

### Train-test split

Let us distribute the data into **training** and **test** datasets using the **train\_test\_split()** function.

```
In [55]: from sklearn.model_selection import train_test_split
```

```
In [56]: #train_test_split
```

```
In [57]: df_train, df_test = train_test_split(inp2, train_size = 0.7, random_state = 100)
```

```
In [58]: df_train.shape, df_test.shape
```

```
Out[58]:
((6215, 157), (2664, 157))
```

Let us separate the dataframes into **X\_train, y\_train, X\_test, y\_test**.

```
In [59]: y_train = df_train.pop("Rating")
X_train = df_train
```

```
In [60]: X_train.head()
```

```
Out[60]:
  Reviews    Size    Installs  Price  Category.AUTO_AND_VEHICLES  Category.BEAUTY  Category.BOOKS_AND_REFERENCE  Category.BUSINESS
1279  7.63627  4900.0  11.512935   0.0                        0                        0                        0
```

1 rows × 156 columns

```
In [61]: y_test = df_test.pop("Rating")
X_test = df_test
```

```
In [62]: X_test.head()
```

```
Out[62]:
  Reviews    Size    Installs  Price  Category.AUTO_AND_VEHICLES  Category.BEAUTY  Category.BOOKS_AND_REFERENCE  Category.BUSINESS
1161  9.329456  24000.0  13.815512   0.0                        0                        0                        0
```

1 rows × 156 columns

## \*\*Regression Algorithms:\*\*

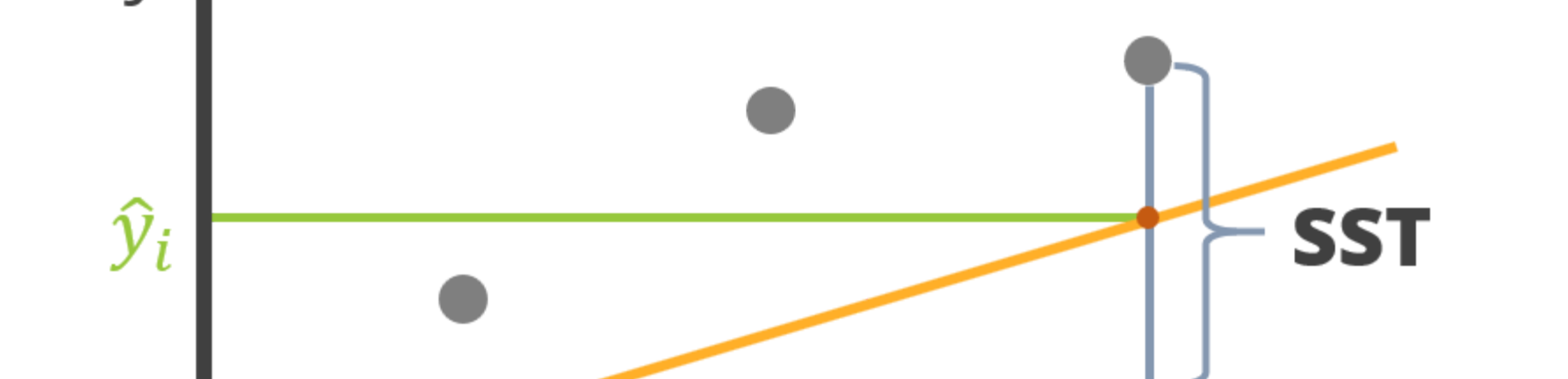
Note: Let us take a look at the theory part before moving on to the training and prediction.

### Types of Regression Algorithms:

- Linear regression
- Multiple linear regression
- Polynomial regression
- Ridge regression
- Lasso regression
- ElasticNet regression

### When to use regression?

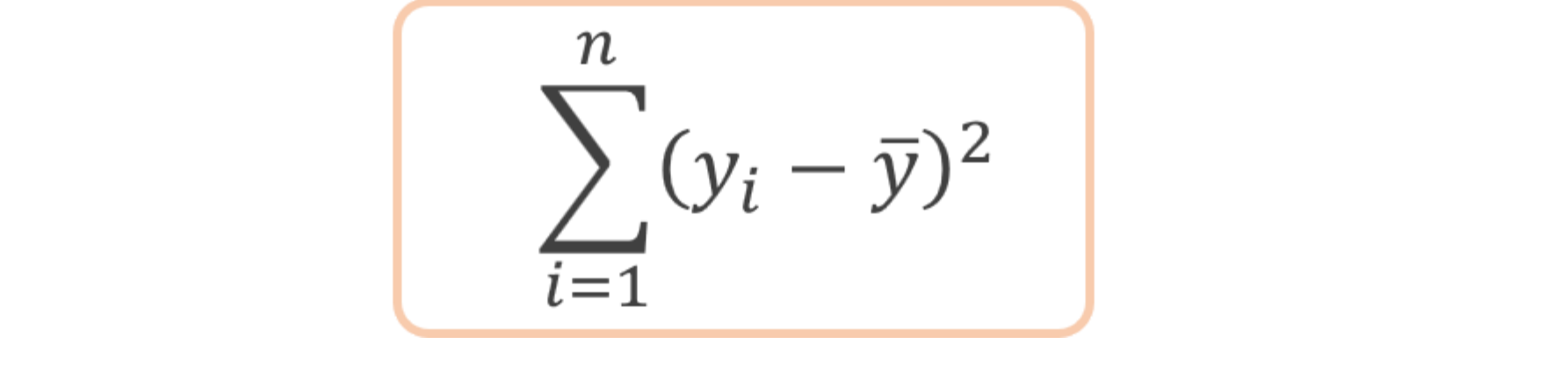
If target variable is a continuous numeric variable (100–2000), then use a regression algorithm.



Note: Let us take a look at the basics of linear regression and then move on to the model building part where we are going to use all the concepts that we saw in previous sessions.

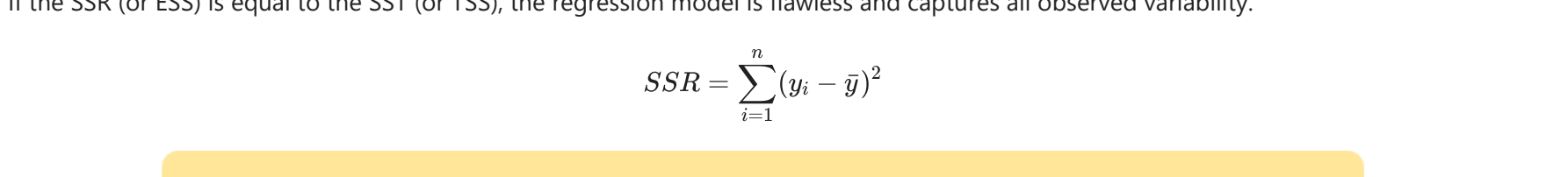
### 1. Linear Regression:

Linear Regression is a statistical model used to predict the relationship between independent and dependent variables denoted by  $x$  and  $y$  respectively.



### 2. Multiple Linear Regression:

Multiple linear regression is a statistical technique used to predict the outcome of a response variable through several explanatory variables and model the relationships between them.



### 3. Polynomial Regression:

Polynomial regression is applied when data is not formed in a straight line. It is used to fit a linear model to non-linear data by creating new features from powers of non-linear features.

Example: Quadratic features

$$\begin{aligned} X' &= X^2 \\ Y &= w_1 X_1 + w_2 X_1^2 + \text{residual error} \\ Y &= w_1 X_1 + w_2 X_1^2 + \text{residual error} \end{aligned}$$



## Types of Model Evaluation Metrics:

### Assumption

Let us consider the following:

- $y_i$  – the observed value
- $\bar{y}$  – the mean value of a sample
- $\hat{y}_i$  – the value estimated by the regression line

### Sum of Squares Total (SST)

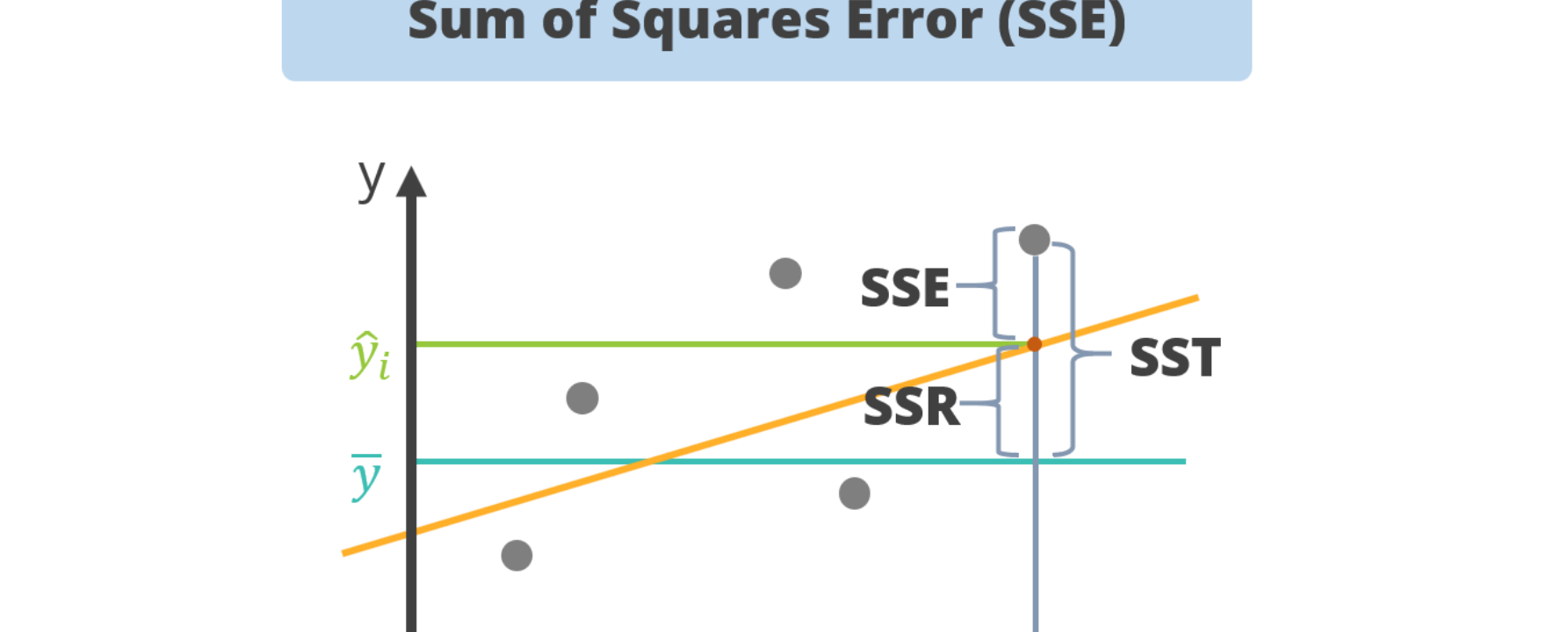
The squared variations between the measured dependent variable and its mean are referred to as the **Sum of Squares Total (SST)** or **Total Sum of Squares (TSS)**.

It's similar to the variation of descriptive statistics in that it's the dispersion of measured variables around the mean.

It is a measure of the dataset's overall variability.

$$SST = SSR + SSE = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 + \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

## Sum of Squares Total (SST)



### Sum of Squares due to Regression (SSR)

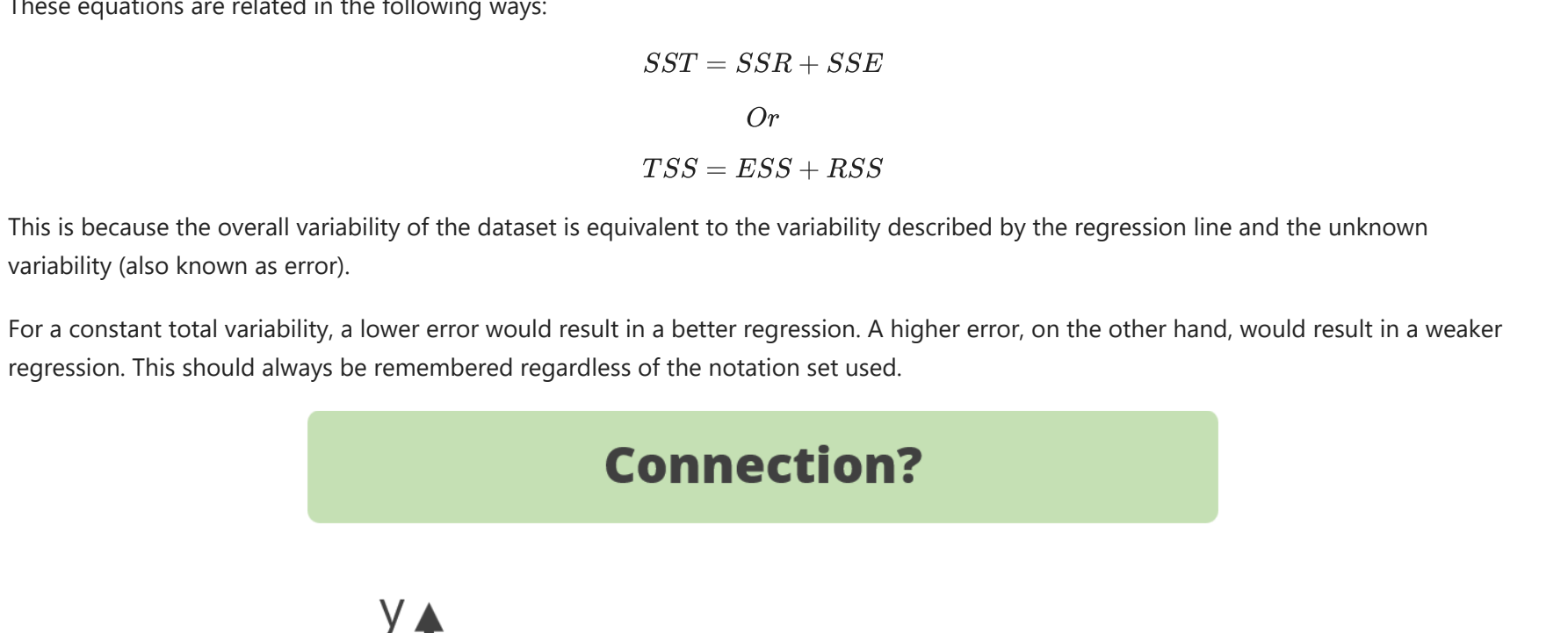
The difference between the predicted value and the dependent variable's mean are referred to as the **Sum of Squares due to Regression (SSR)** or **Explained Sum of Squares (ESS)**.

It can be considered as a metric for describing how well our line fits the data.

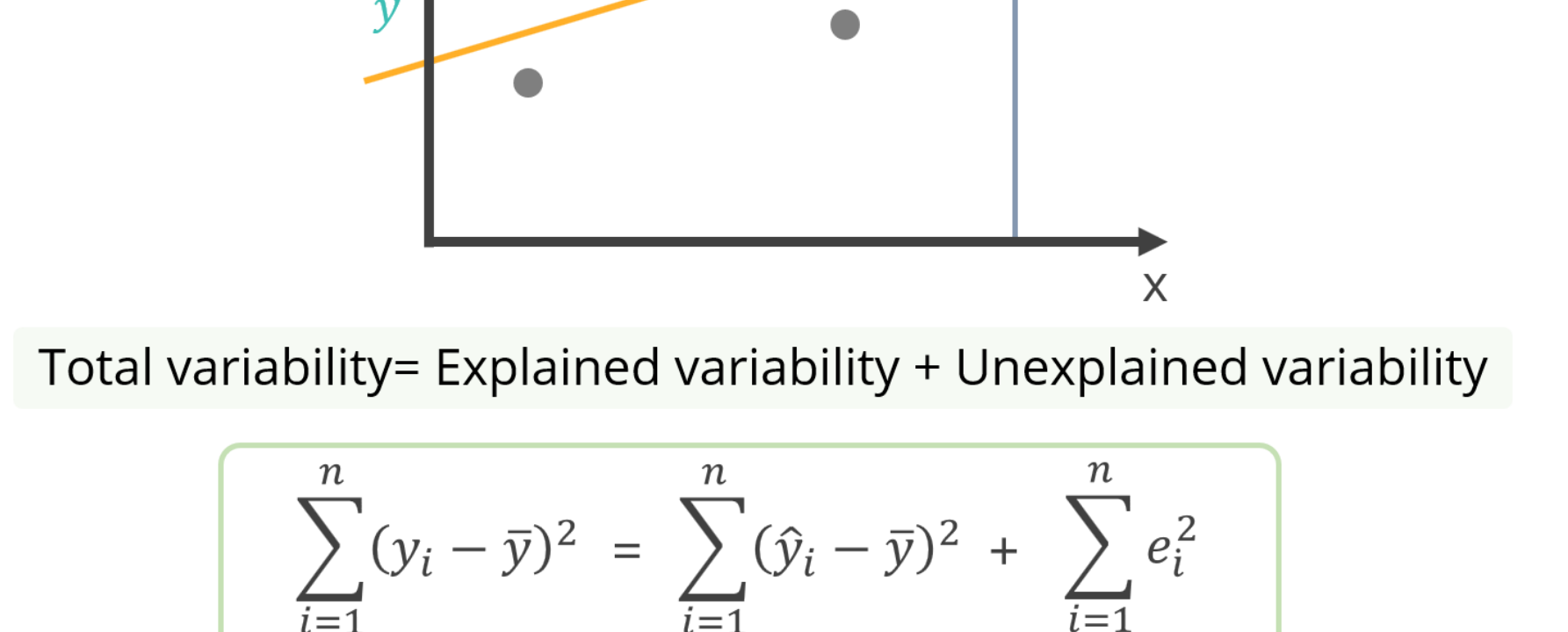
If the SSR (or ESS) is equal to the SST (or TSS), the regression model is flawless and captures all observed variability.

$$SSR = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$$

## Sum of Squares Regression (SSR)



## Measures the explained variability by your line



### Sum of Squares Error (SSE)

The difference between the observed and predicted values are referred to as the **Sum of Squares Error (SSE)** or **Residual Sum of Squares (RSS)**, where **residual** stands for **remaining** or **unexplained**.

This error must be reduced since the smaller it is, the better the regression's estimation power.

$$SSE = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Or

$$SSE = \sum_{i=1}^n e_i^2$$

where,

$$e_i = (y_i - \hat{y}_i)$$

## Sum of Squares Error (SSE)



## Measures the unexplained variability by the regression



### Relation Among SST, SSR, and SSE

Since certain people use these abbreviations in various ways, it can be very confusing.

We use one of two sets of notations for these abbreviations: SST, SSR, and SSE or TSS, ESS, and RSS.

These equations are related in the following ways:

$$SST = SSR + SSE$$

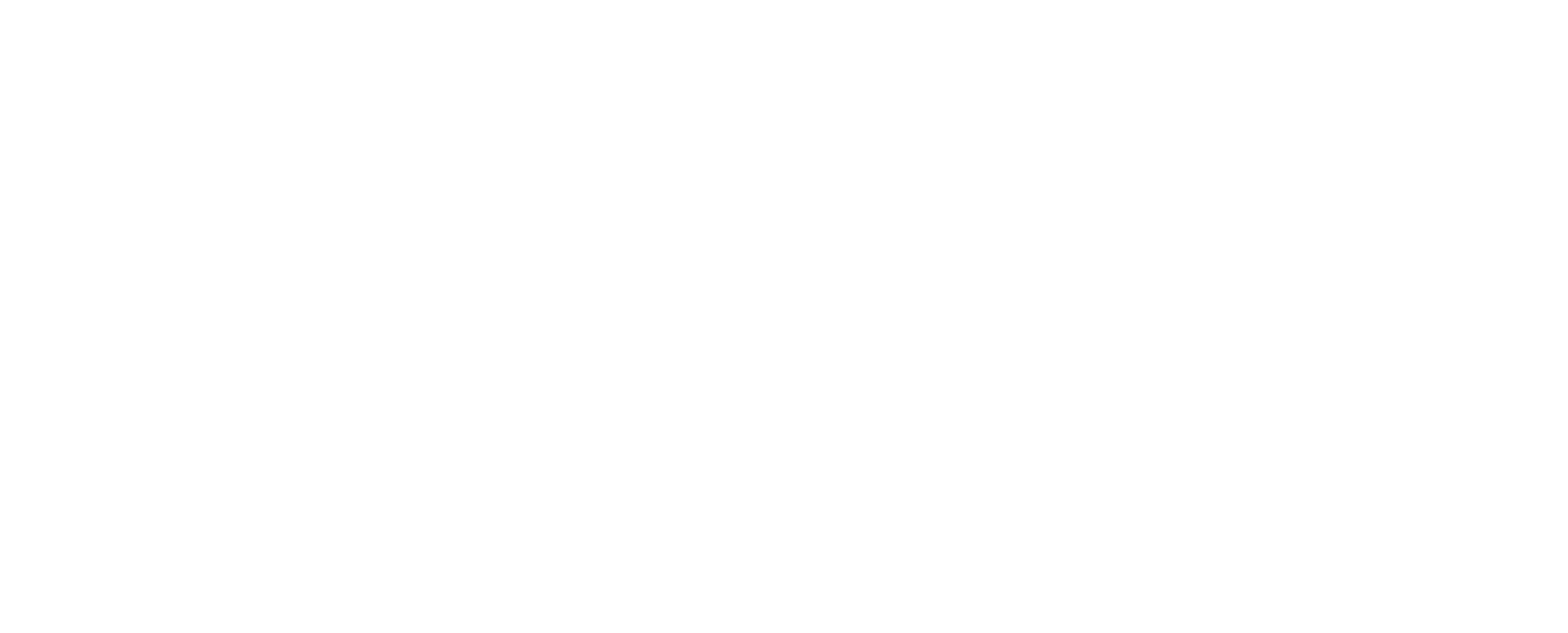
Or

$$TSS = ESS + RSS$$

This is because the overall variability of the dataset is equivalent to the variability described by the regression line and the unknown variability (also known as error).

For a constant total variability, a lower error would result in a better regression. A higher error, on the other hand, would result in a weaker regression. This should always be remembered regardless of the notation set used.

## Connection?



## Total variability= Explained variability + Unexplained variability

$$\sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 + \sum_{i=1}^n e_i^2$$

### R-Square Matrix

The determination coefficient also known as **R2 (R-squared) score** is used for the performance evaluation of a linear regression model.

R2 displays the proportion of data points inside the regression equation line.

A higher R2 value means improved results.

It is calculated as follows:

$$R^2 = 1 - \frac{SSE}{SSR}$$

Or

$$R^2 = 1 - \frac{RSS}{ESS}$$

The highest possible score is 1, which is achieved when the predicted and actual values are the same.

The R2 score is 0 for a baseline model.

In the worst-case scenario, the R2 score can also be negative.

### Import statsmodels Library for Linear Regression

```
In [63]: #Importing the statsmodels library
import statsmodels.api as sm
```

```
In [64]: #Applying linear regression
modell = sm.OLS(y_train, X_train)
```

```
In [65]: modell = modell.fit()
```

```
In [66]: #Finding the summary
modell.summary()
```



<

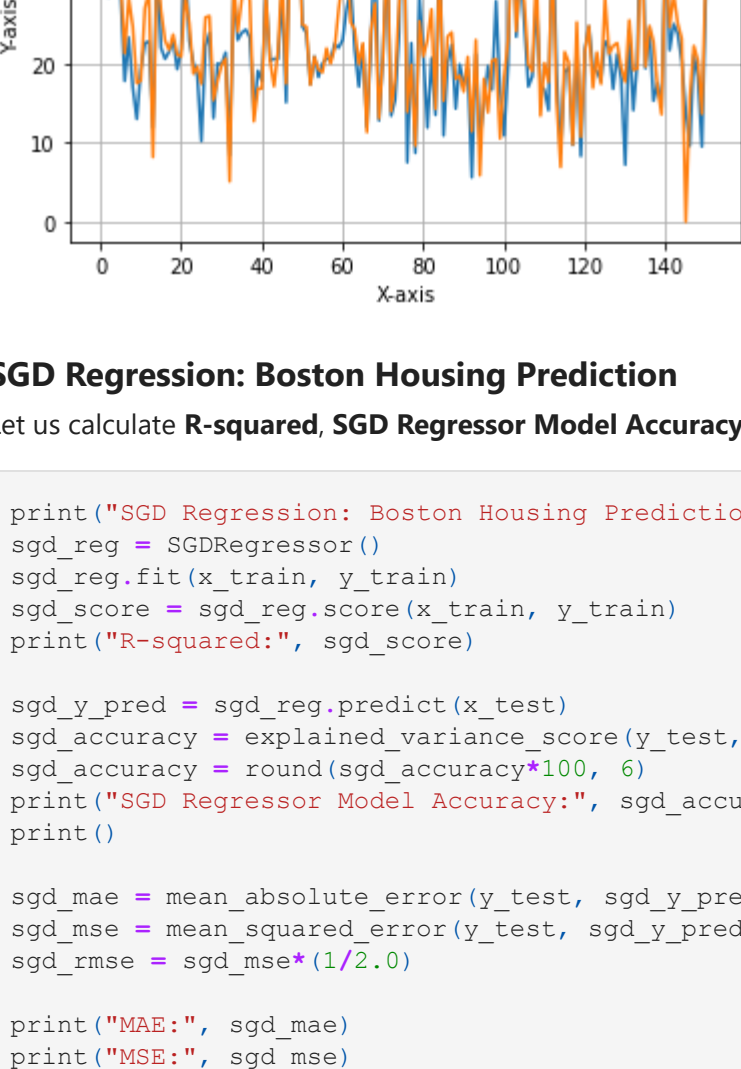






In [106]

```
x_ax = range(len(y_test))
plt.plot(x_ax, y_test, label="original")
plt.plot(x_ax, lin_y_pred, label="predicted")
plt.title("Test vs. Predicted Data")
plt.xlabel("X-axis")
plt.ylabel("Y-axis")
plt.legend(loc="best", fancybox=True, shadow=True)
plt.grid(True)
plt.show()
```



## SGD Regression: Boston Housing Prediction

Let us calculate R-squared, SGD Regressor Model Accuracy, MAE, MSE and RMSE.

In [107]

```
print("SGD Regression: Boston Housing Prediction")
sgd_reg = SGDRegressor()
sgd_reg.fit(x_train, y_train)
sgd_score = sgd_reg.score(x_train, y_train)
print("R-squared:", sgd_score)

sgd_y_pred = sgd_reg.predict(x_test)
sgd_accuracy = explained_variance_score(y_test, sgd_y_pred)
sgd_accuracy = round(sgd_accuracy*100, 6)
print("SGD Regressor Model Accuracy:", sgd_accuracy, "%")
print()

sgd_mae = mean_absolute_error(y_test, sgd_y_pred)
sgd_mse = mean_squared_error(y_test, sgd_y_pred)
sgd_rmse = sgd_mse*(1/2.0)

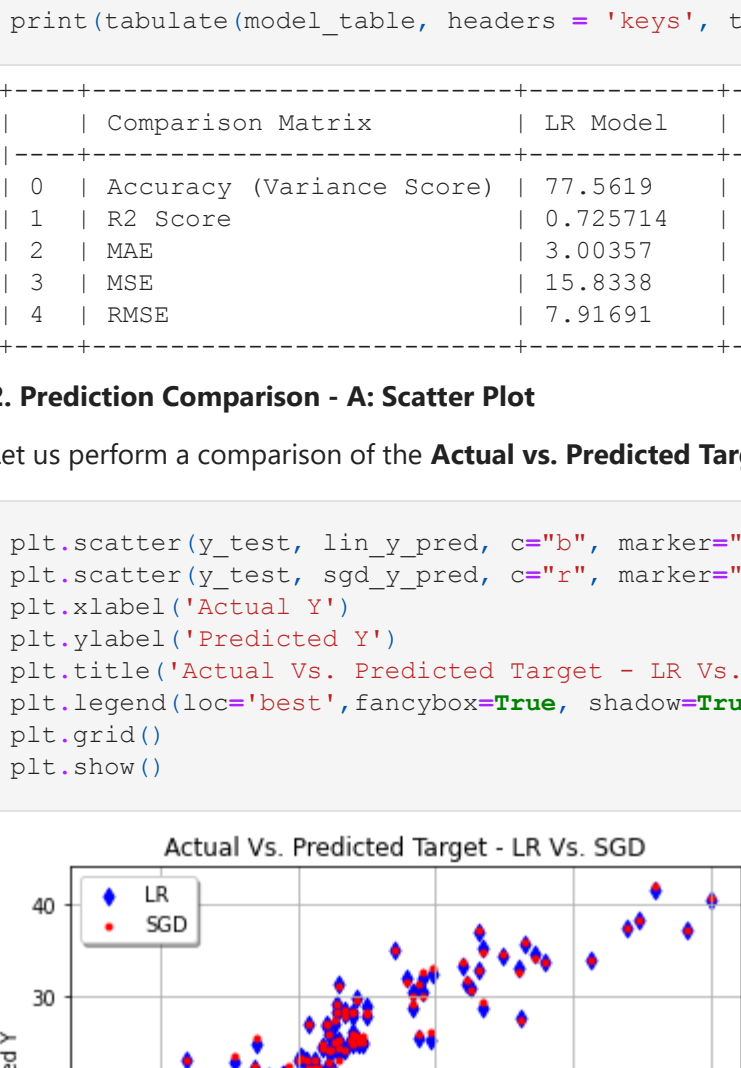
print("MAE:", sgd_mae)
print("MSE:", sgd_mse)
print("RMSE:", sgd_rmse)
print()
```

SGD Regression: Boston Housing Prediction  
R-squared: 0.722884803216978  
SGD Regressor Model Accuracy: 77.132392 %  
  
MAE: 3.026066392810071  
MSE: 16.00535293654782  
RMSE: 8.00267466827391

Let us plot a graph of the Actual vs. Predicted Target.

In [108]

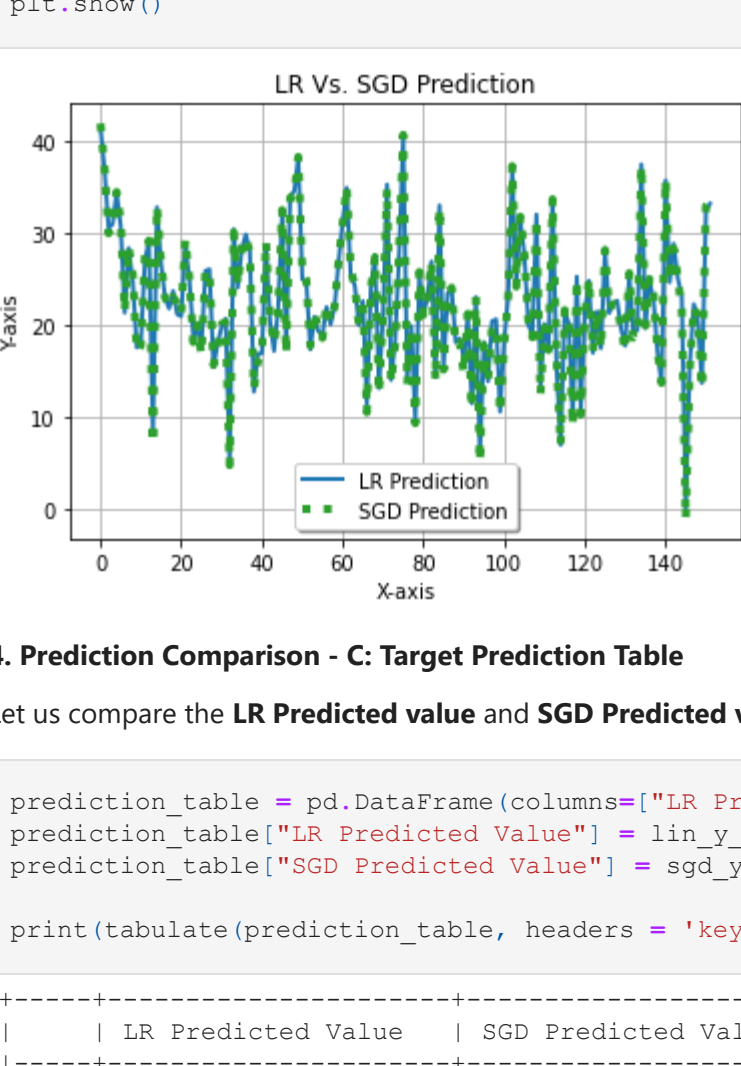
```
plt.scatter(y_test, sgd_y_pred)
plt.grid()
plt.xlabel("Actual Y")
plt.ylabel("Predicted Y")
plt.title("Actual vs. Predicted Target")
plt.show()
```



Let us plot a graph of the Actual vs. Predicted Data.

In [109]

```
x_ax = range(len(sgd_y_pred))
plt.plot(x_ax, y_test, label="original")
plt.plot(x_ax, sgd_y_pred, label="predicted")
plt.title("Test vs. Predicted Data")
plt.xlabel("X-axis")
plt.ylabel("Y-axis")
plt.legend(loc="best", fancybox=True, shadow=True)
plt.grid(True)
plt.show()
```



## Model Comparison

### 1. Evaluation Matrix Comparison

In [110]

```
#Performing evaluation matrix comparison
model_table = pd.DataFrame(columns = ["Comparison Matrix", "LR Model", "SGD Model"])
model_table["Comparison Matrix"] = ["Accuracy (Variance Score)", "R2 Score", "MAE", "RMSE"]
model_table["LR Model"] = [lin_accuracy, lin_score, lin_mae, lin_rmse]
model_table["SGD Model"] = [sgd_accuracy, sgd_score, sgd_mae, sgd_rmse]

print(tabulate(model_table, headers = 'keys', tablefmt = 'psql', numalign="left"))
```

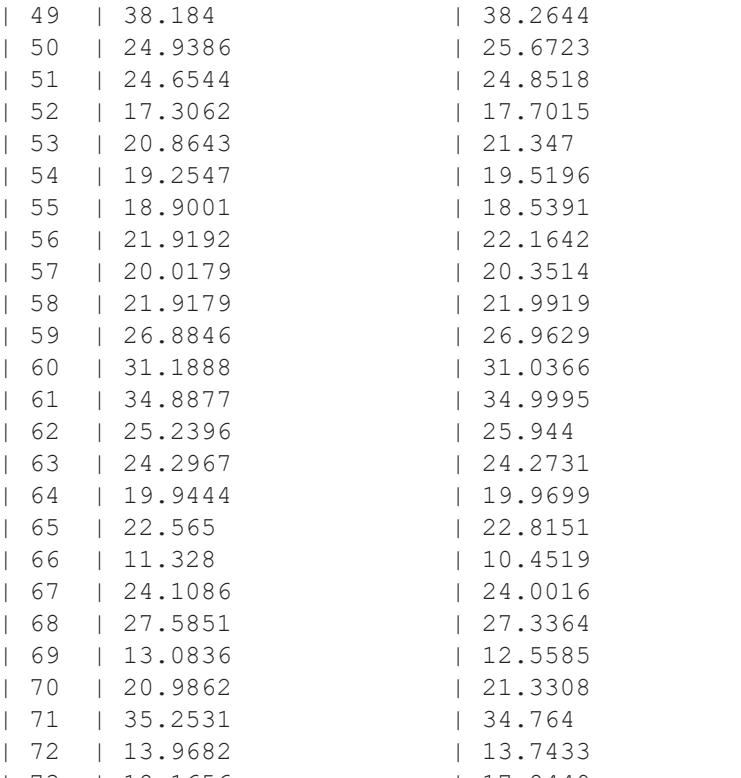
	Comparison Matrix	LR Model	SGD Model
0	Accuracy (Variance Score)	77.5619	77.1324
1	R2 Score	0.723714	0.722788
2	MAE	3.00357	3.02607
3	MSE	15.8338	16.0054
4	RMSE	7.91691	8.00268

### 2. Prediction Comparison - A: Scatter Plot

Let us perform a comparison of the Actual vs. Predicted Target for the Scatter Plot.

In [111]

```
plt.scatter(y_test, lin_y_pred, c="b", marker="d", label="LR")
plt.scatter(y_test, sgd_y_pred, c="r", marker="d", label="SGD")
plt.xlabel("Actual Y")
plt.ylabel("Predicted Y")
plt.title("Actual Vs. Predicted Target - LR Vs. SGD")
plt.legend(loc="best", fancybox=True, shadow=True)
plt.grid()
plt.show()
```

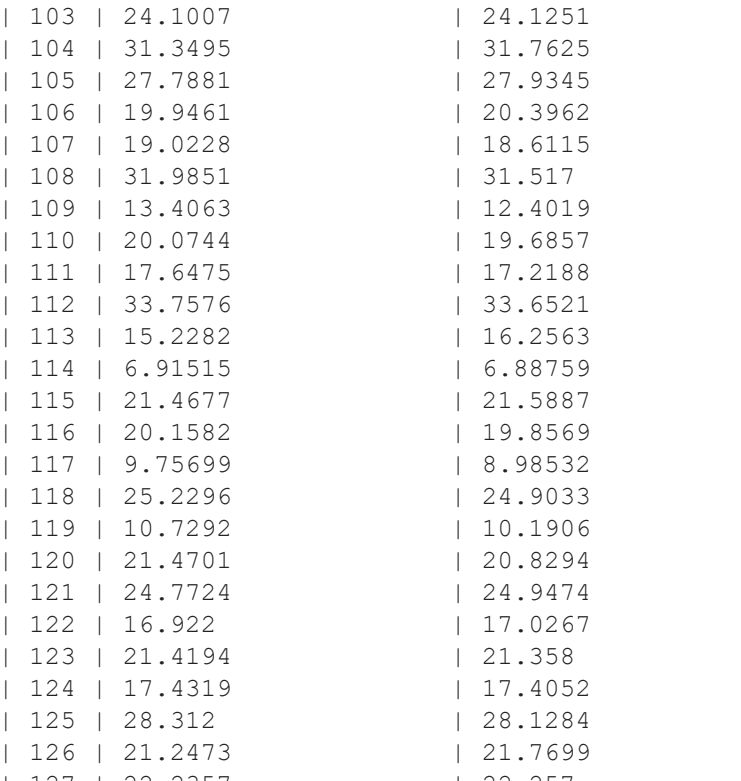


### 3. Prediction Comparison - B: Line Graph

Let us perform a comparison of LR vs. SGD Prediction for the Line Graph.

In [112]

```
x_ax = range(len(sgd_y_pred))
plt.plot(x_ax, lin_y_pred, c="b", linestyle="--", linewidth = 2, label="LR Prediction")
plt.plot(x_ax, sgd_y_pred, c="r", linestyle="--", linewidth = 4, label="SGD Prediction")
plt.xlabel("X-axis")
plt.ylabel("Y-axis")
plt.legend(loc="best", fancybox=True, shadow=True)
plt.grid(True)
plt.show()
```



### 4. Prediction Comparison - C: Target Prediction Table

Let us compare the LR Predicted value and SGD Predicted value.

In [113]

```
prediction_table = pd.DataFrame(columns=["LR Predicted Value", "SGD Predicted Value"])
prediction_table["LR Predicted Value"] = lin_y_pred
prediction_table["SGD Predicted Value"] = sgd_y_pred

print(tabulate(prediction_table, headers = 'keys', tablefmt = 'psql', numalign="left"))
```

	LR Predicted Value	SGD Predicted Value
0	41.4908	41.6855
1	36.8945	37.1855
2	30.5095	29.918
3	30.9245	30.6015
4	34.409	34.4465
5	30.7754	31.3239
6	21.3008	29.2663
7	28.3381	28.3853
8	24.6874	25.4813
9	17.5071	17.1117
10	17.8427	17.7045
11	27.0118	27.3565
12	29.0729	29.5653
13	8.16764	7.0879
14	32.7771	32.8589
15	26.6931	26.5848
16	22.9233	22.9401
17	21.7455	22.097
18	23.6717	22.9054
19	21.1214	22.3216
20	20.9091	21.1535
21	28.743	29.2663
22	25.1269	25.1679
23	18.7371	18.0881
24	19.6255	19.2774
25	17.4443	16.7471
26	25.8447	25.8276
27	26.0272	25.2661
28	15.3838	15.6016
29	17.8436	17.6804
30	19.7745	20.0748
31	20.6521	20.2664
32	5.09075	4.70043
33	30.375	30.4777
34	24.0346	24.3234
35	28.186	28.3361
36	29.8191	29.2664
37	28.0558	27.885
38	12.6971	13.0116
39	16.8389	16.0123
40	16.853	16.7836
41	28.6068	29.0329
42	20.4021	20.6939
43	17.1312	17.9518
44	21.7044	21.9931
45	31.8228	32.4798
46	17.4694	16.845
47	33.9607	33.7991
48	34.5137	34.1487
49	38.184	38.2644
50	24.9386	25.6723
51	24.6544	24.6518
52	17.3062	17.7015
53	20.8643	21.347
54	19.2547	22.0546
55	16.9001	18.5391
56	21.9192	22.1642
57	20.0179	20.2514
58	21.9179	21.9919
59	26.8846	26.9629
60	31.1888	31.2664
61	34.8877	34.9995
62	25.2396	25.944
63	24.2967	24.2791
64	19.9444	19.9699
65	22.565	22.8151
66	11.328	10.4519
67	24.1086	24.0016
68	27.5851	27.3364
69	12.0836	12.5583
70	20.9862	21.3308
71	35.2531	34.764
72	13.9682	13.7423
73	18.1656	17.8449
74	27.7294	28.2534
75	40.3915	40.6661
76	13.8945	13.5083
77	19.862	19.6349
78	9.59307	9.2559
79	25.3634	25.7403
80	20.9301	21.0666
81	22.9832	23.2188
82	26.8731	26.8218
83	14.2712	13.9445
84	32.9805	32.671
85	15.2945	15.0935
86	23.4149	23.2745
87	16.4444	16.0345
88	18.0822	17.7519
89	18.469	18.0179
90	16.4444	16.0345
91	20.8568	21.184
92	11.419	10.9252
93	22.9933	23.015
94	5.85954	5.06437
95	18.0939	17.7409
96	13.8261	13.9644
97	20.3623	19.7046
98	20.5843	19.7939
99	10.5214	10.446
100	18.2128	17.706
101	22.9752	23.0065
102	17.0495	17.2562
103	24.1007	24.1251
104	31.3495	31.7625
105	27.7881	27.6561
106	19.9461	20.3962
107	19.0228	18.6115
108	31.9831	31.517
109	13.4063	12.4019
110	20.0744	19.6857
111	17.4475	17.2188
112	33.7576	33.6521
113	15.2282	16.2563
114	6.9315	6.88759
115	21.4677	21.5887
116	20.1582	19.8569
117	9.75699	8.98532
118	25.2296	24.9033
119	10.7292	10.1906
120	21.4701	20.8294
121	24.7724	24.9474
122	16.922	17.0267
123	21.4194	21.358
124	17.4319	17.4052
125	28.312	28.1284
126	21.2473	21.7699
127	22.2357	22.257
128	22.6046	22.1832
129	19.1826	19.8423
130	17.6914	17.1919
131	25.4865	25.6177
132	19.1826	18.5633
133	19.3867	19.0175
134	37.4434	37.2716
135	19.633	19.8768
136	24.7732	24.9752
137	22.7821	23.3403
138	16.9901	17.1196
139	13.5727	13.6448
140	35.6925	35.7711
141	24.5964	24.1225
142	28.8186	28.2873
143	25.1047	25.2257
144	23.1659	21.6165
145	-0.0163374	-1.13048
146	15.6806	15.4564
147	22.2388	21.888
148	20.4303	20.6844
149	13.6457	13.1583
150	32.3988	32.0267
151	33.174	33.6081

Note: In this topic, we saw the use of the linear regression methods, but in the next topic we will be working on "Logistic Regression".

Powered by [simplilearn](#)