

# LAPTOP PRICE PREDICTOR USING ML

Arav Khandal

Devendra Chouhan

Anjan Das

(20103028)

(20103049)

(20103017)

(Department of Computer Science)

(Department of Computer Science)

Dr. B R Ambedkar national

Dr. B R Ambedkar national

Institute of technology

Institute of technology

Jalandhar, Punjab

Jalandhar, Punjab

Institute of technology

Jalandhar, Punjab

Aravk.cs.20@nitj.ac.in

Anjand.cs.20@nitj.ac.in

Devendra.cs.20@nitj.ac.in

IT industry is in the expansion phase, more and more people are continuously shifting towards IT sector. In this high-tech era, laptop has turned into a necessity. Gone are the days when laptop was considered a luxury item.

However even though people know the specification and the purpose of purchasing a laptop but budget has always been a contributing factor. People tend to buy the best item in the lowest price possible of course with due consideration to the purpose of purchase.

This is why we decided to make a project of predicting the price of laptop based on user's input of its need and specification.

## I. INTRODUCTION AND MOTIVATION

In our daily lives, the laptop has become one of the most important and frequently used items. However, I'm sure you've had to find a personal laptop that meets your requirements? With so many specifications and brand names on the market, it's tough for laptop manufacturers to sell their products and for users to choose which laptop to buy.

Why should we be concerned about this issue? In 2019, it is estimated that 166 million laptops were sold and shipped (International Data Corporation). The average selling price of a laptop is approximately \$629 USD. According to Grand View Research, the global laptop market is estimated to be worth over \$100 billion, with a forecasted market value of 108.9 billion by 2025. Laptops will continue to be in high demand due to rising consumer buying power and increased desire for technologically advanced products. People and manufacturers alike will need to be aware and competitive as laptops become more popular around the world.

## II. Project Motivation:-

1. Determine which of the three kinds of laptop contributes the most in terms of quantity and quality.
2. Research laptop trends across a variety of brands, models, and performance levels.
3. Estimate laptop pricing based on quantity and quality factors such as memory, storage size and kind, CPU and GPU performance, brand, type, and screen size.
4. The KNN Recommendation System recommends a laptop based on its specifications.

**III. AIM :-** Through the proposed project, we aim to create a working machine learning model that streamlines the process of assessing the price of laptops and predict their price .

Machine learning is a very powerful tool when it comes to predicting and analyzing large amounts of data, which can be used to create AI powered systems.

## IV. TERMINOLOGY

1. **Company:-** laptop manufacturing companies manufacture many types and many cost of laptop, Apple, Dell, HP, Acer, Asus, Lenovo etc.
2. **Screen Size:-** Standard laptop models usually come with a built-in **13-inch to 15-inch** wide-

screen LCD with a  $1280 \times 800$  dot or  $1366 \times 768$  dot resolution. This level of resolution is good enough for basic Windows operation. Standard laptop models often include a built-in 13-inch to 15-inch wide-screen LCD with a resolution of 1280 800 dots or 1366 768 dots. Although this resolution is adequate for basic Windows activities, a large, high-resolution computer screen is unquestionably more user-friendly, but it goes without saying that a large, high-resolution computer screen is much more user-friendly.

3. **Screen Resolution:-** Horizontal and vertical pixel counts are used to indicate display resolution. 1366-by-768 (also known as HD) and 1920-by-1080 (sometimes known as Full HD) are the most often utilised resolutions on laptops and 2-in-1 PCs currently (Full HD or 1080p). In our opinion, the most ideal screen resolution for laptops is 1920-by-1080.
4. **CPU(Central Processing Unit):-** All kinds of computing devices such as tablets, PCs, or laptops feature a brain-like unit called the central processing unit or CPU. Your computer's CPU calculates and interprets instructions while you're surfing the web, creating documents, playing games, or running software programs. It's a critical component that your PC can't function without.

We'll look at why a CPU is so crucial and how to keep track of your CPU utilisation to ensure optimal system performance in the sections below.

5. **RAM(Random Access Memory):-** Random-access memory (RAM) stands for random-access memory, but what exactly does that imply? The RAM in your computer is simply short-term memory that stores data as the processor requires it. This should not be mistaken with long-term data stored on your hard disc, which remains accessible even after your computer has been switched off.
6. **GPU(Graphics processing unit):-** A customised processor that was created to speed up the rendering of graphics. GPUs can process a large amount of data at once, making them ideal for machine learning, video editing, and gaming.
7. **Memory:-** Memory is the part of your computer that allows you to store and retrieve data over an extended period of time. Storage is usually in the form of a solid-state drive (SSD) or a hard disc (HDD).

8. **Operating System:-** The operating system (OS) is the software that controls the software, hardware, and resources of a computer. The operating system is required to organise shared functions and to offer a user interface for interacting with software and hardware.

9. **Weight of Laptop:-** Laptops can be as light as 2 pounds or as heavy as 8 pounds. They are usually classified into one of five groups: Ultrabook /Chromebook (2–3 pounds), Ultraportable (2–5 pounds), Thin and Light (3–6 pounds), Desktop Replacement (4 pounds), and Luggables (8 pounds) are the different types of laptops (8 lbs).

## V. PREVIOUS WORK:-

E-commerce websites use a similar type of ML model to sort products based on the user's requirements and availability of products on their website, as well as generate results based on the user's previous search history on the site, providing better results for the user by utilising decision tree algorithms and many other algorithm combinations.

## VI. PROPOSED SOLUTION:-

We have developed a model that predicts the price of the laptop based on the requirement of the user i.e ram , processor , graphics card , screen size , resolution audio quality

The feature act as the independent variable for the model that decides the price of the laptop mainly (price may vary depending on the brand )

### A. Algorithms used:-

- **ColumnTransformer:-** ColumnTransformer is a scikit-learn Python machine learning library class that allows you to apply data preparation transformations selectively.

It allows you to apply a certain transform or series of transforms to only the numerical columns, and another sequence of transforms to only the category columns, for example.

You must give a list of transformers to utilise the ColumnTransformer.

Each transformer is a three-element tuple that specifies the transformer's name, the transform to be applied, and the column indices to which it should be applied.

- **Pipeline:-** The pipeline is a scikit-learn tool for coordinating machine learning activities written in Python.

Pipelines work by connecting a linear succession of data transformations, resulting in a quantifiable modelling process.

The goal is to ensure that all stages of the pipeline, such as training datasets and each fold in the cross-validation process, are confined to the data available for the evaluation.

- **OneHotEncoder:-** On categorical data, most present machine learning methods are ineffective. Rather, categorical data must be translated to numerical data first. One of the strategies used to make this conversion is one-hot encoding. This strategy is typically utilised when applying deep learning techniques to sequential classification tasks.

The representation of categorical variables as binary vectors is known as one-hot encoding. First, the categorical values are converted to integers. Each integer value is then represented as

an all-zero binary vector (except the index of the integer which is marked as 1).

- **LinearRegression:-** The first machine learning method that any data scientist encounters is linear regression. It's a simple model, but it's one that everyone should learn because it's the cornerstone for all other machine learning algorithms.

It is a very effective strategy for determining the elements that determine profitability. It may be used to forecast sales in the following months by looking at previous month's sales data. It may also be used to learn more about how customers behave. By the conclusion of the blog, we will have created a model that looks like the one below, i.e., we will have determined the optimum line to suit the data.

This is the first post in a series on machine learning that I'll be covering. The sheer volume of material regarding machine learning algorithms on the internet might be overwhelming. This site has a dual role for me. It can serve as a starting point for anyone interested in machine learning, as well as a reference for myself.

- **KNeighborsRegressor:-** The K nearest neighbours method stores all available examples and predicts the numerical target using a similarity metric (e.g., distance functions). KNN has been utilised as a non-parametric approach in statistical estimates and pattern recognition since the early 1970s. Algorithm Calculating the average of the numerical goal of the K nearest neighbours is a straightforward implementation of KNN regression. An inverse distance weighted average of the K closest neighbours is another method. The same distance functions are used in KNN regression as in KNN classification.

Only continuous variables are eligible for the following three distance measurements. In the case of categorical variables, the Hamming distance must be used, which is a measure of the number of times similar symbols in two strings of equal length differ.

A single neighbour forecast is just the target value of the nearest neighbour.

- **DecisionTreeRegressor:-** A decision tree is a model of decisions and all of their possible consequences, including outcomes, input costs, and utility, that employs a flowchart-like tree structure.

The supervised learning methods include the decision-tree algorithm. It may be used with both continuous and categorical output variables.

Decision tree regression examines an object's characteristics and trains a model in the shape of a tree to forecast future data and create meaningful continuous output. The output/result is not discrete, in the sense that it is not represented solely by a discrete, known set of numbers or values.

A weather prediction model that forecasts whether or not it will rain on a specific day is an example of a discrete output.

A profit prediction model that specifies the likely profit that may be earned from the sale of a product is an example of continuous output.

- **RandomForestRegressor:-** Every decision tree has a significant variance, but when we mix all of them in parallel, the final variance is minimal since each decision tree is completely trained on that specific sample data, and so the outcome is dependent on numerous decision trees rather than one. In the event of a classification problem, the majority voting classifier is used to determine the final result. The ultimate result in a regression problem is the average of all the outputs. Aggregation is the name for this section.

Random Forest is an ensemble approach that uses several decision trees with a technique called Bootstrap and Aggregation, sometimes known as bagging, to solve both regression and classification problems. Instead than depending on individual decision trees, the main idea is to aggregate numerous decision trees to determine the final outcome.

As a fundamental learning model, Random Forest uses several decision trees. Row and feature sampling are done at random from the dataset to create sample datasets for each model. Bootstrap is the name of this section.

- **GradientBoostingRegressor:-** Gradient boosting is a strong machine learning approach for ensembles.

It's commonly utilised in winning solutions to machine learning contests, such as those on Kaggle, for structured predictive modelling issues like classification and regression on tabular data.

There are several gradient boosting implementations available, including conventional SciPy versions and efficient third-party libraries. Each has its own user interface and even names for the algorithm.

This tutorial will show you how to utilise gradient boosting models in Python for classification and regression.

For the four primary implementations of gradient boosting in Python, standardised code samples are supplied, ready for you to copy-paste and apply in your own predictive modelling project.

- **AdaBoostRegressor:-** Boosting is a type of ensemble machine learning technique that combines the predictions of several weak learners.

A weak learner is a model that is relatively basic but has some ability on the dataset. The AdaBoost (adaptive boosting) algorithm was the first effective attempt to the notion.

As weak learners, the AdaBoost method employs extremely small (one-level) decision trees that are introduced progressively to the ensemble. Each succeeding model in the chain seeks to correct the predictions produced by the model preceding it. This is accomplished by balancing the training dataset to focus more on training cases where previous models failed to predict correctly.

- **ExtraTreesRegressor:-** Extra Trees is a machine learning approach that combines the predictions of many decision trees into a single forecast.

It has something in common with the commonly used random forest algorithm. Although it employs a simpler approach to create the decision trees used as members of the ensemble, it may frequently yield similar or better results than the random forest algorithm.

It's also simple to use, as it just has a few important hyperparameters and logical heuristics for tuning them.

- **SVM:-** SVM stands for Support Vector Machine and is one of the most widely used Supervised Learning algorithms for Classification and Regression issues. However, it is mostly utilised in Machine Learning for Classification difficulties.

The SVM algorithm's purpose is to find the optimum line or decision boundary for categorising n-dimensional space into classes so that additional data points may be readily placed in the proper category in the future. A hyperplane is the name for the optimal choice boundary.

The extreme points/vectors that assist create the hyperplane are chosen via SVM. Support vectors are the extreme instances, and the method is called a Support Vector Machine.

- **XGBRegressor:-** Extreme Gradient Boosting (XGBoost) is an open-source toolkit that implements the gradient boosting technique in an efficient and effective manner.

XGBoost became the go-to approach and frequently the primary component in winning solutions for a variety of issues in machine learning contests shortly after its creation and initial release.

Predicting a numerical value, such as a dollar amount or a height, is an issue in regression predictive modelling. For regression predictive modelling, XGBoost may be utilised directly.

➤

## B. Library used:-

- **numpy:-** NumPy is a Python package that allows you to interact with arrays.

It also provides functions for working with matrices, fourier transforms, and linear algebra.

Travis Oliphant invented NumPy in 2005. It is an open source project that you are free to use. Numerical Python is referred to as NumPy.

We have lists in Python that act as arrays, however they are sluggish to parse. NumPy intends to deliver a 50-fold quicker array object than ordinary Python lists. The array object in NumPy is named ndarray, and it comes with a number of helper methods that make dealing with it a breeze. In data research, when speed

and resources are critical, arrays are widely employed.

- **Pandas:-** Pandas is a Python data analysis package. Wes McKinney created pandas in 2008 to fill a demand for a strong and versatile quantitative analysis tool. It has now evolved to become one of the most used Python tools. It has an incredibly active contributor community.

Pandas is based on two key Python libraries: matplotlib for data visualisation and NumPy for arithmetic operations. Pandas acts as a wrapper around these libraries, enabling you to use fewer lines of code to access many of matplotlib's and NumPy's functions. Pandas'.plot()', for example, integrates numerous matplotlib methods into a single function, allowing you to plot a chart in only a few lines. Most analysts utilised Python for data munging and preprocessing before switching to a more domain specialised language like pandas.

Pandas introduced two new types of data storage objects: Series, which have a list-like structure, and DataFrames, which have a tabular structure, which make analytical jobs easier and eliminate the need to switch tools.

- **Matplotlib and Pyplot:-** Matplotlib is a data visualisation and graphical charting package for Python and its numerical extension NumPy that works across platforms. As a result, it provides an open source alternative to MATLAB. Matplotlib's APIs (Application Programming Interfaces) may also be used to incorporate charts in GUI programmes.

In most cases, a Python matplotlib script is constructed so that only a few lines of code are necessary to produce a visual data plot. Two APIs are overlaid by the matplotlib scripting layer:

Matplotlib is at the top of the pyplot API hierarchy of Python code objects.

pyplot is an OO (Object-Oriented) API collection of objects that can be constructed more easily than pyplot. This allows you to use Matplotlib's backend layers directly.

- **Seaborn:-** Seaborn is a Python module for creating statistical visuals. It is based on matplotlib and tightly interacts with pandas data structures. Seaborn assists you in exploring and comprehending your data. Its charting functions work with dataframes and arrays containing whole datasets, doing the necessary semantic mapping and statistical aggregation internally to generate useful graphs. Its dataset-oriented, declarative API allows you to



concentrate on the meaning of your charts rather than the mechanics of drawing them.

➤ **Sklearn:-** In Python, Scikit-learn (Sklearn) is the most usable and robust machine learning package. It uses a Python consistency interface to deliver a set of fast machine learning and statistical modelling capabilities, such as classification, regression, clustering, and dimensionality reduction. NumPy, SciPy, and Matplotlib are used extensively in this Python-based toolkit.

➤ **Pickle:-** For serialising and de-serializing Python object structures, the pickle package is utilised. Pickling, serialisation, flattening, or marshalling is the process of converting any type of Python object (list, dict, etc.) into byte streams (0s and 1s). Unpickling is a technique that turns the byte stream (produced by pickling) back into Python objects.

Pickling and unpickling are widely used in real-world scenarios because they allow us to effortlessly move data from one server/system to another and then save it in a file or database.

It's best not to unpickle data from an untrustworthy source since it might compromise your security. The pickle module, on the other hand, has no method of knowing or raising an alarm when pickling harmful data.

### C. Solution:-

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
df = pd.read_csv('laptop_data.csv')
df.head()
df.shape
df.info()
df.duplicated().sum()
df.drop(columns=['Unnamed: 0'], inplace=True)
df.head()
df['Ram'] = df['Ram'].str.replace('GB', '')
df['Weight'] = df['Weight'].str.replace('kg', '')
df.head()
df['Ram'] = df['Ram'].astype('int32')
df['Weight'] = df['Weight'].astype('float32')
df.info()
import seaborn as sns
```

```
sns.distplot(df['Price'])
df['Company'].value_counts().plot(kind='bar')
sns.barplot(x=df['Company'], y=df['Price'])
plt.xticks(rotation='vertical')
plt.show()
df['TypeName'].value_counts().plot(kind='bar')
sns.barplot(x=df['TypeName'], y=df['Price'])
plt.xticks(rotation='vertical')
plt.show()
sns.distplot(df['Inches'])
sns.scatterplot(x=df['Inches'], y=df['Price'])
df['ScreenResolution'].value_counts()
df['Touchscreen'] = df['ScreenResolution'].apply(lambda x: 1 if 'Touchscreen' in x else 0)
df.sample(5)
df['Touchscreen'].value_counts().plot(kind='bar')
sns.barplot(x=df['Touchscreen'], y=df['Price'])
df['Ips'] = df['ScreenResolution'].apply(lambda x: 1 if 'IPS' in x else 0)
df.head()
df['Ips'].value_counts().plot(kind='bar')
sns.barplot(x=df['Ips'], y=df['Price'])
new = df['ScreenResolution'].str.split('x', n=1, expand=True)
df['X_res'] = new[0]
df['Y_res'] = new[1]
df.sample(5)
df['X_res'] = df['X_res'].str.replace(',', '').str.findall(r'(\d+\.\d+)').apply(lambda x: x[0])
df.head()
df['X_res'] = df['X_res'].astype('int')
df['Y_res'] = df['Y_res'].astype('int')
df.info()
df.corr()['Price']
df['ppi'] = (((df['X_res']**2) + (df['Y_res']**2))**0.5/df['Inches']).astype('float')
df.corr()['Price']
df.drop(columns=['ScreenResolution'], inplace=True)
```

```

df.head()
df.drop(columns=['Inches','X_res','Y_res'],inplace=True)
df.head()
df['Cpu'].value_counts()
df['Cpu Name'] = df['Cpu'].apply(lambda
    x: " ".join(x.split()[0:3]))
df.head()
def fetch_processor(text):
    if text == 'Intel Core i7' or text
    == 'Intel Core i5' or text == 'Intel Co
    re i3':
        return text
    else:
        if text.split()[0] == 'Intel':
            return 'Other Intel Process
or'
        else:
            return 'AMD Processor'
df['Cpu brand'] = df['Cpu Name'].apply(
fetch_processor)
df.head()
df['Cpu brand'].value_counts().plot(kin
d='bar')
sns.barplot(x=df['Cpu brand'],y=df['Pri
ce'])
plt.xticks(rotation='vertical')
plt.show()
df.drop(columns=['Cpu','Cpu Name'],inpl
ace=True)
df.head()
df['Ram'].value_counts().plot(kind='bar
')
sns.barplot(x=df['Ram'],y=df['Price'])
plt.xticks(rotation='vertical')
plt.show()
df['Memory'].value_counts()
df['Memory'] = df['Memory'].astype(str)
.replace('\.0', '', regex=True)
df["Memory"] = df["Memory"].str.replace
('GB', '')
df["Memory"] = df["Memory"].str.replace
('TB', '000')
new = df["Memory"].str.split("+", n = 1
, expand = True)

df["first"] = new[0]
df["first"] = df["first"].str.strip()

df["second"] = new[1]

```

```

df["Layer1HDD"] = df["first"].apply(lam
bda x: 1 if "HDD" in x else 0)
df["Layer1SSD"] = df["first"].apply(lam
bda x: 1 if "SSD" in x else 0)
df["Layer1Hybrid"] = df["first"].apply(
lambda x: 1 if "Hybrid" in x else 0)
df["Layer1Flash_Storage"] = df["first"]
.apply(lambda x: 1 if "Flash Storage" i
n x else 0)

df['first'] = df['first'].str.replace(r
'\D', '')

df["second"].fillna("0", inplace = True
)

df["Layer2HDD"] = df["second"].apply(la
mbda x: 1 if "HDD" in x else 0)
df["Layer2SSD"] = df["second"].apply(la
mbda x: 1 if "SSD" in x else 0)
df["Layer2Hybrid"] = df["second"].apply
(lambda x: 1 if "Hybrid" in x else 0)
df["Layer2Flash_Storage"] = df["second"
].apply(lambda x: 1 if "Flash Storage"
in x else 0)

df['second'] = df['second'].str.replace
(r'\D', '')

df["first"] = df["first"].astype(int)
df["second"] = df["second"].astype(int)

df["HDD"] = (df["first"]*df["Layer1HDD"]+
df["second"]*df["Layer2HDD"])
df["SSD"] = (df["first"]*df["Layer1SSD"]+
df["second"]*df["Layer2SSD"])
df["Hybrid"] = (df["first"]*df["Layer1Hyb
rid"]+df["second"]*df["Layer2Hybrid"])
df["Flash_Storage"] = (df["first"]*df["La
yer1Flash_Storage"]+df["second"]*df["La
yer2Flash_Storage"])

df.drop(columns=['first', 'second', 'La
yer1HDD', 'Layer1SSD', 'Layer1Hybrid',
'Layer1Flash_Storage', 'Layer2HDD',
'Layer2SSD', 'Layer2Hybrid',
'Layer2Flash_Storage'],inplace=Tr
ue)
df.sample(5)

```

```

df.drop(columns=['Memory'],inplace=True)
df.head()
df.corr()['Price']
df.drop(columns=['Hybrid','Flash_Storage'],inplace=True)
df.head()
df['Gpu'].value_counts()
df['Gpu brand'] = df['Gpu'].apply(lambda x:x.split()[0])
df.head()
df['Gpu brand'].value_counts()
df = df[df['Gpu brand'] != 'ARM']
df['Gpu brand'].value_counts()
sns.barplot(x=df['Gpu brand'],y=df['Price'],estimator=np.median)
plt.xticks(rotation='vertical')
plt.show()
df.drop(columns=['Gpu'],inplace=True)
df.head()
df['OpSys'].value_counts()
sns.barplot(x=df['OpSys'],y=df['Price'])
plt.xticks(rotation='vertical')
plt.show()
def cat_os(inp):
    if inp == 'Windows 10' or inp == 'Windows 7' or inp == 'Windows 10 S':
        return 'Windows'
    elif inp == 'macOS' or inp == 'Mac OS X':
        return 'Mac'
    else:
        return 'Others/No OS/Linux'
df['os'] = df['OpSys'].apply(cat_os)
df.head()
df.drop(columns=['OpSys'],inplace=True)
sns.barplot(x=df['os'],y=df['Price'])
plt.xticks(rotation='vertical')
plt.show()
sns.distplot(df['Weight'])
sns.scatterplot(x=df['Weight'],y=df['Price'])
df.corr()['Price']
sns.heatmap(df.corr())
sns.distplot(np.log(df['Price']))
X = df.drop(columns=['Price'])
y = np.log(df['Price'])
X
y

```

```

from sklearn.model_selection import train_test_split
X_train,X_test,y_train,y_test = train_test_split(X,y,test_size=0.15,random_state=2)
X_train
from sklearn.compose import ColumnTransformer
from sklearn.pipeline import Pipeline
from sklearn.preprocessing import OneHotEncoder
from sklearn.metrics import r2_score,mean_absolute_error
from sklearn.linear_model import LinearRegression,Ridge,Lasso
from sklearn.neighbors import KNeighborsRegressor
from sklearn.tree import DecisionTreeRegressor
from sklearn.ensemble import RandomForestRegressor,GradientBoostingRegressor,AdaBoostRegressor,ExtraTreesRegressor
from sklearn.svm import SVR
from xgboost import XGBRegressor
step1 = ColumnTransformer(transformers=[
    ('col_tnf',OneHotEncoder(sparse=False,drop='first'),[0,1,7,10,11])
],remainder='passthrough')
step2 = LinearRegression()
pipe = Pipeline([
    ('step1',step1),
    ('step2',step2)
])
pipe.fit(X_train,y_train)
y_pred = pipe.predict(X_test)
print('R2 score',r2_score(y_test,y_pred))
print('MAE',mean_absolute_error(y_test,y_pred))
step1 = ColumnTransformer(transformers=[
    ('col_tnf',OneHotEncoder(sparse=False,drop='first'),[0,1,7,10,11])
],remainder='passthrough')

```



```

],remainder='passthrough')

step2 = Ridge(alpha=10)

pipe = Pipeline([
    ('step1', step1),
    ('step2', step2)
])

pipe.fit(X_train, y_train)

y_pred = pipe.predict(X_test)

print('R2 score', r2_score(y_test, y_pred))
print('MAE', mean_absolute_error(y_test, y_pred))
step1 = ColumnTransformer(transformers=[
    ('col_tnf', OneHotEncoder(sparse=False, drop='first'), [0, 1, 7, 10, 11])
], remainder='passthrough')

step2 = Lasso(alpha=0.001)

pipe = Pipeline([
    ('step1', step1),
    ('step2', step2)
])

pipe.fit(X_train, y_train)

y_pred = pipe.predict(X_test)

print('R2 score', r2_score(y_test, y_pred))
print('MAE', mean_absolute_error(y_test, y_pred))
step1 = ColumnTransformer(transformers=[
    ('col_tnf', OneHotEncoder(sparse=False, drop='first'), [0, 1, 7, 10, 11])
], remainder='passthrough')

step2 = KNeighborsRegressor(n_neighbors=3)

pipe = Pipeline([
    ('step1', step1),
    ('step2', step2)
])

```

```

])

pipe.fit(X_train, y_train)

y_pred = pipe.predict(X_test)

print('R2 score', r2_score(y_test, y_pred))
print('MAE', mean_absolute_error(y_test, y_pred))
step1 = ColumnTransformer(transformers=[
    ('col_tnf', OneHotEncoder(sparse=False, drop='first'), [0, 1, 7, 10, 11])
], remainder='passthrough')

step2 = DecisionTreeRegressor(max_depth=8)

pipe = Pipeline([
    ('step1', step1),
    ('step2', step2)
])

pipe.fit(X_train, y_train)

y_pred = pipe.predict(X_test)

print('R2 score', r2_score(y_test, y_pred))
print('MAE', mean_absolute_error(y_test, y_pred))

step1 = ColumnTransformer(transformers=[
    ('col_tnf', OneHotEncoder(sparse=False, drop='first'), [0, 1, 7, 10, 11])
], remainder='passthrough')

step2 = RandomForestRegressor(n_estimators=100,
                                random_state=3,
                                max_samples=0.5,
                                max_features=0.75,
                                max_depth=15)

```

```

pipe = Pipeline([
    ('step1', step1),
    ('step2', step2)
])

pipe.fit(X_train, y_train)

y_pred = pipe.predict(X_test)

print('R2 score', r2_score(y_test, y_pred))
print('MAE', mean_absolute_error(y_test, y_pred))
step1 = ColumnTransformer(transformers=[
    ('col_tnf', OneHotEncoder(sparse=False, drop='first'), [0, 1, 7, 10, 11])
], remainder='passthrough')

step2 = ExtraTreesRegressor(n_estimators=100,
                             random_state=3,
                             max_samples=0.5,
                             max_features=0.75,
                             max_depth=15)

pipe = Pipeline([
    ('step1', step1),
    ('step2', step2)
])

pipe.fit(X_train, y_train)

y_pred = pipe.predict(X_test)

print('R2 score', r2_score(y_test, y_pred))
print('MAE', mean_absolute_error(y_test, y_pred))
step1 = ColumnTransformer(transformers=[
    ('col_tnf', OneHotEncoder(sparse=False, drop='first'), [0, 1, 7, 10, 11])
], remainder='passthrough')

```

```

step2 = AdaBoostRegressor(n_estimators=15, learning_rate=1.0)

pipe = Pipeline([
    ('step1', step1),
    ('step2', step2)
])

pipe.fit(X_train, y_train)

y_pred = pipe.predict(X_test)

print('R2 score', r2_score(y_test, y_pred))
print('MAE', mean_absolute_error(y_test, y_pred))

step1 = ColumnTransformer(transformers=[
    ('col_tnf', OneHotEncoder(sparse=False, drop='first'), [0, 1, 7, 10, 11])
], remainder='passthrough')

step2 = GradientBoostingRegressor(n_estimators=500)

pipe = Pipeline([
    ('step1', step1),
    ('step2', step2)
])

pipe.fit(X_train, y_train)

y_pred = pipe.predict(X_test)

print('R2 score', r2_score(y_test, y_pred))
print('MAE', mean_absolute_error(y_test, y_pred))
step1 = ColumnTransformer(transformers=[
    ('col_tnf', OneHotEncoder(sparse=False, drop='first'), [0, 1, 7, 10, 11])
], remainder='passthrough')

step2 = XGBRegressor(n_estimators=45, max_depth=5, learning_rate=0.5)

pipe = Pipeline([

```

```

        ('step1', step1),
        ('step2', step2)
    ])

pipe.fit(X_train, y_train)

y_pred = pipe.predict(X_test)

print('R2 score', r2_score(y_test, y_pred))
print('MAE', mean_absolute_error(y_test, y_pred))

from sklearn.ensemble import VotingRegressor, StackingRegressor

step1 = ColumnTransformer(transformers=[
    ('col_tnf', OneHotEncoder(sparse=False, drop='first'), [0, 1, 7, 10, 11])
], remainder='passthrough')

rf = RandomForestRegressor(n_estimators=350, random_state=3, max_samples=0.5, max_features=0.75, max_depth=15)
gbdt = GradientBoostingRegressor(n_estimators=100, max_features=0.5)
xgb = XGBRegressor(n_estimators=25, learning_rate=0.3, max_depth=5)
et = ExtraTreesRegressor(n_estimators=100, random_state=3, max_samples=0.5, max_features=0.75, max_depth=10)

step2 = VotingRegressor([
    ('rf', rf), ('gbdt', gbdt), ('xgb', xgb), ('et', et)
], weights=[5, 1, 1, 1])

pipe = Pipeline([
    ('step1', step1),
    ('step2', step2)
])

pipe.fit(X_train, y_train)

y_pred = pipe.predict(X_test)

print('R2 score', r2_score(y_test, y_pred))

```

```

print('MAE', mean_absolute_error(y_test, y_pred))

step1 = ColumnTransformer(transformers=[
    ('col_tnf', OneHotEncoder(sparse=False, drop='first'), [0, 1, 7, 10, 11])
], remainder='passthrough')

step2 = AdaBoostRegressor(n_estimators=15, learning_rate=1.0)

pipe = Pipeline([
    ('step1', step1),
    ('step2', step2)
])

pipe.fit(X_train, y_train)

y_pred = pipe.predict(X_test)

print('R2 score', r2_score(y_test, y_pred))
print('MAE', mean_absolute_error(y_test, y_pred))

step1 = ColumnTransformer(transformers=[
    ('col_tnf', OneHotEncoder(sparse=False, drop='first'), [0, 1, 7, 10, 11])
], remainder='passthrough')

step2 = GradientBoostingRegressor(n_estimators=500)

pipe = Pipeline([
    ('step1', step1),
    ('step2', step2)
])

pipe.fit(X_train, y_train)

y_pred = pipe.predict(X_test)

print('R2 score', r2_score(y_test, y_pred))
print('MAE', mean_absolute_error(y_test, y_pred))

```

```

step1 = ColumnTransformer(transformers=
[
    ('col_tnf',OneHotEncoder(sparse=False,drop='first'),[0,1,7,10,11])
],remainder='passthrough')

```

```

step2 = XGBRegressor(n_estimators=45,max_depth=5,learning_rate=0.5)

```

```

pipe = Pipeline([
    ('step1',step1),
    ('step2',step2)
])

```

```

pipe.fit(X_train,y_train)

```

```

y_pred = pipe.predict(X_test)

```

```

print('R2 score',r2_score(y_test,y_pred))
print('MAE',mean_absolute_error(y_test,y_pred))

```

```

from sklearn.ensemble import VotingRegressor,StackingRegressor

```

```

step1 = ColumnTransformer(transformers=
[
    ('col_tnf',OneHotEncoder(sparse=False,drop='first'),[0,1,7,10,11])
],remainder='passthrough')

```

```

rf = RandomForestRegressor(n_estimators=350,random_state=3,max_samples=0.5,max_features=0.75,max_depth=15)
gbdt = GradientBoostingRegressor(n_estimators=100,max_features=0.5)
xgb = XGBRegressor(n_estimators=25,learning_rate=0.3,max_depth=5)
et = ExtraTreesRegressor(n_estimators=100,random_state=3,max_samples=0.5,max_features=0.75,max_depth=10)

```

```

step2 = VotingRegressor([('rf', rf), ('gbdt', gbdt), ('xgb',xgb), ('et',et)],weights=[5,1,1,1])

```

```

pipe = Pipeline([
    ('step1',step1),
    ('step2',step2)
])

```

```

pipe.fit(X_train,y_train)

```

```

y_pred = pipe.predict(X_test)

```

```

print('R2 score',r2_score(y_test,y_pred))
print('MAE',mean_absolute_error(y_test,y_pred))

```

```

from sklearn.ensemble import VotingRegressor,StackingRegressor

```

```

step1 = ColumnTransformer(transformers=
[
    ('col_tnf',OneHotEncoder(sparse=False,drop='first'),[0,1,7,10,11])
],remainder='passthrough')

```

```

estimators = [
    ('rf', RandomForestRegressor(n_estimators=350,random_state=3,max_samples=0.5,max_features=0.75,max_depth=15)),
    ('gbdt',GradientBoostingRegressor(n_estimators=100,max_features=0.5)),
    ('xgb', XGBRegressor(n_estimators=25,learning_rate=0.3,max_depth=5))
]

```

```

step2 = StackingRegressor(estimators=estimators, final_estimator=Ridge(alpha=100))

```

```

pipe = Pipeline([
    ('step1',step1),
    ('step2',step2)
])

```

```

pipe.fit(X_train,y_train)

```

```

y_pred = pipe.predict(X_test)

```

```

print('R2 score',r2_score(y_test,y_pred))

```

```
print('MAE',mean_absolute_error(y_test,
y_pred))

import pickle

pickle.dump(df,open('df.pkl','wb'))
pickle.dump(pipe,open('pipe.pkl','wb'))
df
X_train
```

## VII. Experimental Setup:-

For implementation of machine learning we have using python.

➤**Python:-** Python is a popular high-level programming language for general-purpose applications. Guido van Rossum designed it in 1991, and the Python Software Foundation developed it. Its syntax lets programmers to express concepts in fewer lines of code, and it was primarily built with code readability in mind.

➤**Why :-** Python's extensive library ecosystem allows it to be used in almost every IT field, and it makes it easier for programmers to design things.

Python was created with readability in mind, making it straightforward to learn for beginners.

Its adaptability allows it to be used with other languages.

Python may run on a variety of operating systems (Windows, Mac, Linux, Raspberry Pi, etc).

Python has a syntax that is quite close to English.

Python's syntax allows programmers to build programmes in less lines than other programming languages.

Python is an interpreter language, meaning that code may be run immediately after it is written.

Prototyping may thus be completed quite rapidly.

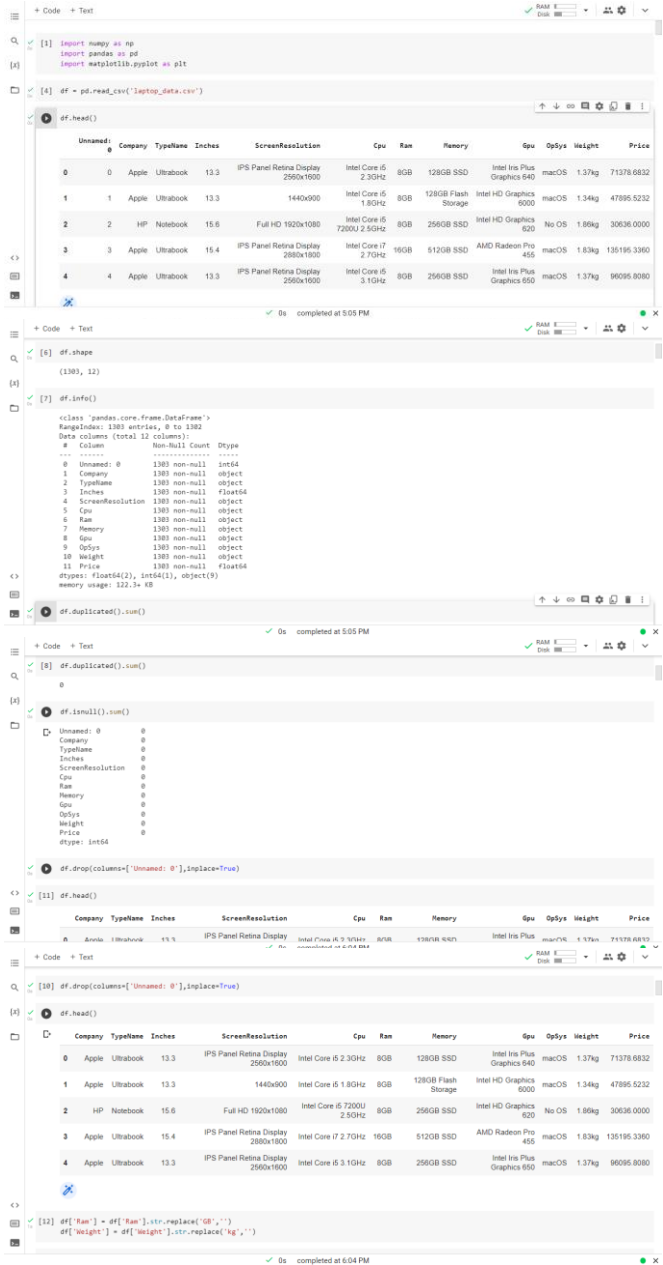
Python can be approached in three ways: procedural, object-oriented, and functional.

➤ **Google colab:-** Google Colab is a Google-provided freemium tool. For deep learning technologies, Google Colab is an excellent tool. Google Colab provides a jupyter notebook that runs entirely in the cloud. Most importantly, there is no need to set up the environment, and the notebook you create can be edited simultaneously. Many major machine learning libraries are also available in Google Colab.

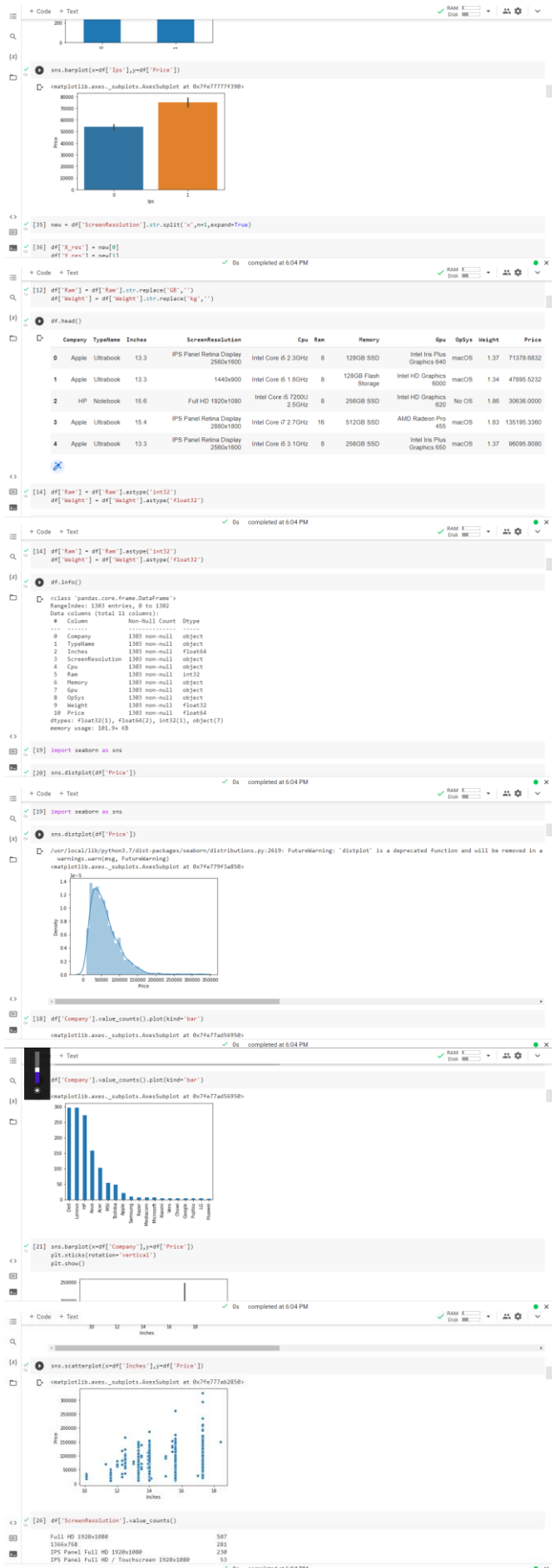
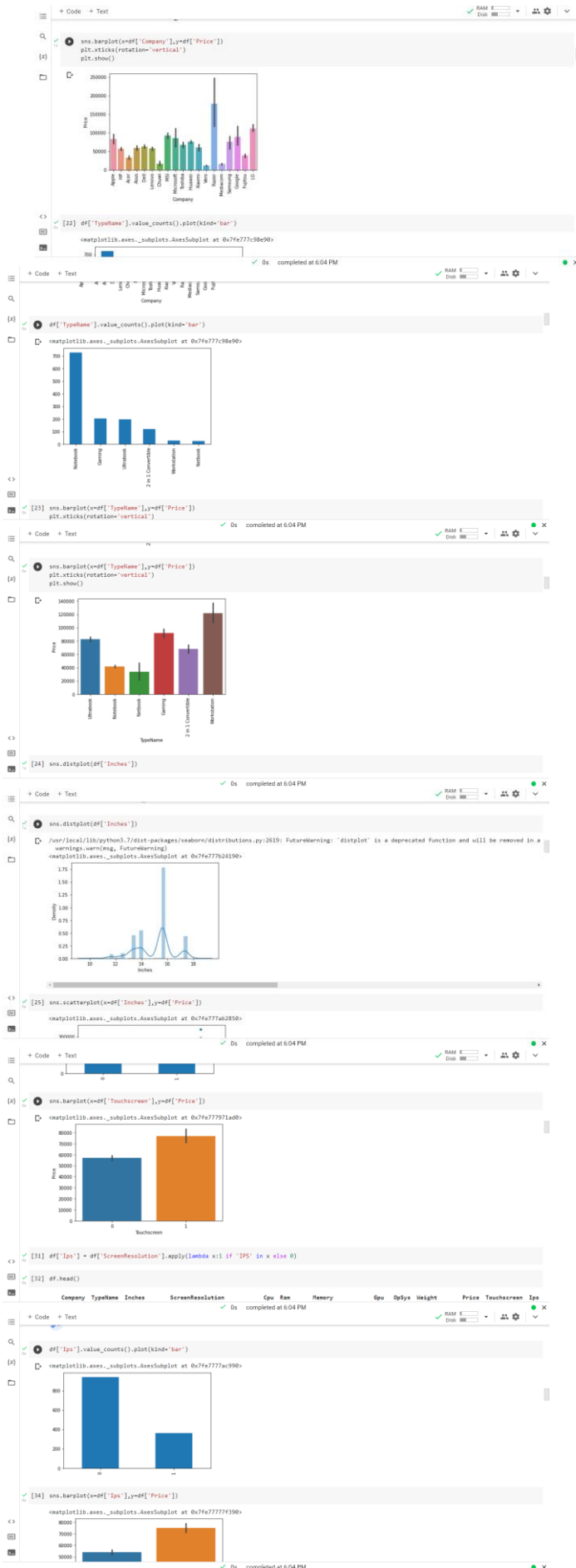
➤ **Why:-** No installation is required: one of the benefits of using Google Colab is that you do not need

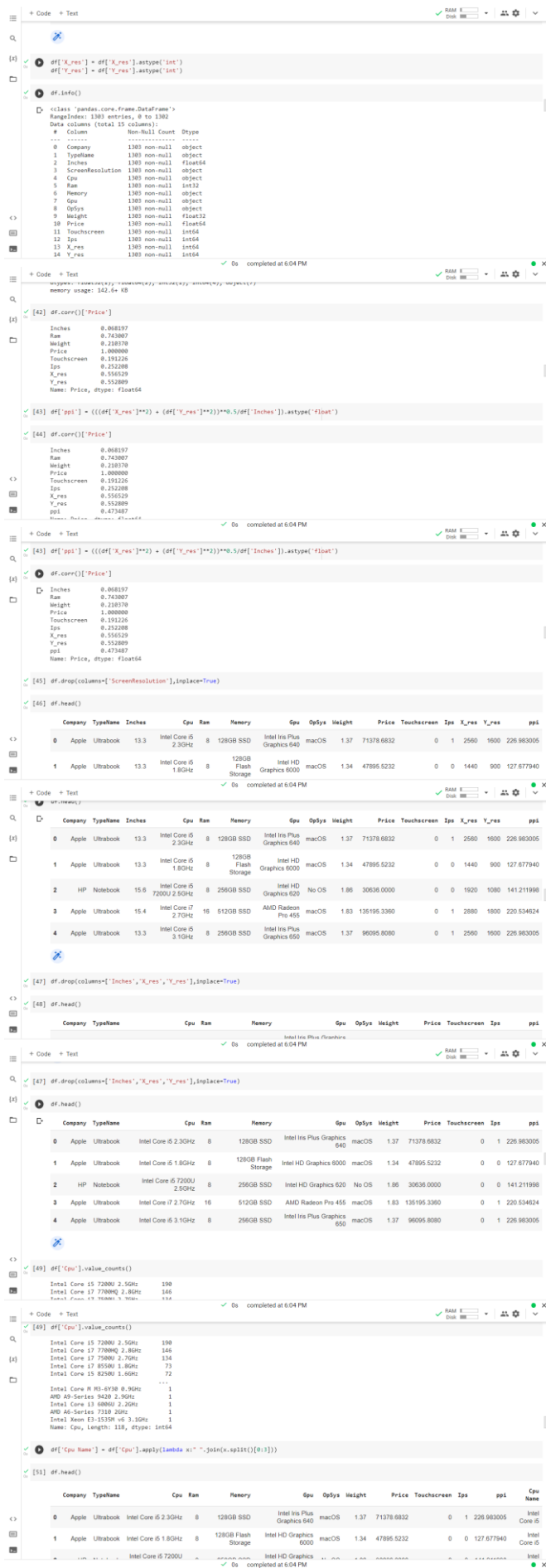
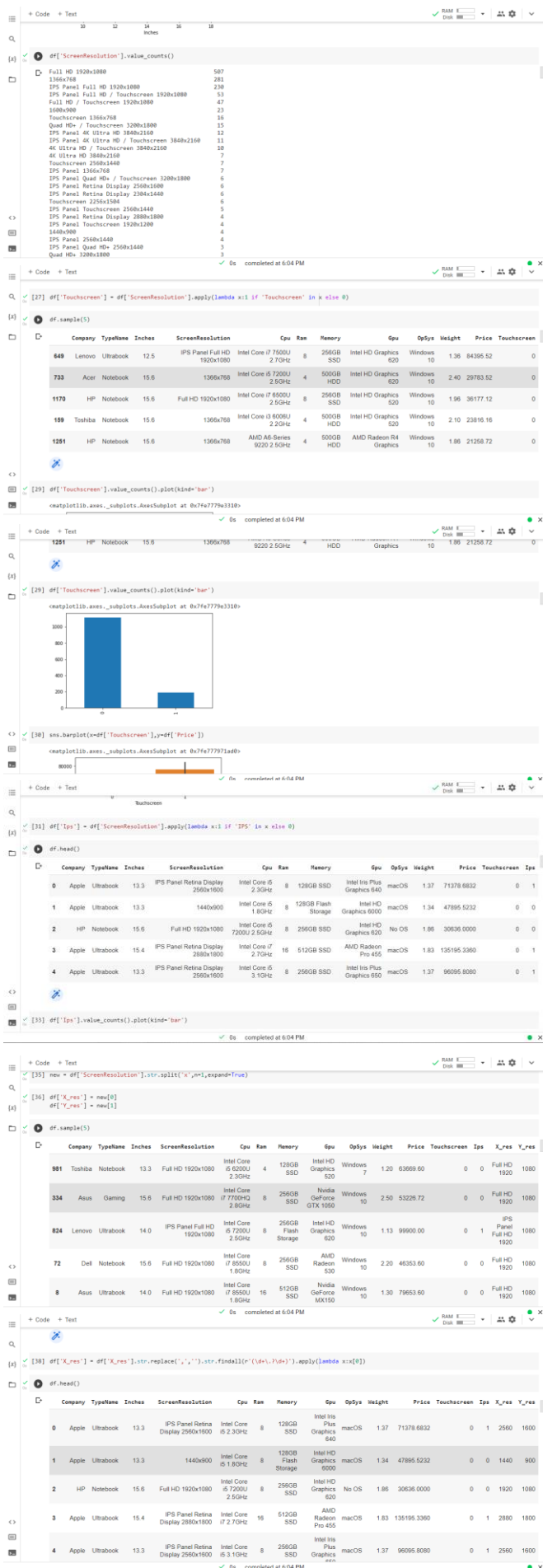
to install it on your local machine; it runs entirely on the cloud, making it accessible anywhere there is an internet connection. Because it runs on a cloud-based local system, it has no effect on the Google Colab's performance. We only require an internet connection and a browser. Google Colab also has a collaboration function that allows numerous developers to work on the same project at the same time.

## VIII. Result:-









```
+ Code + Text
[52] def fetch_processor(text):
    if text == 'Intel Core i7' or text == 'Intel Core i5' or text == 'Intel Core i3':
        return text
    else:
        if text.split()[0] == 'Intel':
            return 'Other Intel Processor'
        else:
            return 'AMD Processor'

[53] df['Cpu brand'] = df['Cpu Name'].apply(fetch_processor)

df.head()

Company TypeName Cpu Ram Memory Gpu OpSys Weight Price Touchscreen Ips ppi Cpu brand
0 Apple Ultrabook Intel Core i5 8 128GB SSD Intel Iris Plus Graphics 640 macOS 1.37 71378.6832 0 1 226.983005 Intel Core i5
1 Apple Ultrabook Intel Core i5 8 128GB Flash Storage Intel HD Graphics 6000 macOS 1.34 47895.5232 0 0 127.677940 Intel Core i5
2 HP Notebook Intel Core i5 8 256GB SSD Intel HD Graphics 620 No OS 1.86 30636.0000 0 0 141.211998 Intel Core i5
3 Apple Ultrabook Intel Core i7 2.7GHz 16 512GB SSD AMD Radeon Pro 455 macOS 1.83 135195.3360 0 1 220.534624 Intel Core i7

[54] df['Cpu brand'].value_counts().plot(kind='bar')
<matplotlib.axes._subplots.AxesSubplot at 0x7f777663b10>

sns.barplot(x=df['Cpu brand'], y=df['Price'])
plt.xticks(rotation='vertical')
plt.show()

[56] sns.barplot(x=df['Cpu brand'], y=df['Price'])
plt.xticks(rotation='vertical')
plt.show()

[57] df.drop(columns=['Cpu', 'Cpu Name'], inplace=True)

df.head()

Company TypeName Ram Memory Gpu OpSys Weight Price Touchscreen Ips ppi Cpu brand
0 Apple Ultrabook 8 128GB SSD Intel Iris Plus Graphics 640 macOS 1.37 71378.6832 0 1 226.983005 Intel Core i5
1 Apple Ultrabook 8 128GB Flash Storage Intel HD Graphics 6000 macOS 1.34 47895.5232 0 0 127.677940 Intel Core i5
2 HP Notebook 8 256GB SSD Intel HD Graphics 620 No OS 1.86 30636.0000 0 0 141.211998 Intel Core i5
3 Apple Ultrabook 16 512GB SSD AMD Radeon Pro 455 macOS 1.83 135195.3360 0 1 220.534624 Intel Core i7
4 Apple Ultrabook 8 256GB SSD Intel Iris Plus Graphics 650 macOS 1.37 96095.8080 0 1 226.983005 Intel Core i5

[59] df['Ram'].value_counts().plot(kind='bar')
<matplotlib.axes._subplots.AxesSubplot at 0x7f7775da90>

[60] sns.barplot(x=df['Ram'], y=df['Price'])
plt.xticks(rotation='vertical')
plt.show()

[61] df['Memory'].value_counts()

256GB SSD 412
1TB HDD 223
512GB HDD 112
```

```
+ Code + Text
[62] df['Memory'] = df['Memory'].astype(str).replace('\n', '', regex=True)
df['Memory'] = df['Memory'].str.replace('GB', '')
df['Memory'] = df['Memory'].str.replace('TB', '900')
new = df['Memory'].str.split('-', n=1, expand=True)

df['first'] = new[0]
df['first'] = df['first'].str.strip()

df['second'] = new[1]

df['LayerHDD'] = df['first'].apply(lambda x: 1 if 'HDD' in x else 0)
df['LayerSSD'] = df['first'].apply(lambda x: 1 if 'SSD' in x else 0)
df['LayerHybrid'] = df['first'].apply(lambda x: 1 if 'Hybrid' in x else 0)
df['LayerFlash_Storage'] = df['first'].apply(lambda x: 1 if 'Flash Storage' in x else 0)

df['first'] = df['first'].str.replace('-', '')
df['second'] = df['second'].str.replace('-', '')

df['LayerHDD'] = df['second'].apply(lambda x: 1 if 'HDD' in x else 0)
df['LayerSSD'] = df['second'].apply(lambda x: 1 if 'SSD' in x else 0)
df['LayerHybrid'] = df['second'].apply(lambda x: 1 if 'Hybrid' in x else 0)
df['LayerFlash_Storage'] = df['second'].apply(lambda x: 1 if 'Flash Storage' in x else 0)

df['second'] = df['second'].str.replace('-', '')

df['HDD'] = df['first'] & df['second']
df['SSD'] = df['first'] & df['second']
df['Hybrid'] = df['first'] & df['second']
df['Flash_Storage'] = df['first'] & df['second']

df.drop(columns=['first', 'second', 'LayerHDD', 'LayerSSD', 'LayerHybrid', 'LayerFlash_Storage', 'LayerHDD', 'LayerSSD', 'LayerHybrid', 'LayerFlash_Storage'], inplace=True)

./usr/local/lib/python3.7/dist-packages/ipykernel_launcher.py:16: FutureWarning: The default value of regex will change from True to False in a fut
./usr/local/lib/python3.7/dist-packages/ipykernel_launcher.py:25: FutureWarning: The default value of regex will change from True to False in a fut

[63] df.sample(5)

Company TypeName Ram Memory Gpu OpSys Weight Price Touchscreen Ips ppi Cpu brand HDD SSD Hybrid Flash_Storage
908 HP Notebook 8 512 SSD Intel HD Graphics 620 Windows 10 1.26 79014.24 0 0 176.232574 Intel Core i7 0 512 0 0
189 Dell Notebook 8 128 SSD + 1000 HDD AMD Radeon S50 Windows 10 2.80 67808.80 0 0 127.336675 Intel Core 1000 128 0 0
89 Asus Gaming 12 1000 Nvidia GeForce GTX 1050 Linux 3.00 50562.72 0 0 127.336675 Intel Core i7 0 0 0 0
133 Acer Notebook 4 256 SSD Nvidia GeForce MX150 Windows 10 2.20 30476.16 0 0 141.211998 Intel Core i3 0 256 0 0
364 Lenovo Notebook 8 256 SSD Nvidia GeForce 920MX No OS 2.20 26586.72 0 0 100.454870 Intel Core i3 0 256 0 0

[64] df.drop(columns=['Memory'], inplace=True)

df.head()

Company TypeName Ram Gpu OpSys Weight Price Touchscreen Ips ppi Cpu brand HDD SSD Hybrid Flash_Storage
0 Apple Ultrabook 8 Intel Iris Plus Graphics 640 macOS 1.37 71378.6832 0 1 226.983005 Intel Core i5 0 128 0 0
1 Apple Ultrabook 8 Intel HD Graphics 6000 macOS 1.34 47895.5232 0 0 127.677940 Intel Core i5 0 0 0 128
2 HP Notebook 8 Intel HD Graphics 620 No OS 1.86 30636.0000 0 0 141.211998 Intel Core i5 0 256 0 0
3 Apple Ultrabook 16 AMD Radeon Pro 455 macOS 1.83 135195.3360 0 1 220.534624 Intel Core i7 0 512 0 0
4 Apple Ultrabook 8 Intel Iris Plus Graphics 650 macOS 1.37 96095.8080 0 1 226.983005 Intel Core i5 0 256 0 0

[66] df.corr()['Price']

Price
Ram 0.743807
Weight 0.228378
Price 1.000000
Touchscreen 0.191226
Ips 0.252288
ppi 0.471487
HDD -0.895441
SSD 0.578799
Hybrid 0.807985
Flash_Storage -0.848511
Name: Price, dtype: float64

[67] df.drop(columns=['Hybrid', 'Flash_Storage'], inplace=True)

[68] df.head()

Company TypeName Ram Gpu OpSys Weight Price Touchscreen Ips ppi Cpu brand HDD SSD
0 Apple Ultrabook 8 Intel Iris Plus Graphics 640 macOS 1.37 71378.6832 0 1 226.983005 Intel Core i5 0 128
1 Apple Ultrabook 8 Intel HD Graphics 6000 macOS 1.34 47895.5232 0 0 127.677940 Intel Core i5 0 0
2 HP Notebook 8 Intel HD Graphics 620 No OS 1.86 30636.0000 0 0 141.211998 Intel Core i5 0 256
3 Apple Ultrabook 16 AMD Radeon Pro 455 macOS 1.83 135195.3360 0 1 220.534624 Intel Core i7 0 512
```

+ Code + Text

RAM 8

Disk 88

Q

[69] df['Gpu'].value\_counts()

Intel HD Graphics 620 281  
Intel HD Graphics 520 185  
Intel UHD Graphics 620 68  
Nvidia GeForce GTX 1050 66  
Nvidia GeForce GTX 1060 48  
...  
AMD Radeon R5 520 1  
AMD Radeon R7 1  
Intel HD Graphics 540 1  
AMD Radeon 540 1  
ARM Mali T860 MP4 1  
Name: Gpu, Length: 110, dtype: int64

✓ 0s completed at 6:04 PM

Q

[70] df['Gpu brand'] = df['Gpu'].apply(lambda x:x.split()[0])

✓ 0s completed at 6:04 PM

Q

[71] df.head()

	Company	TypeName	Ram	Gpu	OpSys	Weight	Price	Touchscreen	Ips	ppi	Cpu brand	HDD	SSD	Gpu brand
0	Apple	Ultrabook	8	Intel Iris Plus Graphics 640	macOS	1.37	71378.6632	0	1	226.983005	Intel Core i5	0	128	Intel
1	Apple	Ultrabook	8	Intel HD Graphics 6000	macOS	1.34	47895.5232	0	0	127.677940	Intel Core i5	0	0	Intel
2	HP	Notebook	8	Intel HD Graphics 620	No OS	1.86	30636.0000	0	0	141.211998	Intel Core i5	0	256	Intel
3	Apple	Ultrabook	16	AMD Radeon Pro 455	macOS	1.83	135195.3360	0	1	220.534624	Intel Core i7	0	512	AMD

✓ 0s completed at 6:04 PM

Q

[72] df['Gpu brand'].value\_counts()

Intel 722  
Nvidia 400  
AMD 180  
ARM 1  
Name: Gpu brand, dtype: int64

✓ 0s completed at 6:04 PM

Q

[73] df = df[df['Gpu brand'] != 'ARM']

✓ 0s completed at 6:04 PM

Q

[74] df['Gpu brand'].value\_counts()

Intel 722  
Nvidia 400  
AMD 180  
Name: Gpu brand, dtype: int64

✓ 0s completed at 6:04 PM

Q

[75] sns.barplot(x=df['Gpu brand'], y=df['Price'], estimator=np.median)  
plt.xticks(rotation='vertical')  
plt.show()

✓ 0s completed at 6:04 PM

Q

[76] df.drop(columns=['Gpu'], inplace=True)

✓ 0s completed at 6:04 PM

Q

[77] df.head()

	Company	TypeName	Ram	OpSys	Weight	Price	Touchscreen	Ips	ppi	Cpu brand	HDD	SSD	Gpu brand
0	Apple	Ultrabook	8	macOS	1.37	71378.6632	0	1	226.983005	Intel Core i5	0	128	Intel
1	Apple	Ultrabook	8	macOS	1.34	47895.5232	0	0	127.677940	Intel Core i5	0	0	Intel
2	HP	Notebook	8	No OS	1.86	30636.0000	0	0	141.211998	Intel Core i5	0	256	Intel
3	Apple	Ultrabook	16	macOS	1.83	135195.3360	0	1	220.534624	Intel Core i7	0	512	AMD
4	Apple	Ultrabook	8	macOS	1.37	96095.8080	0	1	226.983005	Intel Core i5	0	256	Intel

✓ 0s completed at 6:04 PM

Q

[78] df['OpSys'].value\_counts()

Windows 10 1873  
No OS 66  
Linux 62  
Windows 7 45  
Chrome OS 26  
macOS 13  
Mac OS X 8  
Windows 10 S 8  
Android 2  
Name: OpSys, dtype: int64

✓ 0s completed at 6:04 PM

Q

[79] sns.barplot(x=df['OpSys'], y=df['Price'], estimator=np.median)  
plt.xticks(rotation='vertical')  
plt.show()

✓ 0s completed at 6:04 PM

Q

[80] df.drop(columns=['OpSys'], inplace=True)

✓ 0s completed at 6:04 PM

Q

[81] df['os'] = df['OpSys'].apply(cat\_os)

✓ 0s completed at 6:04 PM

Q

[82] df.head()

	Company	TypeName	Ram	OpSys	Weight	Price	Touchscreen	Ips	ppi	Cpu brand	HDD	SSD	Gpu brand	os
0	Apple	Ultrabook	8	macOS	1.37	71378.6632	0	1	226.983005	Intel Core i5	0	128	Intel	Mac
1	Apple	Ultrabook	8	macOS	1.34	47895.5232	0	0	127.677940	Intel Core i5	0	0	Intel	Mac
2	HP	Notebook	8	No OS	1.86	30636.0000	0	0	141.211998	Intel Core i5	0	256	Intel	Others/No OS/Linux
3	Apple	Ultrabook	16	macOS	1.83	135195.3360	0	1	220.534624	Intel Core i7	0	512	AMD	Mac
4	Apple	Ultrabook	8	macOS	1.37	96095.8080	0	1	226.983005	Intel Core i5	0	256	Intel	Mac

✓ 0s completed at 6:04 PM

Q

[83] df.drop(columns=['OpSys'], inplace=True)

✓ 0s completed at 6:04 PM

+ Code + Text

RAM 8

Disk 88

Q

[83] df.drop(columns=['OpSys'], inplace=True)

✓ 0s completed at 6:04 PM

Q

[84] sns.barplot(x=df['os'], y=df['Price'])  
plt.xticks(rotation='vertical')  
plt.show()

✓ 0s completed at 6:04 PM

Q

[85] sns.distplot(df['Weight'])

✓ 0s completed at 6:04 PM

Q

[86] sns.scatterplot(x=df['Weight'], y=df['Price'])

✓ 0s completed at 6:04 PM

Q

[87] df.corr()['Price']

Price	0.742905
Weight	0.209867
Price	1.000000
Touchscreen	0.192017
Ips	0.253328
ppi	0.475368
HDD	-0.096891
SSD	0.270658

✓ 0s completed at 6:04 PM

Q

[88] sns.heatmap(df.corr())

✓ 0s completed at 6:04 PM

Q

[89] sns.distplot(np.log(df['Price']))

✓ 0s completed at 6:04 PM

Q

[90] X = df.drop(columns=['Price'])  
y = np.log(df['Price'])

✓ 0s completed at 6:04 PM

Q

[91] X

	Company	TypeName	Ram	Weight	Touchscreen	Ips	ppi	Cpu brand	HDD	SSD	Gpu brand	os
0	Apple	Ultrabook	8	1.37	0	1	226.983005	Intel Core i5	0	128	Intel	Mac
1	Apple	Ultrabook	8	1.34	0	0	127.677940	Intel Core i5	0	0	Intel	Mac
2	HP	Notebook	8	1.86	0	0	141.211998	Intel Core i5	0	256	Intel	Others/No OS/Linux
3	Apple	Ultrabook	16	1.83	0	1	220.534624	Intel Core i7	0	512	AMD	Mac
4	Apple	Ultrabook	8	1.37	0	1	226.983005	Intel Core i5	0	256	Intel	Mac

✓ 0s completed at 6:04 PM

Q

[92] y

1302 rows x 12 columns

✓ 0s completed at 6:04 PM

```
+ Code + Text
RAM 8 GB
Disk 100 GB

[4] y
0 11.175755
1 10.776777
2 10.29931
3 11.814476
4 11.47381
...
1298 10.433899
1299 11.288115
1300 9.489283
1301 10.614219
1302 9.886358
Name: Price, Length: 1302, dtype: float64

[93] from sklearn.model_selection import train_test_split
X_train,X_test,y_train,y_test = train_test_split(X,y,test_size=0.15,random_state=2)

[94] X_train
Company Typeplane Ram Height Touchscreen Ips ppi Cpu brand HDD SSD Gpu brand os
183 Toshiba Notebook 8 2.00 0 0 100.454670 Intel Core i5 0 128 Intel Windows
1141 MSI Gaming 8 2.40 0 0 141.211998 Intel Core i7 1000 128 Nvidia Windows
1049 Asus Netbook 4 1.20 0 0 135.094211 Other Intel Processor 0 0 Intel Others/No OS/Linux
1020 Dell 2 in 1 Convertible 4 2.08 1 1 141.211998 Intel Core i3 1000 0 Intel Windows

+ Code + Text
[95] from sklearn.compose import ColumnTransformer
from sklearn.pipeline import Pipeline
from sklearn.preprocessing import OneHotEncoder
from sklearn.metrics import r2_score,mean_absolute_error

[96] from sklearn.linear_model import LinearRegression,Ridge,Lasso
from sklearn.neighbors import KNeighborsRegressor
from sklearn.tree import DecisionTreeRegressor
from sklearn.ensemble import RandomForestRegressor,GradientBoostingRegressor,AdaBoostRegressor,ExtraTreeRegressor
from sklearn.svm import SVC
from xgboost import XGBRegressor

Linear regression
[ ] cell hidden

Ridge Regression
[ ] step1 = ColumnTransformer(transformers=[
('col_tf',OneHotEncoder(sparse=False,drop='first'),[0,1,7,10,11])
],remainder='passthrough')
step2 = Ridge(alpha=10)
pipe = Pipeline([
('step1',step1),
('step2',step2)
])
pipe.fit(X_train,y_train)
y_pred = pipe.predict(X_test)
print('R2 score',r2_score(y_test,y_pred))
print('MAE',mean_absolute_error(y_test,y_pred))
R2 score 0.8073277468418649
MAE 0.2101782797424092

+ Code + Text
[98] step1 = ColumnTransformer(transformers=[
('col_tf',OneHotEncoder(sparse=False,drop='first'),[0,1,7,10,11])
],remainder='passthrough')
step2 = LinearRegression()
pipe = Pipeline([
('step1',step1),
('step2',step2)
])
pipe.fit(X_train,y_train)
y_pred = pipe.predict(X_test)
print('R2 score',r2_score(y_test,y_pred))
print('MAE',mean_absolute_error(y_test,y_pred))
R2 score 0.8073277468418649
MAE 0.2101782797424092

Ridge Regression
[ ] cell hidden

+ Code + Text
[99] step1 = ColumnTransformer(transformers=[
('col_tf',OneHotEncoder(sparse=False,drop='first'),[0,1,7,10,11])
],remainder='passthrough')
step2 = Ridge(alpha=10)
pipe = Pipeline([
('step1',step1),
('step2',step2)
])
pipe.fit(X_train,y_train)
y_pred = pipe.predict(X_test)
print('R2 score',r2_score(y_test,y_pred))
print('MAE',mean_absolute_error(y_test,y_pred))
R2 score 0.8127331031118009
MAE 0.2092680242582973

Lasso Regression
[ ] cell hidden

+ Code + Text
[100] step1 = ColumnTransformer(transformers=[
('col_tf',OneHotEncoder(sparse=False,drop='first'),[0,1,7,10,11])
],remainder='passthrough')
step2 = Lasso(alpha=0.001)
pipe = Pipeline([
('step1',step1),
('step2',step2)
])
pipe.fit(X_train,y_train)
y_pred = pipe.predict(X_test)
print('R2 score',r2_score(y_test,y_pred))
print('MAE',mean_absolute_error(y_test,y_pred))
R2 score 0.8071853945317105
MAE 0.2111431613472365

+ Code + Text
[101] from sklearn.neighbors import KNeighborsRegressor
step1 = ColumnTransformer(transformers=[
('col_tf',OneHotEncoder(sparse=False,drop='first'),[0,1,7,10,11])
],remainder='passthrough')
step2 = KNeighborsRegressor(n_neighbors=5)
pipe = Pipeline([
('step1',step1),
('step2',step2)
])
pipe.fit(X_train,y_train)
y_pred = pipe.predict(X_test)
print('R2 score',r2_score(y_test,y_pred))
print('MAE',mean_absolute_error(y_test,y_pred))
R2 score 0.802198468448555
MAE 0.191971672121116
```

```
+ Code + Text
RAM 8 GB
Disk 100 GB

Decision Tree
[ ] cell hidden

[ ] step1 = ColumnTransformer(transformers=[
('col_tf',OneHotEncoder(sparse=False,drop='first'),[0,1,7,10,11])
],remainder='passthrough')
step2 = DecisionTreeRegressor(max_depth=4)
pipe = Pipeline([
('step1',step1),
('step2',step2)
])
pipe.fit(X_train,y_train)
y_pred = pipe.predict(X_test)
print('R2 score',r2_score(y_test,y_pred))
print('MAE',mean_absolute_error(y_test,y_pred))
R2 score 0.842866862998522
MAE 0.18885641910288333

SVM
[ ] cell hidden

[ ] step1 = ColumnTransformer(transformers=[
('col_tf',OneHotEncoder(sparse=False,drop='first'),[0,1,7,10,11])
],remainder='passthrough')
step2 = SVC(kernel='rbf',C=10000,epsilon=0.1)
pipe = Pipeline([
('step1',step1),
('step2',step2)
])
pipe.fit(X_train,y_train)
y_pred = pipe.predict(X_test)
print('R2 score',r2_score(y_test,y_pred))
print('MAE',mean_absolute_error(y_test,y_pred))
R2 score 0.888118092274576
MAE 0.202305567294315

Random Forest
[ ] cell hidden

[ ] step1 = ColumnTransformer(transformers=[
('col_tf',OneHotEncoder(sparse=False,drop='first'),[0,1,7,10,11])
],remainder='passthrough')
step2 = RandomForestRegressor(n_estimators=100,
random_state=1,
max_samples=0.5,
max_features=0.75,
max_depth=15)
pipe = Pipeline([
('step1',step1),
('step2',step2)
])
pipe.fit(X_train,y_train)
y_pred = pipe.predict(X_test)
print('R2 score',r2_score(y_test,y_pred))
print('MAE',mean_absolute_error(y_test,y_pred))
R2 score 0.8871603178182488
MAE 0.20471912980481796

+ Code + Text
[102] cell hidden

AdaBoost
[ ] cell hidden

[ ] step1 = ColumnTransformer(transformers=[
('col_tf',OneHotEncoder(sparse=False,drop='first'),[0,1,7,10,11])
],remainder='passthrough')
step2 = AdaBoostRegressor(n_estimators=15,learning_rate=1.0)
pipe = Pipeline([
('step1',step1),
('step2',step2)
])
pipe.fit(X_train,y_train)
y_pred = pipe.predict(X_test)
print('R2 score',r2_score(y_test,y_pred))
print('MAE',mean_absolute_error(y_test,y_pred))
R2 score 0.7862862278529815
MAE 0.23471912980481796

+ Code + Text
[103] step1 = ColumnTransformer(transformers=[
('col_tf',OneHotEncoder(sparse=False,drop='first'),[0,1,7,10,11])
],remainder='passthrough')
step2 = XGBRegressor(n_estimators=45,max_depth=5,learning_rate=0.5)
pipe = Pipeline([
('step1',step1),
('step2',step2)
])
pipe.fit(X_train,y_train)
y_pred = pipe.predict(X_test)
print('R2 score',r2_score(y_test,y_pred))
print('MAE',mean_absolute_error(y_test,y_pred))
[09:26:43] WARNING: /workspace/src/objective/regression_obj_cu::reg_linear is now deprecated in favor of reg_squarederror.
R2 score 0.8817734558243
MAE 0.14682681912400914

Voting Regressor
[ ] cell hidden

[ ] from sklearn.ensemble import VotingRegressor,StackingRegressor
step1 = ColumnTransformer(transformers=[
('col_tf',OneHotEncoder(sparse=False,drop='first'),[0,1,7,10,11])
],remainder='passthrough')
estimators = [
('rf', RandomForestRegressor(n_estimators=50,random_state=1,max_samples=0.5,max_features=0.75,max_depth=15)),
('gbt', GradientBoostingRegressor(n_estimators=100,max_features=0.5)),
('xgb', XGBRegressor(n_estimators=15,learning_rate=0.1,max_depth=5))
]
step2 = StackingRegressor(estimators=estimators,final_estimator=Ridge(alpha=100))
pipe = Pipeline([
('step1',step1),
('step2',step2)
])
pipe.fit(X_train,y_train)
```



