

## A.Importing the necessary Libraries

```
In [1]: 1
2 import pandas as pd
3 import numpy as np
4 import matplotlib.pyplot as plt
5 import seaborn as sns
6 import warnings
7 from sklearn.model_selection import train_test_split
8
9 # Ignore all warnings
10 warnings.filterwarnings("ignore")
```

## B.Data Loading

```
In [2]: 1 df = pd.read_csv("laptopPrice.csv")
2 df.head(5)
```

Out[2]:

	brand	processor_brand	processor_name	processor_gnrtn	ram_gb	ram_type	ssd	hdd	
0	ASUS	Intel	Core i3	10th	4 GB	DDR4	0 GB	1024 GB	v
1	Lenovo	Intel	Core i3	10th	4 GB	DDR4	0 GB	1024 GB	v
2	Lenovo	Intel	Core i3	10th	4 GB	DDR4	0 GB	1024 GB	v
3	ASUS	Intel	Core i5	10th	8 GB	DDR4	512 GB	0 GB	v
4	ASUS	Intel	Celeron Dual	Not Available	4 GB	DDR4	0 GB	512 GB	v

## C.EDA and Data Preprocessing

```
In [3]: 1 print(f"The total number of latops given are : {df.shape[0]}")
2 print(f"The number of feature columns are: {df.shape[1]}")
```

The total number of latops given are : 823  
The number of feature columns are: 19

In [4]:

```
1 """Quick summary of dataframe"""
2 df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 823 entries, 0 to 822
Data columns (total 19 columns):
#   Column                Non-Null Count  Dtype
---  -
0   brand                 823 non-null    object
1   processor_brand       823 non-null    object
2   processor_name        823 non-null    object
3   processor_gnrtn       823 non-null    object
4   ram_gb                823 non-null    object
5   ram_type              823 non-null    object
6   ssd                   823 non-null    object
7   hdd                   823 non-null    object
8   os                    823 non-null    object
9   os_bit                823 non-null    object
10  graphic_card_gb       823 non-null    object
11  weight                823 non-null    object
12  warranty              823 non-null    object
13  Touchscreen           823 non-null    object
14  msoffice              823 non-null    object
15  Price                 823 non-null    int64
16  rating                823 non-null    object
17  Number of Ratings     823 non-null    int64
18  Number of Reviews     823 non-null    int64
dtypes: int64(3), object(16)
memory usage: 122.3+ KB
```

In [5]:

```
1 """
2 Duplicate rows:
3 1. Check for duplicate rows
4 2. If present drop them
5 """
6
7 print(f"The Number of Duplicated rows are: {df.duplicated().sum()}")
```

The Number of Duplicated rows are: 21

In [6]:

```
1 """Duplicated rows"""
2 df[df.duplicated()].tail(3)
```

Out[6]:

	brand	processor_brand	processor_name	processor_gnrtn	ram_gb	ram_type	ssd	hdd
<b>605</b>	APPLE	M1	M1	10th	8 GB	DDR4	256 GB	0 GB
<b>616</b>	APPLE	M1	M1	10th	8 GB	DDR4	512 GB	0 GB
<b>622</b>	APPLE	M1	M1	10th	16 GB	DDR4	1024 GB	0 GB

```
In [7]: 1 df.drop_duplicates(inplace=True)
```

```
In [8]: 1 """Shape of dataset after dropping the duplicate rows"""  
2 df.shape
```

Out[8]: (802, 19)

```
In [9]: 1 """Checking for null values if any available"""  
2 df.isnull().sum()
```

```
Out[9]: brand                0  
processor_brand             0  
processor_name              0  
processor_gnrtn             0  
ram_gb                     0  
ram_type                   0  
ssd                        0  
hdd                        0  
os                         0  
os_bit                     0  
graphic_card_gb            0  
weight                     0  
warranty                   0  
Touchscreen                0  
msoffice                   0  
Price                      0  
rating                     0  
Number of Ratings          0  
Number of Reviews          0  
dtype: int64
```

```
In [10]: 1 """ This is showing no null values are present but while going through data  
2 this column contains 'Not Available' values instead of Null so will try to  
3  
4 df["processor_gnrtn"].replace("Not Available", None , inplace=True)
```

```
In [11]: 1 df.isnull().sum()
```

```
Out[11]: brand                0
processor_brand              0
processor_name               0
processor_gnrtn             224
ram_gb                      0
ram_type                    0
ssd                         0
hdd                         0
os                          0
os_bit                      0
graphic_card_gb             0
weight                      0
warranty                    0
Touchscreen                 0
msoffice                    0
Price                       0
rating                      0
Number of Ratings           0
Number of Reviews           0
dtype: int64
```

```
In [12]: 1 print(f"For {round((224/df.shape[0])*100,3)}% of data values are not present")
2         """Based on the analysis, processor_gnrtn column with a high percentage of missing values is dropped"""
3 df.drop("processor_gnrtn", axis=1, inplace=True)
```

For 27.93% of data values are not present

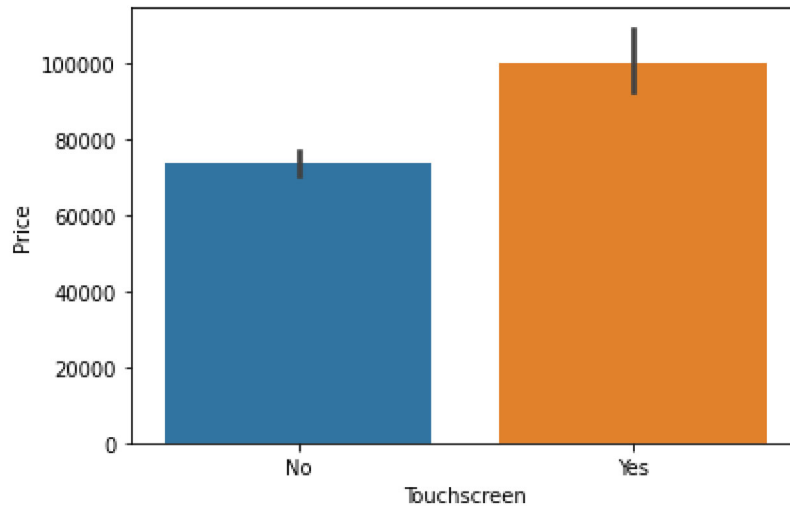
## Touch Screen

```
In [13]: 1 df["Touchscreen"].value_counts()
```

```
Out[13]: No      706
Yes        96
Name: Touchscreen, dtype: int64
```

```
In [14]: 1 sns.barplot(x=df['Touchscreen'],y=df['Price'])
```

```
Out[14]: <Axes: xlabel='Touchscreen', ylabel='Price'>
```



We can see from the plot, that the average price of laptop is higher with Touchscreen.

```
In [15]: 1 # converting this Touchscreen object coulumn in to int
2 df["Touchscreen"] = df["Touchscreen"].apply(lambda x: 1 if x=="Yes" else 0)
```

## Warranty

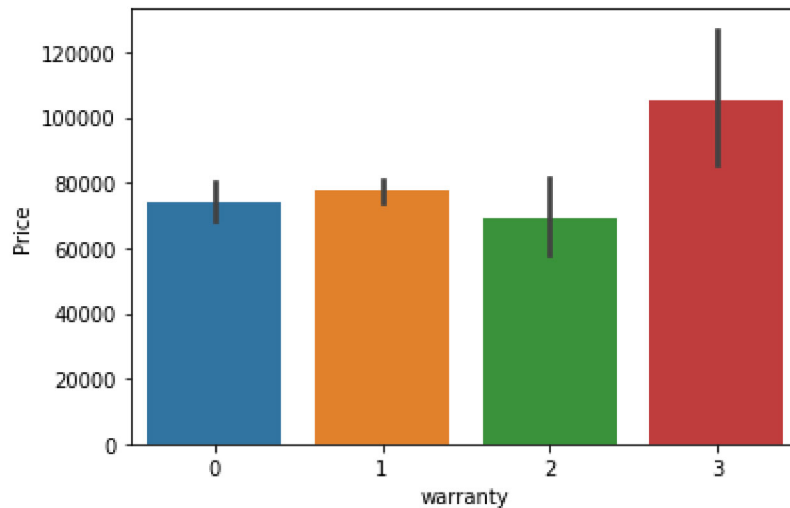
```
In [16]: 1 df["warranty"].value_counts()
```

```
Out[16]: 1 year          498
No warranty       268
2 years           23
3 years           13
Name: warranty, dtype: int64
```

```
In [17]: 1 """
2 We will transform this column into int
3 No warranty    0
4 1 year        1
5 2 years        2
6 3 years        3
7 """
8 def warranty_tranform(x):
9     if x== "No warranty": return 0
10    elif x == "1 year": return 1
11    elif x == "2 years": return 2
12    else: return 3
13
14 df["warranty"] = df["warranty"].apply(warranty_tranform)
```

```
In [18]: 1 """ We will check the affect of warranty duration on Laptop Price"""  
2 sns.barplot(x=df["warranty"],y=df['Price'])
```

Out[18]: <Axes: xlabel='warranty', ylabel='Price'>



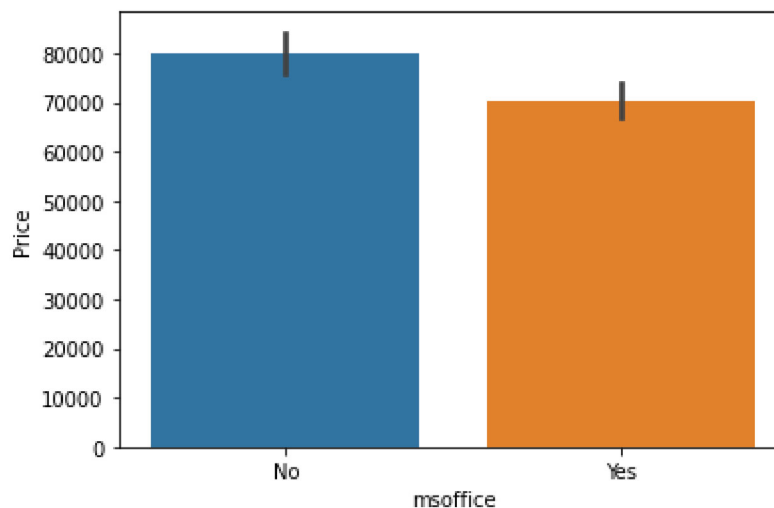
## MSOffice

```
In [19]: 1 df["msoffice"].value_counts()
```

Out[19]: No 522  
Yes 280  
Name: msoffice, dtype: int64

```
In [20]: 1 sns.barplot(x=df["msoffice"],y=df['Price'])
```

Out[20]: <Axes: xlabel='msoffice', ylabel='Price'>



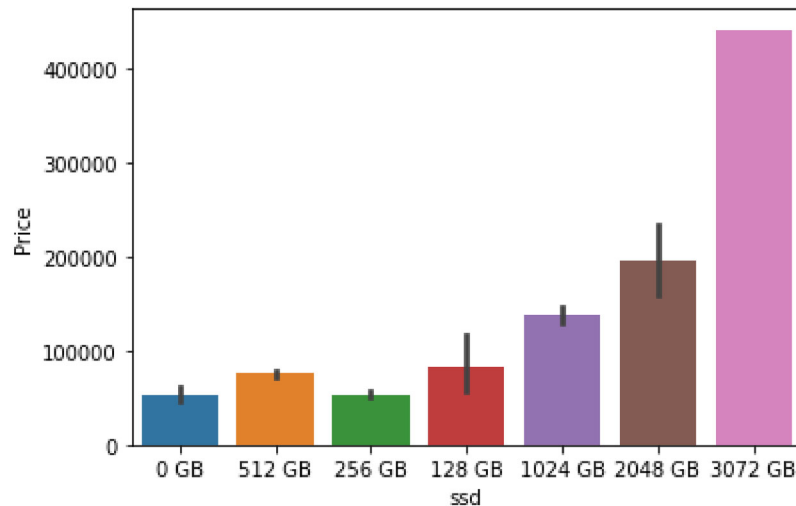
The bar plot indicates that laptops with MS Office have a lower average price compared to laptops without MS Office.

```
In [21]: 1 # converting this msoffice object coulumn in to int
          2 df["msoffice"] = df["msoffice"].apply(lambda x: 1 if x=="Yes" else 0)
```

## SSD

```
In [22]: 1 """ Comparing SSD with Price """
          2 sns.barplot(x=df["ssd"],y=df['Price'])
```

Out[22]: <Axes: xlabel='ssd', ylabel='Price'>



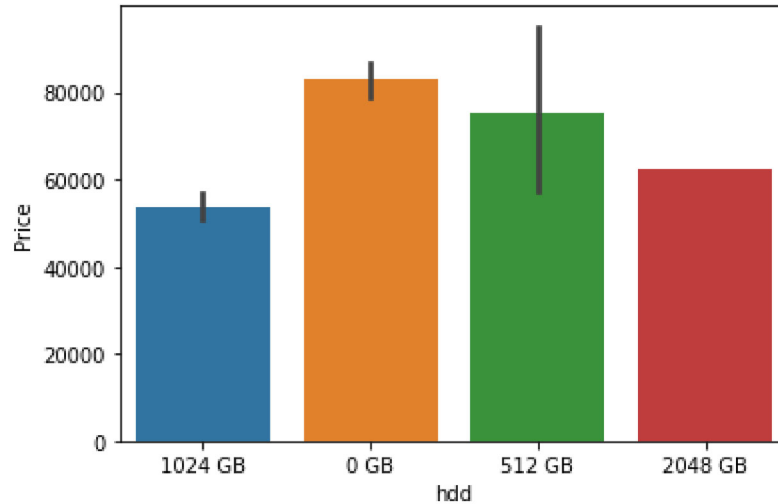
It is observed from the plot, with increase in SSD the laptop price increases.

```
In [23]: 1 """ Remove GB from SSD column values and convert it from object to int."""
          2 df["ssd"] = df["ssd"].str.replace("GB", "")
          3 df["ssd"] = df["ssd"].astype('int')
```

## HDD

```
In [24]: 1 """ Comparing HDD with Price """  
2 sns.barplot(x=df["hdd"],y=df['Price'])
```

Out[24]: <Axes: xlabel='hdd', ylabel='Price'>

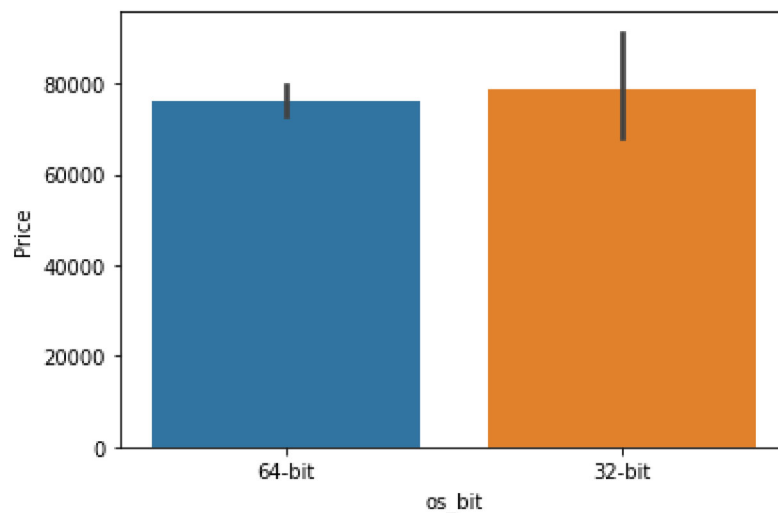


```
In [25]: 1 """ Remove GB from SSD column values and convert it from object to int."""  
2 df["hdd"] = df["hdd"].str.replace("GB","")  
3 df["hdd"] = df["hdd"].astype('int')
```

## OS-Bit

```
In [26]: 1 """ Comparing OS-Bit with Price """  
2 sns.barplot(x=df["os_bit"],y=df['Price'])
```

Out[26]: <Axes: xlabel='os\_bit', ylabel='Price'>



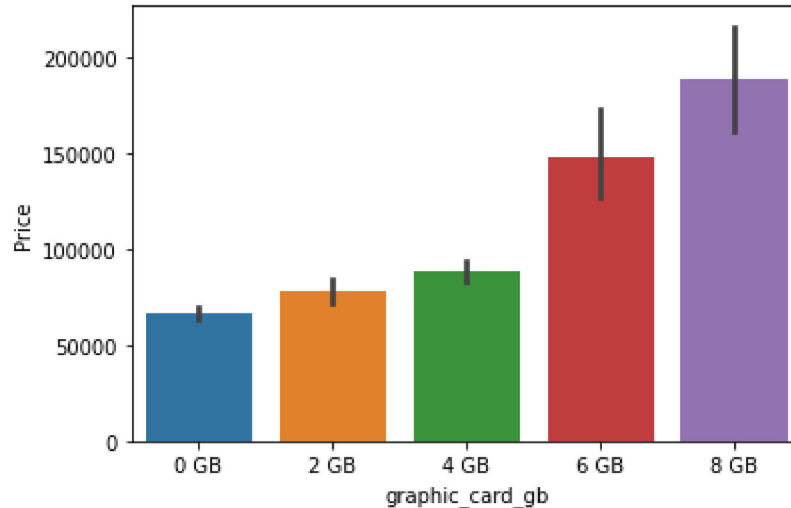
There is almost no effect of OS-Bit on price of Laptop



## Graphic\_Card\_GB

```
In [27]: 1 """ Comparing graphic_card_gb with Price """
        2 sns.barplot(x=df["graphic_card_gb"],y=df['Price'])
```

```
Out[27]: <Axes: xlabel='graphic_card_gb', ylabel='Price'>
```



It can be clearly seen that as graphic\_card\_gb size increases laptop price increases linearly

```
In [28]: 1 """ Remove GB from 'graphic_card_gb' column values and convert it from ob
        2 df["graphic_card_gb"] = df["graphic_card_gb"].str.replace("GB","")
        3 df["graphic_card_gb"] = df["graphic_card_gb"].astype('int')
```

```
In [29]: 1 df.corr()['Price']
```

```
Out[29]: ssd          0.628734
         hdd         -0.251266
         graphic_card_gb  0.467499
         warranty     0.057953
         Touchscreen   0.191227
         msoffice      -0.103783
         Price         1.000000
         Number of Ratings -0.152553
         Number of Reviews -0.156791
         Name: Price, dtype: float64
```

In [30]:

1 df

Out[30]:

	brand	processor_brand	processor_name	ram_gb	ram_type	ssd	hdd	os	os_bi
0	ASUS	Intel	Core i3	4 GB	DDR4	0	1024	Windows	64-bi
1	Lenovo	Intel	Core i3	4 GB	DDR4	0	1024	Windows	64-bi
2	Lenovo	Intel	Core i3	4 GB	DDR4	0	1024	Windows	64-bi
3	ASUS	Intel	Core i5	8 GB	DDR4	512	0	Windows	32-bi
4	ASUS	Intel	Celeron Dual	4 GB	DDR4	0	512	Windows	64-bi
...	...	...	...	...	...	...	...	...	...
818	ASUS	AMD	Ryzen 9	4 GB	DDR4	1024	0	Windows	64-bi
819	ASUS	AMD	Ryzen 9	4 GB	DDR4	1024	0	Windows	64-bi
820	ASUS	AMD	Ryzen 9	4 GB	DDR4	1024	0	Windows	64-bi
821	ASUS	AMD	Ryzen 9	4 GB	DDR4	1024	0	Windows	64-bi
822	Lenovo	AMD	Ryzen 5	8 GB	DDR4	512	0	DOS	64-bi

802 rows × 10 columns



## Rating

In [31]:

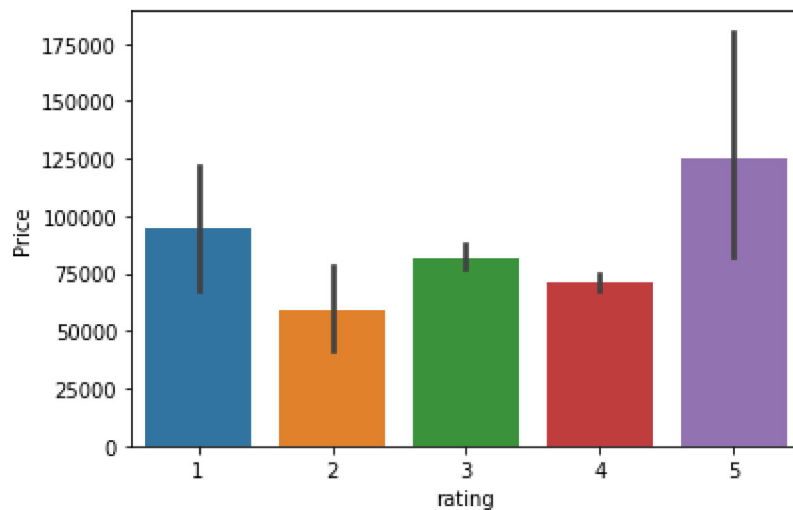
```

1 """ Rating values are in range of 1 to 5 i.e. worst to good. Remove stars
2 pattern = '|'.join(['stars', 'star'])
3 df['rating'] = df['rating'].str.replace(pattern, '')
4 df["rating"] = df["rating"].astype('int')

```

```
In [32]: 1 """ Comparing Rating with Price """  
2 sns.barplot(x=df["rating"],y=df['Price'])
```

Out[32]: <Axes: xlabel='rating', ylabel='Price'>



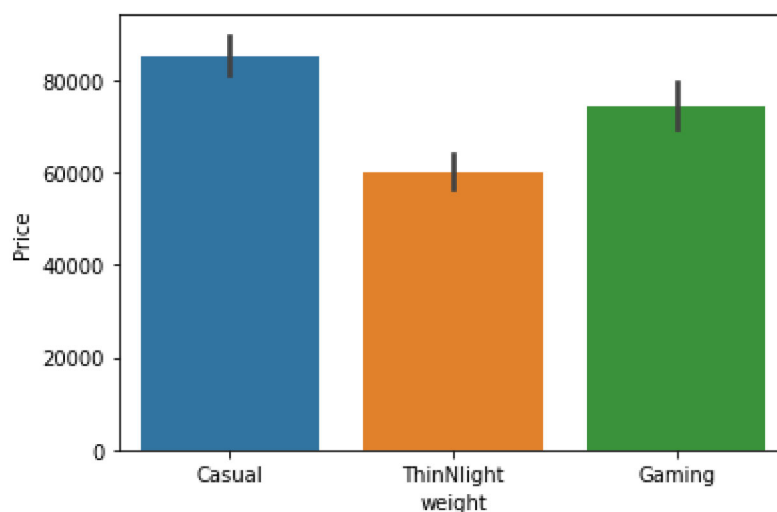
## Weight

```
In [33]: 1 df["weight"].value_counts()
```

Out[33]: Casual 509  
ThinNlight 254  
Gaming 39  
Name: weight, dtype: int64

```
In [34]: 1 """ Comparing weight with Price """  
2 sns.barplot(x=df["weight"],y=df['Price'])
```

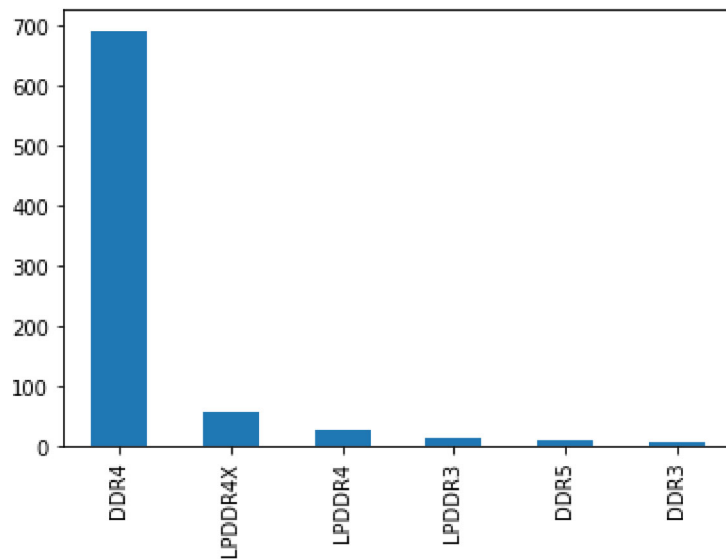
Out[34]: <Axes: xlabel='weight', ylabel='Price'>



## Ram\_type

```
In [35]: 1 df['ram_type'].value_counts().plot(kind='bar')
```

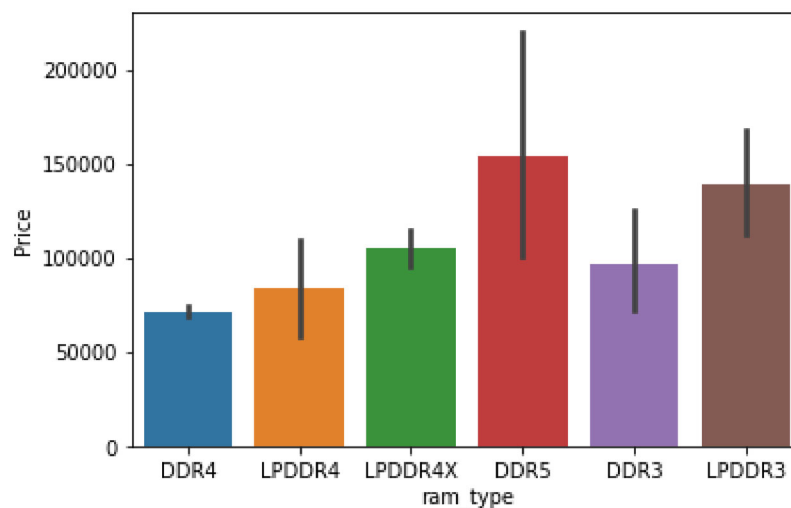
```
Out[35]: <Axes: >
```



observation: most of the laptops are of DDR4 type

```
In [36]: 1 """ Comparing ram_type with Price """  
2 sns.barplot(x=df["ram_type"],y=df['Price'])
```

```
Out[36]: <Axes: xlabel='ram_type', ylabel='Price'>
```



Based on the observed plot, it is evident that the "ram\_type" column has an impact on the price of laptops. Considering the nature of this categorical column as ordinal, we can handle it accordingly.

```
In [37]: 1 # Define the mapping for ordinal encoding
2 mapping = {
3     'DDR4': 1,
4     'LPDDR4X': 4,
5     'LPDDR4': 2,
6     'LPDDR3': 5,
7     'DDR5': 6,
8     'DDR3': 3,
9 }
10 # Perform ordinal encoding
11 df['ram_type'] = df['ram_type'].map(mapping)
12
```

```
In [38]: 1 df.corr()["Price"]
```

```
Out[38]: ram_type      0.316037
ssd      0.628734
hdd     -0.251266
graphic_card_gb  0.467499
warranty  0.057953
Touchscreen  0.191227
msoffice  -0.103783
Price      1.000000
rating    -0.040564
Number of Ratings -0.152553
Number of Reviews -0.156791
Name: Price, dtype: float64
```

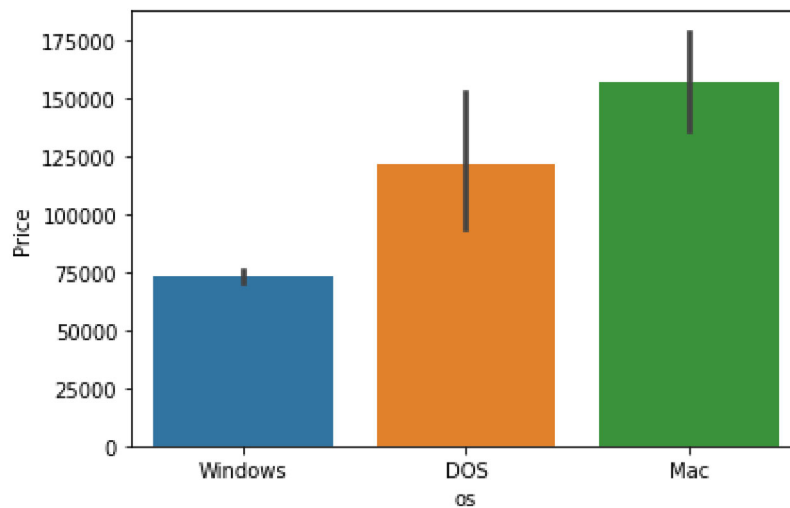
## OS

```
In [39]: 1 df["os"].value_counts()
```

```
Out[39]: Windows      763
Mac          23
DOS          16
Name: os, dtype: int64
```

```
In [40]: 1 """ Comparing os with Price """
        2 sns.barplot(x=df["os"],y=df['Price'])
```

Out[40]: <Axes: xlabel='os', ylabel='Price'>



```
In [ ]: 1
```

```
In [41]: 1 df
```

Out[41]:

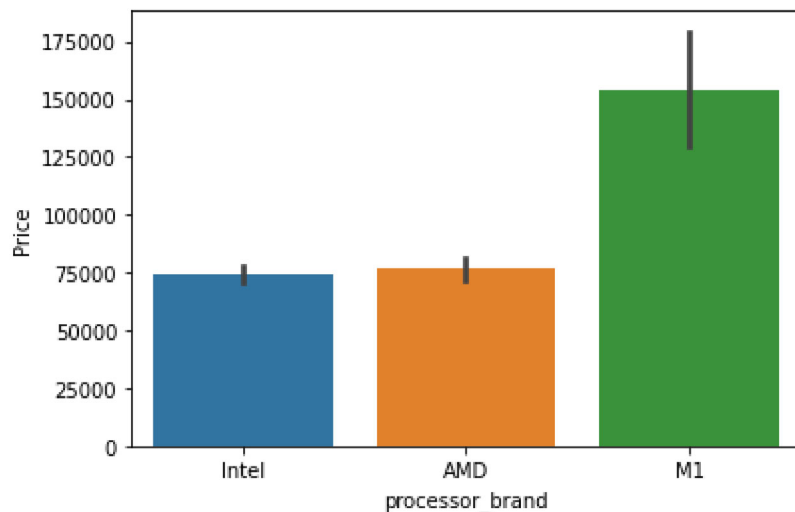
	brand	processor_brand	processor_name	ram_gb	ram_type	ssd	hdd	os	os_bi
0	ASUS	Intel	Core i3	4 GB	1	0	1024	Windows	64-bi
1	Lenovo	Intel	Core i3	4 GB	1	0	1024	Windows	64-bi
2	Lenovo	Intel	Core i3	4 GB	1	0	1024	Windows	64-bi
3	ASUS	Intel	Core i5	8 GB	1	512	0	Windows	32-bi
4	ASUS	Intel	Celeron Dual	4 GB	1	0	512	Windows	64-bi
...	...	...	...	...	...	...	...	...	...
818	ASUS	AMD	Ryzen 9	4 GB	1	1024	0	Windows	64-bi
819	ASUS	AMD	Ryzen 9	4 GB	1	1024	0	Windows	64-bi
820	ASUS	AMD	Ryzen 9	4 GB	1	1024	0	Windows	64-bi
821	ASUS	AMD	Ryzen 9	4 GB	1	1024	0	Windows	64-bi
822	Lenovo	AMD	Ryzen 5	8 GB	1	512	0	DOS	64-bi

802 rows × 10 columns



```
In [42]: 1 """ Comparing processor_brand with Price """  
2 sns.barplot(x=df["processor_brand"],y=df['Price'])
```

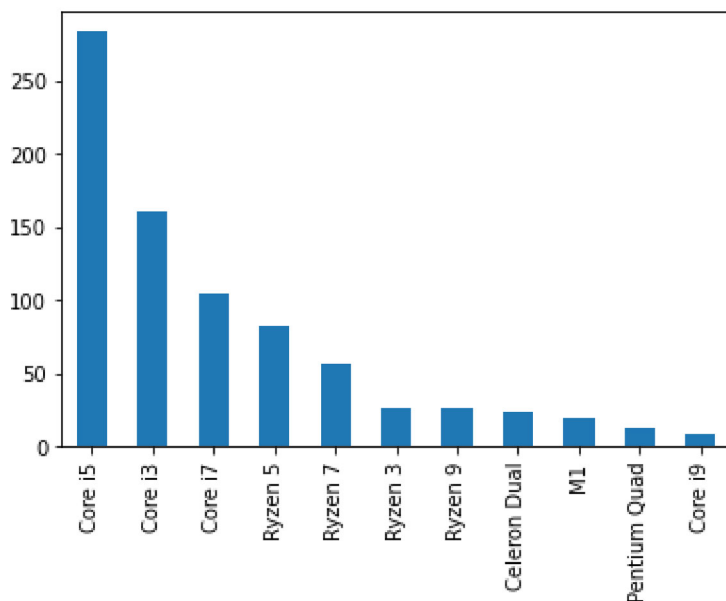
Out[42]: <Axes: xlabel='processor\_brand', ylabel='Price'>



## Processor Name

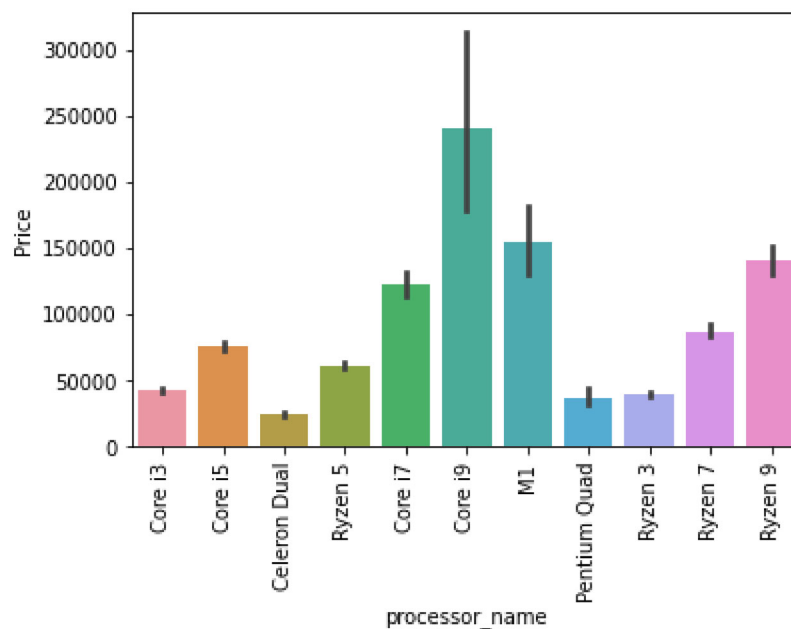
```
In [43]: 1 """Laptop data distribution based on processor_name"""  
2 df['processor_name'].value_counts().plot(kind='bar')
```

Out[43]: <Axes: >



```
In [44]: 1 """ Comparing processor_name with Price """
2 sns.barplot(x=df["processor_name"],y=df['Price'])
3 plt.xticks(rotation = 90)
```

```
Out[44]: (array([ 0,  1,  2,  3,  4,  5,  6,  7,  8,  9, 10]),
 [Text(0, 0, 'Core i3'),
  Text(1, 0, 'Core i5'),
  Text(2, 0, 'Celeron Dual'),
  Text(3, 0, 'Ryzen 5'),
  Text(4, 0, 'Core i7'),
  Text(5, 0, 'Core i9'),
  Text(6, 0, 'M1'),
  Text(7, 0, 'Pentium Quad'),
  Text(8, 0, 'Ryzen 3'),
  Text(9, 0, 'Ryzen 7'),
  Text(10, 0, 'Ryzen 9')])
```



```
In [45]: 1 df["processor_name"].value_counts()
```

```
Out[45]: Core i5      284
Core i3      161
Core i7      104
Ryzen 5       82
Ryzen 7       56
Ryzen 3       26
Ryzen 9       26
Celeron Dual  23
M1            19
Pentium Quad  13
Core i9        8
Name: processor_name, dtype: int64
```

The bar plot clearly indicates that the "Processor Brand" feature is an ordinal categorical variable, as the version of the processor directly impacts its price.



In [46]:

```

1  """
2  The processor_name feature column will be transformed into an ordinal cate
3  """
4  def processor_name_transform(x):
5      if x == "Celeron Dual":
6          return 1
7      elif x == "Pentium Quad" or x == "Ryzen 3":
8          return 2
9      elif x == "Core i3":
10         return 3
11         elif x == "Ryzen 5":
12             return 4
13         elif x == "Core i5":
14             return 5
15         elif x == "Ryzen 7":
16             return 6
17         elif x == "Core i7":
18             return 7
19         elif x == "Ryzen 9":
20             return 8
21         elif x == "M1":
22             return 9
23         elif x == "Core i9":
24             return 10
25
26
27  df["processor_name"] = df["processor_name"].apply(processor_name_transform)

```

## Ram

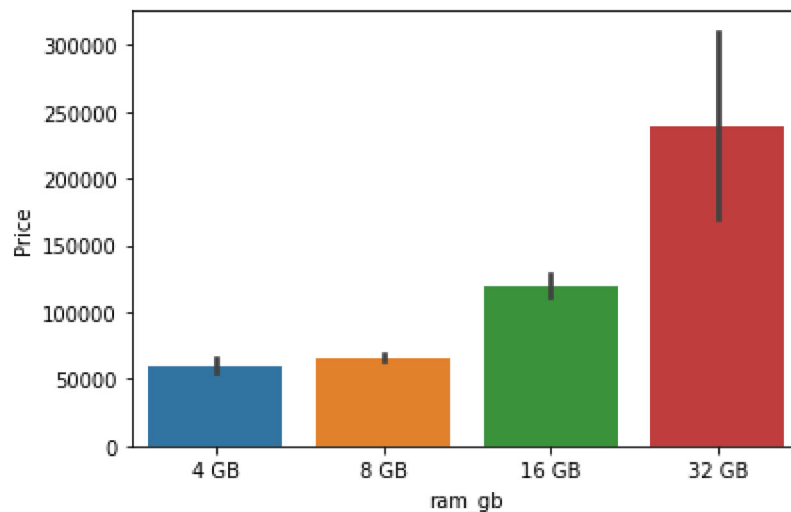
In [47]:

```

1  """ Comparing ram_gb with Price """
2  sns.barplot(x=df["ram_gb"],y=df['Price'])

```

Out[47]: &lt;Axes: xlabel='ram\_gb', ylabel='Price'&gt;



Here also, the bar plot clearly indicates that the "Ram Size" feature is an ordinal categorical variable, as the size of the ram directly impacts its price.

```
In [48]: 1 """ Remove GB from ram_gb column values and convert it from object to int
2 df["ram_gb"] = df["ram_gb"].str.replace("GB", "")
3 df["ram_gb"] = df["ram_gb"].astype('int')
```

```
In [49]: 1 df.corr()["Price"]
```

```
Out[49]: processor_name    0.745720
ram_gb    0.516454
ram_type  0.316037
ssd       0.628734
hdd      -0.251266
graphic_card_gb  0.467499
warranty    0.057953
Touchscreen 0.191227
msoffice   -0.103783
Price      1.000000
rating     -0.040564
Number of Ratings -0.152553
Number of Reviews -0.156791
Name: Price, dtype: float64
```

```
In [50]: 1 df
```

```
Out[50]:
```

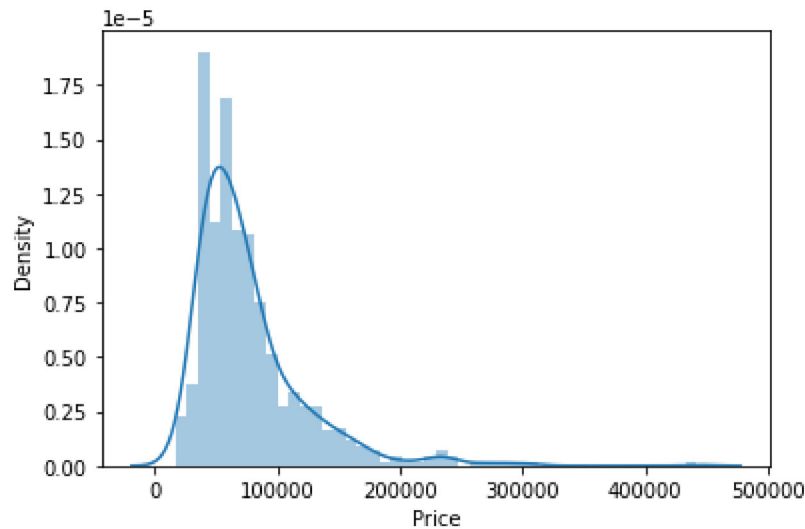
	brand	processor_brand	processor_name	ram_gb	ram_type	ssd	hdd	os	os_bi
0	ASUS	Intel	3	4	1	0	1024	Windows	64-bi
1	Lenovo	Intel	3	4	1	0	1024	Windows	64-bi
2	Lenovo	Intel	3	4	1	0	1024	Windows	64-bi
3	ASUS	Intel	5	8	1	512	0	Windows	32-bi
4	ASUS	Intel	1	4	1	0	512	Windows	64-bi
...	...	...	...	...	...	...	...	...	..
818	ASUS	AMD	8	4	1	1024	0	Windows	64-bi
819	ASUS	AMD	8	4	1	1024	0	Windows	64-bi
820	ASUS	AMD	8	4	1	1024	0	Windows	64-bi
821	ASUS	AMD	8	4	1	1024	0	Windows	64-bi
822	Lenovo	AMD	4	8	1	512	0	DOS	64-bi

802 rows × 10 columns



```
In [51]: 1 """Checking the distribution of price column"""  
2 sns.distplot(df['Price'])
```

Out[51]: <Axes: xlabel='Price', ylabel='Density'>



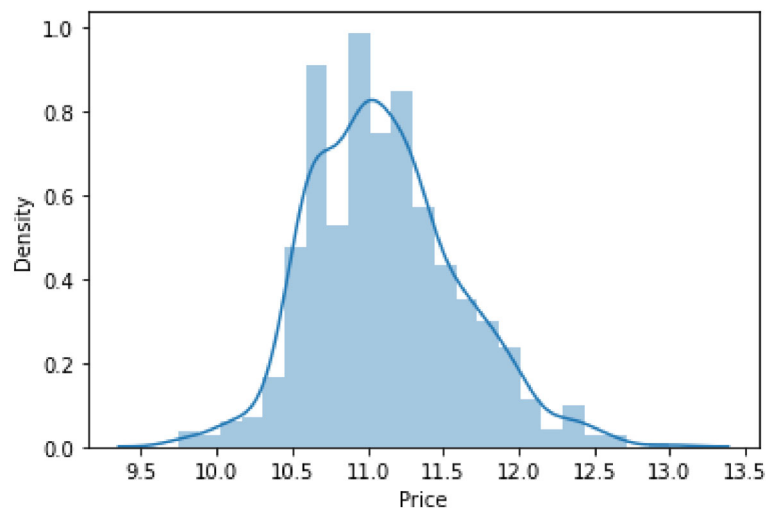
The distribution plot reveals a right-skewed pattern, indicating a higher concentration of data towards the lower values.

Applying a log transformation to a right-skewed distribution can help in achieving a more symmetrical or normally distributed shape, which is often preferred in statistical modeling and analysis.

This transformation can be beneficial because many statistical techniques assume normality, allowing for more accurate and reliable analysis results.

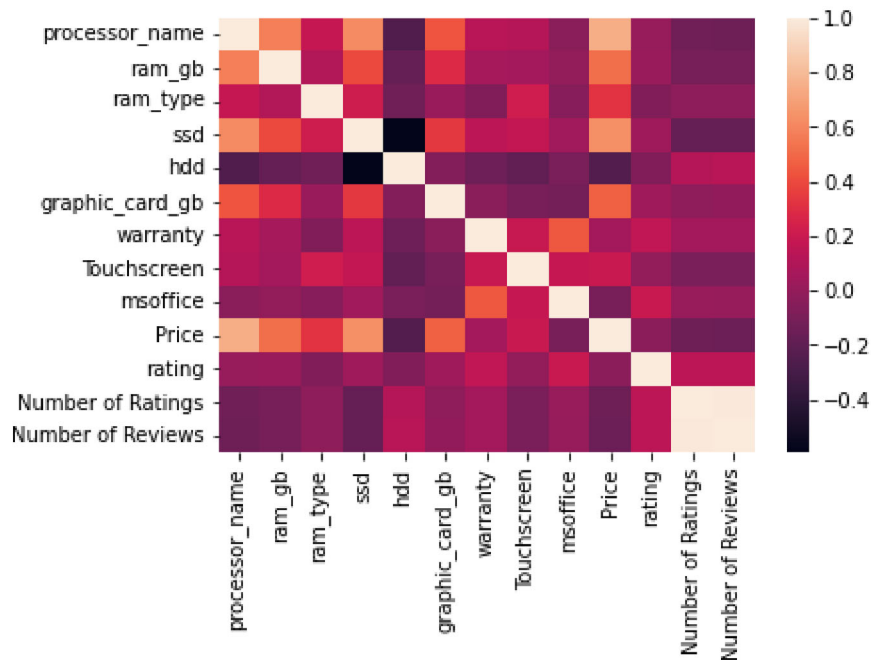
```
In [52]: 1 sns.distplot(np.log(df['Price']))
```

Out[52]: <Axes: xlabel='Price', ylabel='Density'>



```
In [53]: 1 """Heatmap Analysis"""
          2 sns.heatmap(df.corr())
```

Out[53]: <Axes: >



The varying intensities of the heatmap highlight the strength and direction of these correlations, aiding in identifying key factors that influence the pricing of laptops.

The features having high correlations with the price of laptop are:

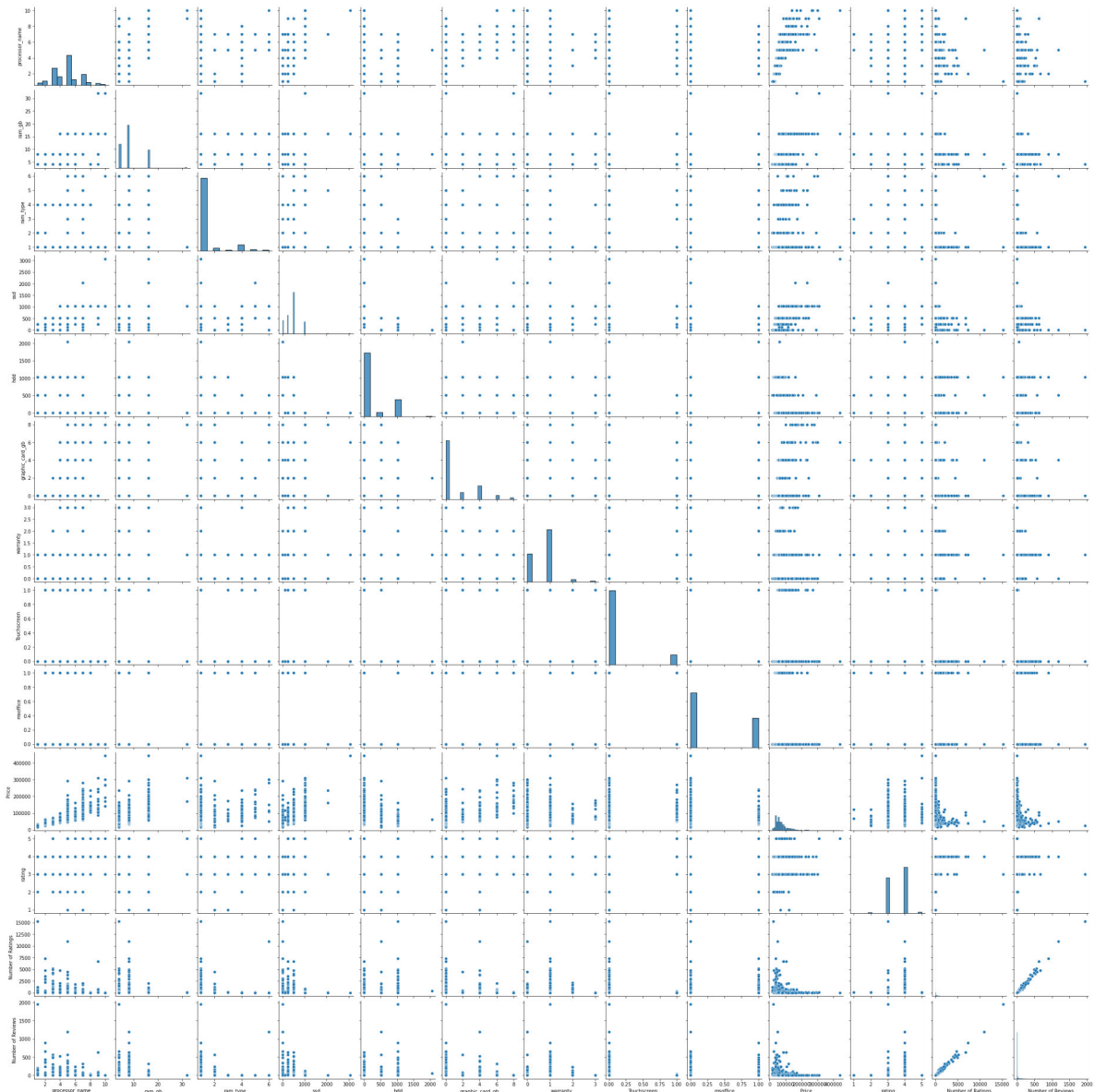
1. processor\_name
2. ram\_gb
3. ssd
4. ram\_type
5. graphic\_card\_gb

```
In [54]: 1 df.corr()["Price"]
```

```
Out[54]: processor_name    0.745720
ram_gb      0.516454
ram_type    0.316037
ssd         0.628734
hdd        -0.251266
graphic_card_gb 0.467499
warranty     0.057953
Touchscreen  0.191227
msoffice    -0.103783
Price       1.000000
rating      -0.040564
Number of Ratings -0.152553
Number of Reviews -0.156791
Name: Price, dtype: float64
```

```
In [66]: 1 sns.pairplot(df)
```

```
Out[66]: <seaborn.axisgrid.PairGrid at 0x19c46942e20>
```



## D. Input Variables and Target separation and Splitting into train and test data

```
In [55]: 1 """ Separating the input variables and Target into X and y repectively"""
2 X = df.drop(columns=['Price'])
3 y = np.log(df['Price'])
```

```
In [56]: 1 """Splitting our dataset into Train and Validation"""
2 X_train,X_test,y_train,y_test = train_test_split(X,y,test_size=0.15,random
```

## E.Applying Machine Learning Model for Laptop Price

## E1. Linear Regression Model

```
In [57]: 1 from sklearn.compose import ColumnTransformer
2 from sklearn.preprocessing import OneHotEncoder
3 from sklearn.metrics import r2_score, mean_absolute_error
4 from sklearn.linear_model import LinearRegression
5 from sklearn.neighbors import KNeighborsRegressor
6 from sklearn.tree import DecisionTreeRegressor
7 from sklearn.model_selection import RandomizedSearchCV, GridSearchCV
8 from sklearn.ensemble import GradientBoostingRegressor
9 from scipy.stats import randint
```

```
In [58]: 1 """Converting Nominal Categorical Fetures into One Hot Encodings"""
2
3 col_tnf = ColumnTransformer(transformers=[
4     ('col_tnf', OneHotEncoder(sparse=False, drop='first'), [0,1,7,8,10])
5 ], remainder='passthrough')
6
7 X_train_transformed = col_tnf.fit_transform(X_train)
8 X_test_transformed = col_tnf.transform(X_test)
```

```
In [59]: 1 LinearRegressionModel = LinearRegression()
2
3 # fit and prediction
4 LinearRegressionModel.fit(X_train_transformed, y_train)
5 y_pred = LinearRegressionModel.predict(X_test_transformed)
6
7 # Calculate R2 score and MAE
8 print('R2 score:', r2_score(y_test, y_pred))
9 print('MAE:', mean_absolute_error(y_test, y_pred))
```

R2 score: 0.8492439253870319

MAE: 0.1459487523971702

## E2. KNN Regression Model

```
In [60]: 1 KNNModel = KNeighborsRegressor()
2 # Create a parameter grid for randomized search
3 param_grid = {
4     'n_neighbors': np.arange(1, 20),
5 }
6 # Randomized Search CV
7 rs_cv = RandomizedSearchCV(
8     estimator=KNNModel,
9     param_distributions=param_grid,
10    scoring='r2',
11    n_iter=10, # Number of parameter settings that are sampled
12    cv=5,      # Number of folds in cross-validation
13    random_state=42
14 )
15
16 # fit and prediction
17 rs_cv.fit(X_train_transformed, y_train)
18 y_pred = rs_cv.predict(X_test_transformed)
19
20 # Calculate R2 score and MAE
21 print('R2 score:', r2_score(y_test, y_pred))
22 print('MAE:', mean_absolute_error(y_test, y_pred))
```

R2 score: 0.5743328756286931

MAE: 0.2554921623082943

## E3.Decision Tree Regression Model

```
In [61]: 1 dtr = DecisionTreeRegressor()
2 # Create a parameter grid for randomized search
3 param_grid = {
4     'max_depth': np.arange(1, 10),
5     'min_samples_split': np.arange(2, 11),
6     'min_samples_leaf': np.arange(1, 11)
7 }
8 # Randomized Search CV
9 rs_cv = RandomizedSearchCV(
10     estimator=dtr,
11     param_distributions=param_grid,
12     scoring='r2',
13     n_iter=10, # Number of parameter settings that are sampled
14     cv=5,      # Number of folds in cross-validation
15     random_state=42
16 )
17 rs_cv.fit(X_train_transformed, y_train)
18 best_params = rs_cv.best_params_
19 dtr_best = DecisionTreeRegressor(**best_params)
20
21 # fit and prediction
22 dtr_best.fit(X_train_transformed, y_train)
23 y_pred = dtr_best.predict(X_test_transformed)
24
25 # Calculate R2 score and MAE
26 print('R2 score:', r2_score(y_test, y_pred))
27 print('MAE:', mean_absolute_error(y_test, y_pred))
28 print('Best parameters:', best_params)
29
```

R2 score: 0.8573330917319021

MAE: 0.1454562949028369

Best parameters: {'min\_samples\_split': 4, 'min\_samples\_leaf': 4, 'max\_depth': 7}



## E4. Gradient Boosting Model

```
In [62]: 1 param_grid = {
2         'n_estimators': randint(100, 1000),
3         'learning_rate': [0.1, 0.01, 0.001],
4         'max_depth': randint(3, 6)
5     }
6
7 GradientBoostingRegressorModel = GradientBoostingRegressor()
8
9 rs_cv = RandomizedSearchCV(GradientBoostingRegressorModel, param_distrib
10 rs_cv.fit(X_train_transformed, y_train)
11 best_params = rs_cv.best_params_
12 best_model = GradientBoostingRegressor(**best_params)
13
14 # fit and prediction
15 best_model.fit(X_train_transformed, y_train)
16 y_pred = best_model.predict(X_test_transformed)
17
18 # Calculate R2 score and MAE
19 print('R2 score:', r2_score(y_test, y_pred))
20 print('MAE:', mean_absolute_error(y_test, y_pred))
21 print('Best parameters:', best_params)
22
```

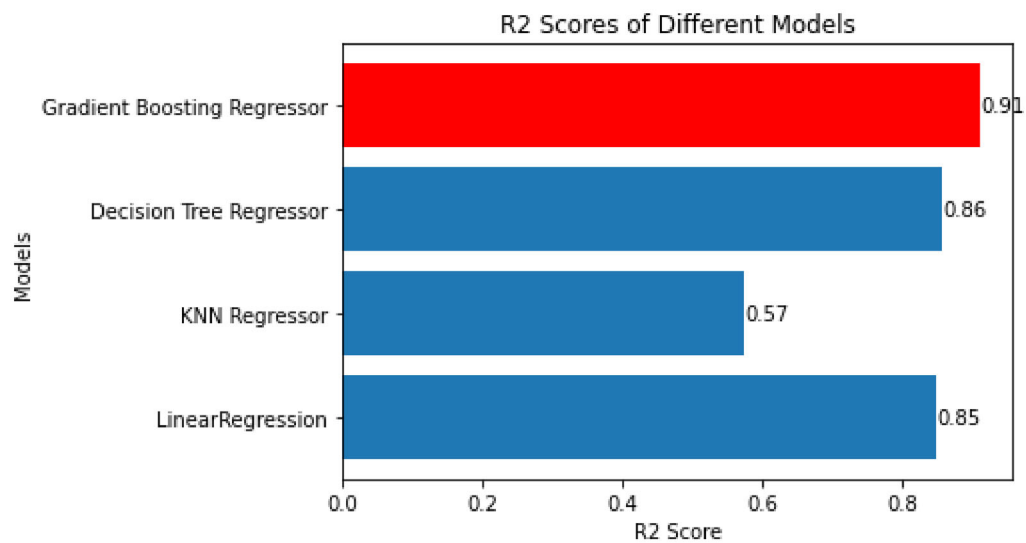
R2 score: 0.9107112564822367

MAE: 0.11991196744780869

Best parameters: {'learning\_rate': 0.01, 'max\_depth': 4, 'n\_estimators': 869}

## F. Plots of Different tried models results for analysis.

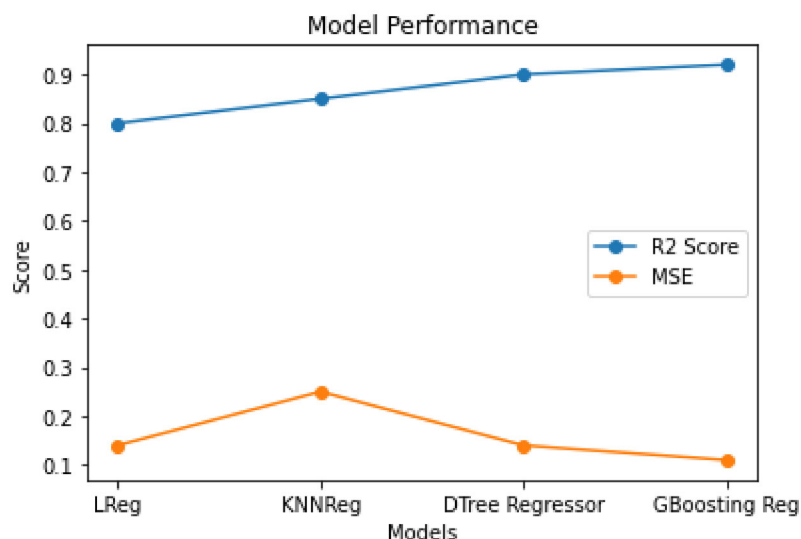
```
In [63]: 1 r2_scores = [0.8492, 0.5743, 0.8573, 0.9116]
2 model_names = ['LinearRegression', 'KNN Regressor', 'Decision Tree Regressor', 'Gradient Boosting Regressor']
3 highest_index = np.argmax(r2_scores)
4 plt.barh(model_names, r2_scores)
5 for i, score in enumerate(r2_scores):
6     plt.text(score, i, f'{score:.2f}', ha='left', va='center')
7 plt.barh(highest_index, r2_scores[highest_index], color='red')
8
9 plt.xlabel('R2 Score')
10 plt.ylabel('Models')
11 plt.title('R2 Scores of Different Models')
12 plt.show()
13
```



**The R2 score measures the goodness of fit of a regression model. It ranges from 0 to 1, A score of 1 indicates that the model perfectly predicts the observed data.**

**As we can see from the plot Gradient Boosting Regressor model is performing best amongst the 4 models with highest R2 : 0.91**

```
In [64]: 1 import matplotlib.pyplot as plt
2 r2_scores = [0.8, 0.85, 0.9, 0.92]
3 mse_values = [0.14, 0.25, 0.14, 0.11]
4 model_names = ['LReg', 'KNNReg', 'DTree Regressor', 'GBoosting Reg']
5 plt.plot(model_names, r2_scores, marker='o', label='R2 Score')
6 plt.plot(model_names, mse_values, marker='o', label='MSE')
7
8 plt.xlabel('Models')
9 plt.ylabel('Score')
10 plt.title('Model Performance')
11 plt.legend()
12 plt.show()
13
```



```
In [65]: 1 # G. Inference
```

We have selected Linear Regression (LR) as our first model because of its simplicity and interpretability, it is giving an  $r^2$  around 0.8. LR models performance is limited when dealing with complex relationships or non-linearities in the data.

To get more better results we have opted KNN Regressor as our second model, as it can capture non-linear patterns by considering the proximity of data points. Despite attempting various values for the `n_neighbors` parameter in the range of "1-20" for the KNN Regressor, we were still unable to achieve satisfactory results. It's possible that it might be struggling with high-dimensional data or large datasets due to computational limitations.

As decision trees have the ability to capture non-linear relationships, handle complex interactions between variables, be robust to outliers, and handle mixed data types effectively. We selected it as our next model. After performing Randomized Search CV on the Decision Tree Regressor (DTR) model, the best parameter combination that we got the best results was:

`min_samples_split: 4 min_samples_leaf: 4 max_depth: 7`

Type *Markdown* and LaTeX:  $\alpha^2$