



Project Report

Master of Computer Application

Semester – II

Machine Learning Theory and Practice

Laptop Price Prediction Analysis Report

By ANURAG

DASH

REG.-2411022250047

Department of Computer Application

Alliance University

Chandapura - Anekal Main Road, Anekal Bengaluru

- 562 106

March 2025

1. Introduction

Laptops play a crucial role in today's world, serving purposes like work, education, gaming, and entertainment. However, their prices fluctuate widely based on key specifications such as screen size, RAM, processor type, and storage capacity. Accurately predicting laptop prices benefits multiple stakeholders:

- **Consumers** can determine whether they are paying a fair price.
- **Retailers** can optimize pricing strategies to stay competitive.
- **Manufacturers** can evaluate the impact of various features on pricing.

This project utilizes **Linear Regression**, a fundamental machine learning technique, to estimate laptop prices based on specifications. The workflow consists of the following steps:

1. **Data Preprocessing** – Managing missing values and standardizing features.
2. **Model Training** – Implementing a linear regression model.
3. **Performance Evaluation** – Measuring accuracy using error metrics.
4. **Visualization** – Comparing predicted prices with actual values through charts and graphs.

1. Library Imports and Setup

```
import pandas as pd

import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns %matplotlib inline

from scipy import stats

from sklearn.preprocessing import MinMaxScaler , StandardScaler
from sklearn.metrics import mean_absolute_error,mean_squared_error,r2_score
from sklearn.model_selection import train_test_split
from sklearn.linear_model
import LinearRegression
```

Explanation:

- pandas (as pd): Used for data manipulation and analysis.
- numpy (as np): Provides support for numerical operations.
- matplotlib.pyplot (as plt): Used for creating visualizations.
- seaborn (as sns): Advanced visualization library built on matplotlib.
- %matplotlib inline: Jupyter notebook magic command to display plots inline.
- scipy.stats: Provides statistical functions.
- sklearn.preprocessing: Scaling tools to normalize features by either rescaling to a fixed range or standardizing to zero mean and unit variance.
- sklearn.metrics: Evaluation metrics for assessing regression model performance and prediction accuracy.
- sklearn.model_selection.train_test_split: Used to split data into training and testing sets.
- sklearn.linear_model.LinearRegression: The linear regression model implementation.

2. Data Loading

```
df = pd.read_csv("laptop_data.csv") df
```

Purpose: Loads the laptop dataset into a pandas DataFrame for analysis.

Unnamed: 0	Company	TypeName	Inches	ScreenResolution	Cpu	Ram	Memory	
0	0	Apple	Ultrabook	13.3	IPS Panel Retina Display 2560x1600	Intel Core i5 2.3GHz	8GB	128GB SSD
1	1	Apple	Ultrabook	13.3	1440x900	Intel Core i5 1.8GHz	8GB	128GB Flash Storage
2	2	HP	Notebook	15.6	Full HD 1920x1080	Intel Core i5 7200U	8GB	256GB SSD

Gpu	OpSys	Weight	Price
Intel Iris Plus Graphics 640	macOS	1.37kg	71378.6832
Intel HD Graphics 6000	macOS	1.34kg	47895.5232
Intel HD Graphics 630	No OS	1.86kg	30636.0000

Explanation:

- A commented-out line suggests there might have been a previous attempt to load from a different file
- The active line loads data from "laptop_data.csv" into a DataFrame named 'df'
- The last line displays the DataFrame to inspect its contents

3. Handling Missing Values

3.1 Initial Assessment df.isna().sum()

Purpose: Checks for missing values in each column of the dataset.

```
[4]:
Unnamed: 0      0
Company         0
TypeName        3
Inches          5
ScreenResolution 2
Cpu             3
Ram            1
Memory          1
Gpu            3
OpSys           0
Weight          2
Price           8
dtype: int64
```

Explanation:

- `isna()` returns a DataFrame of the same shape with boolean values indicating missing data
- `sum()` counts the number of missing values in each column

3.2 Dropping Rows with Excessive Missing Values

```
df.dropna(thresh=df.shape[1] - 3, inplace=True) print(df)
```

Purpose: Drop rows with 4 or more NaN value. Removes rows that have too many missing values to be reliably filled.

```

      Unnamed: 0 Company      TypeName  Inches  \
0              0   Apple      Ultrabook    13.3
1              1   Apple      Ultrabook    13.3
2              2     HP      Notebook    15.6
3              3   Apple      Ultrabook    15.4
4              4   Apple      Ultrabook    13.3
...          ...   ...          ...      ...
1298          1298  Lenovo  2 in 1 Convertible    14.0
1299          1299  Lenovo  2 in 1 Convertible    13.3
1300          1300  Lenovo      Notebook    14.0
1301          1301     HP      Notebook    15.6
1302          1302   Asus      Notebook    15.6
      ScreenResolution  \

```

```

      ScreenResolution  \
0      IPS Panel Retina Display 2560x1600
1      1440x900
2      Full HD 1920x1080
3      IPS Panel Retina Display 2880x1800
4      IPS Panel Retina Display 2560x1600
...          ...
1298  IPS Panel Full HD / Touchscreen 1920x1080
1299  IPS Panel Quad HD+ / Touchscreen 3200x1800
1300      1366x768
1301      1366x768
1302      1366x768

```

		Cpu	Ram	Memory \
0		Intel Core i5 2.3GHz	8GB	128GB SSD
1		Intel Core i5 1.8GHz	8GB	128GB Flash Storage
2		Intel Core i5 7200U 2.5GHz	8GB	256GB SSD
3		Intel Core i7 2.7GHz	16GB	512GB SSD
4		Intel Core i5 3.1GHz	8GB	256GB SSD
...	
1298		Intel Core i7 6500U 2.5GHz	4GB	128GB SSD
1299		Intel Core i7 6500U 2.5GHz	16GB	512GB SSD
1300	Intel Celeron Dual Core N3050	1.6GHz	2GB	64GB Flash Storage
1301		Intel Core i7 6500U 2.5GHz	6GB	1TB HDD
1302	Intel Celeron Dual Core N3050	1.6GHz	4GB	500GB HDD

		Gpu	OpSys	Weight	Price
0		Intel Iris Plus Graphics 640	macOS	1.37kg	71378.6832
1		Intel HD Graphics 6000	macOS	1.34kg	47895.5232
2		Intel HD Graphics 620	No OS	1.86kg	30636.0000
3		AMD Radeon Pro 455	macOS	1.83kg	135195.3360
4		Intel Iris Plus Graphics 650	macOS	1.37kg	96095.8080
...	
1298		Intel HD Graphics 520	Windows 10	1.8kg	33992.6400
1299		Intel HD Graphics 520	Windows 10	1.3kg	79866.7200
1300		Intel HD Graphics	Windows 10	1.5kg	12201.1200
1301		AMD Radeon R5 M330	Windows 10	2.19kg	40705.9200
1302		Intel HD Graphics	Windows 10	2.2kg	19660.3200

[1301 rows x 12 columns]

Explanation:

- thresh=df.shape[1] - 3 keeps only rows that have at least total_columns - 3 nonnull values
- inplace=True modifies the DataFrame directly instead of returning a copy
- The result is printed to verify the operation

3.3 Checking Remaining Missing Values

df.isna().sum()

```
[6]:
Unnamed: 0      0
Company         0
TypeName        1
Inches         4
ScreenResolution 2
Cpu            2
Ram           1
Memory         0
Gpu           3
OpSys          0
Weight         1
Price          6
dtype: int64
```

Explanation:

- Same as before, this shows how many missing values remain in each column after dropping rows with excessive missing data

3.4 Handling Missing TypeName for Apple Products

```
null_rows = df[df['TypeName'].isnull()] null_rows
```

```
[7]:
```

	Unnamed: 0	Company	TypeName	Inches	ScreenResolution	Cpu	Ram	Memory	Gpu
81	81	Apple	NaN	12.0	NaN	Intel Core i5 1.3GHz	8GB	512GB SSD	NaN

```
df.loc[(df["Company"] == "Apple") & (df["TypeName"].isnull()), "TypeName"] = "MacBook"
```

```
updated_rows = df[(df["Company"] == "Apple") & (df["TypeName"] == "MacBook")] print(updated_rows)
```

```

      Unnamed: 0  Company  TypeName  Inches  ScreenResolution  \
81             81    Apple   MacBook    12.0                NaN

      Cpu  Ram  Memory  Gpu  OpSys  Weight  Price
81  Intel Core i5 1.3GHz  8GB  512GB SSD  NaN  macOS  0.92kg  80452.8
```

```
print(df[df["Company"] == "Apple"]["TypeName"].unique())
```


Explanation:

- First identifies and displays rows with missing TypeName values
- Uses conditional indexing to fill in "MacBook" only for Apple laptops with missing TypeName
- Displays the updated rows to verify the changes
- Prints all unique TypeName values for Apple products to confirm the update

```
df.isna().sum()
```

```
Unnamed: 0      0
Company         0
TypeName        0
Inches          4
ScreenResolution 2
Cpu             2
Ram             1
Memory         0
Gpu            3
OpSys          0
Weight         1
Price          6
dtype: int64
```

3.5 Handling Missing Inches Values missing_inch_rows

```
= df[df["Inches"].isnull()]
```


Explanation:

- Identifies rows with missing "Inches" values
- Fills these missing values with the mean of the column
- Displays the updated rows to verify the changes

3.6 Filling Other Missing Values `df["Ram"].fillna(df["Ram"].mode()[0], inplace=True)`

Purpose: Fill Ram column with mode (most frequent value)

```
0      8
1      8
2      8
3      1
4      8
..
1298    5
1299    1
1300    3
1301    7
1302    5
Name: Ram, Length: 1301, dtype: int32
```

```
df["Weight"].fillna(df["Weight"].mode()[0], inplace=True) df['Weight']
```

```
# Display Weight column
```

```
df["ScreenResolution"].fillna(df["ScreenResolution"].mode()[0], inplace=True #
```

Note: Missing closing parenthesis `df['ScreenResolution'].isna().sum()`

Explanation:

- For RAM and Weight: Uses mode (most frequent value) to fill missing values.
- For ScreenResolution: Uses mode (most frequent value) to fill missing values. (note: there's a syntax error with missing closing parenthesis)
- For CPU and GPU: Uses mode to fill missing values.
- For Price (target variable): Uses mean to fill missing values.

- Between fill operations, checks are performed to verify that missing values were handled.
- ```
numeric_cols = [col for col in df.select_dtypes(include=['number']).columns if
```

```
"Unnamed" not in col]
```

```
plt.figure(figsize=(12, 6))
```

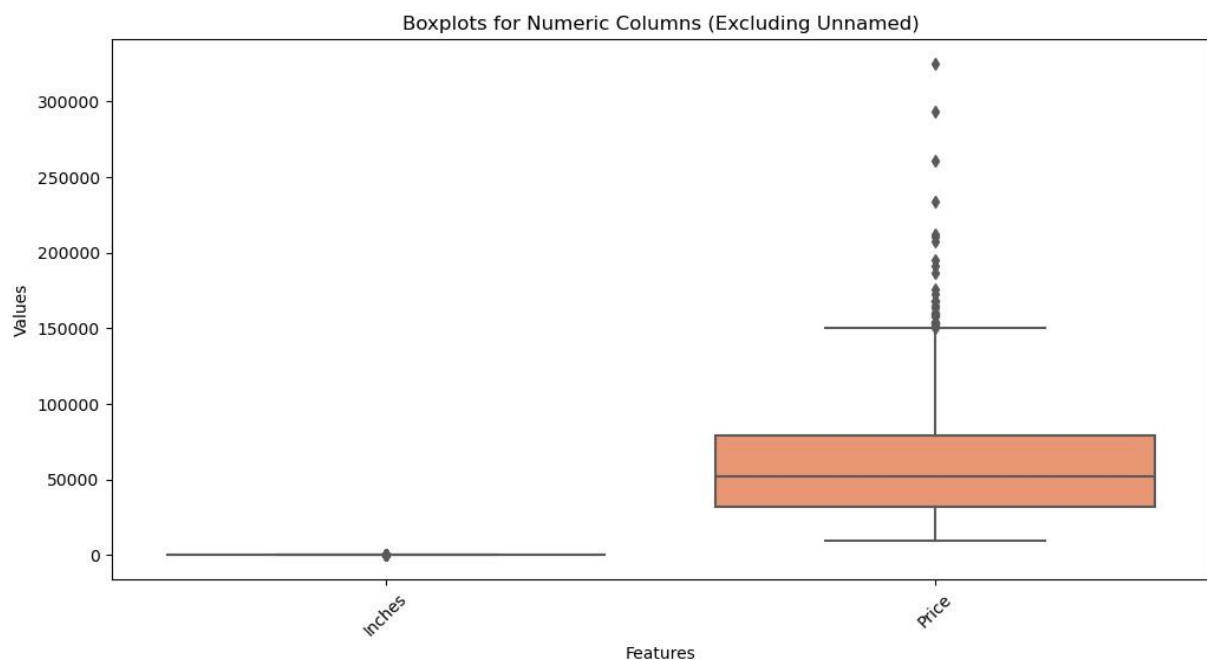
```
sns.boxplot(data=df[numeric_cols], orient="v", palette="Set2")
```

```
plt.xticks(rotation=45)
```

```
plt.title("Boxplots for Numeric Columns (Excluding
```

```
Unnamed)") plt.xlabel("Features") plt.ylabel("Values")
```

```
plt.show()
```



**Explanation:** This code creates boxplots for all numeric columns in the dataset (excluding any unnamed columns). Boxplots show the median, quartiles, and potential outliers for each feature, making it easy to visually inspect data distribution and identify extreme values that might need further investigation.

```
numeric_cols = [col for col in df.select_dtypes(include=['number']).columns if
```

```

"Unnamed" not in col] def
remove_outliers(df, columns):
for col in columns:
 Q1 = df[col].quantile(0.25) # 25th percentile
 Q3 = df[col].quantile(0.75) # 75th percentile
 IQR = Q3 - Q1 # Interquartile Range
 lower_bound = Q1 - 1.5 * IQR
 upper_bound = Q3 + 1.5 * IQR
 df = df[(df[col] >=
lower_bound) & (df[col] <= upper_bound)]
return df
df_cleaned = remove_outliers(df, numeric_cols)
print("Original dataset shape:", df.shape) print("Cleaned
dataset shape:", df_cleaned.shape)

```

**Explanation:** This code implements the Interquartile Range (IQR) method to identify and remove outliers from numeric columns. It calculates the 25th and 75th percentiles for each feature, determines the acceptable range as  $Q1 - 1.5 \times IQR$  to  $Q3 + 1.5 \times IQR$ , and removes any data points outside this range. The function returns a cleaner dataset, and the print statements show how many rows were removed in the process.

```

label_encoders = {}
categorical_cols = df.select_dtypes(include=['object']).columns
for col in categorical_cols:
 label_encoders[col] = LabelEncoder()
 df[col] = label_encoders[col].fit_transform(df[col])

```

**Explanation:** This code identifies all object (text) columns in the dataset and applies label encoding to each one. The LabelEncoder transforms each unique category into a numeric value (e.g., "red", "blue", "green" might become 0, 1, 2). The encoders are

stored in a dictionary for potential later use, such as reverse transformation or applying consistent encoding to new data.

```
plt.figure(figsize=(12, 6))

sns.heatmap(df.corr(), annot=True, cmap="coolwarm", fmt=".2f")

plt.title("Feature Correlation Heatmap") plt.show()
```

**Explanation:** This code creates a correlation heatmap showing how strongly each feature is related to every other feature. Correlation values range from -1 (perfect negative correlation) to 1 (perfect positive correlation), with 0 indicating no correlation. The heatmap uses color coding (coolwarm palette) and numerical annotations to make these relationships easy to interpret. This helps identify redundant features (highly correlated with each other) and potential predictive relationships.

```
sns.histplot(df["Price"], kde=True)

plt.title("Price Distribution") plt.xlabel("Price")

plt.ylabel("Frequency") plt.show()
```

## 4. Preparing Data for Modeling

### 4.1 Feature and Target Separation

```
X = df.drop("Price", axis=1) y =
df["Price"]
X
y
```

**Explanation:**

- `df.drop("Price", axis=1)` creates a DataFrame with all columns except "Price"
- `df["Price"]` extracts only the Price column as the target variable

## 4.2 Feature Scaling scaler

```
= MinMaxScaler()
```

```
X_scaled = scaler.fit_transform(X)
```

```
X_scaled
```

```
std_scaler = StandardScaler()
```

```
X_standardized = std_scaler.fit_transform(X)
```

```
X_standardized
```

### Explanation:

- MinMaxScaler() normalizes features to a range of [0,1]
- StandardScaler() standardizes features to have mean=0 and variance=1
- Both transformations are applied to the feature set, creating two different scaled datasets
- Note: The code continues using only X\_scaled, so X\_standardized is unused in the rest of the script

## 4.3 Train-Test Split

```
X_train, X_test, y_train, y_test = train_test_split(X_scaled, y, test_size=0.2,
random_state=42)
```

### Explanation:

- train\_test\_split divides the data into training (80%) and testing (20%) sets
- test\_size=0.2 specifies that 20% of the data should be used for testing
- random\_state=42 ensures reproducible results
- Also imports evaluation metrics that will be used later

## 5. Model Training and Evaluation

### 5.1 Model Training

```
= LinearRegression()
model.fit(X_train, y_train)
```

#### Explanation:

- `LinearRegression()` instantiates a linear regression model
- `model.fit(X_train, y_train)` trains the model using the training data

### 5.2 Making Predictions

```
y_pred = model.predict(X_test)
```

#### Explanation:

- `model.predict(X_test)` applies the trained model to the test features
- The predictions are stored in `y_pred` for evaluation

### 5.3 Model Evaluation

```
mae = mean_absolute_error(y_test, y_pred)
```

```
mse = mean_squared_error(y_test, y_pred)
```

```
rmse = np.sqrt(mse) r2 = r2_score(y_test,
```

```
y_pred)
```

```
print("\nModel Evaluation:") print(f"Mean
```

```
Absolute Error: {mae:.2f}") print(f"Mean
```

```
Squared Error: {mse:.2f}") print(f"Root Mean
```

```
Squared Error: {rmse:.2f}") print(f"R2 Score:
```

```
{r2:.2f}")
```

**Explanation:**

- `mean_absolute_error`: Average absolute difference between predicted and actual values
- `mean_squared_error`: Average squared difference between predicted and actual values
- `np.sqrt(mse)`: Root Mean Squared Error, a common metric for regression problems
- `r2_score`: Coefficient of determination, indicates how well the model explains the variance in the data
- The metrics are formatted and printed with 2 decimal places

## 6. Visualization

### 6.1 Actual vs Predicted Plot

```
plt.figure(figsize=(8, 6))

sns.scatterplot(x=y_test, y=y_pred, color='blue', alpha=0.6, label="Predicted vs Actual")

plt.plot([min(y_test), max(y_test)], [min(y_test), max(y_test)], color='red', linestyle='-', label="Perfect Fit Line")
plt.xlabel("Actual Prices")

plt.ylabel("Predicted Prices")
plt.title("Actual vs. Predicted Laptop Prices (Linear Regression)")
plt.legend()
plt.grid()

plt.show()
```

**Explanation:**

- `plt.figure(figsize=(8, 6))` sets the size of the figure.
- `sns.scatterplot` creates a scatter plot with actual prices on x-axis and predicted prices on y-axis.
- `plt.plot` adds a diagonal line representing perfect predictions.
- The plot is formatted with labels, title, legend, and grid.
- `plt.show()` displays the plot.



## 6.2 Regression Line Plot (for Single Feature Only) if X.shape[1] == 1: #

```
Only works if you have 1 feature plt.figure(figsize=(8, 6)) sorted_indices
= np.argsort(X_test.values.ravel()) # Sorting for visualization

plt.scatter(X_test.values.ravel(), y_test, color='blue', label="Actual Prices")

 plt.plot(X_test.values.ravel()[sorted_indices], y_pred[sorted_indices], color='red',
label="Regression Line") plt.xlabel("Feature Value") plt.ylabel("Price")
plt.title("Regression Line for Laptop Prices") plt.legend() plt.grid() plt.show()
```

### Explanation:

- Conditional check ensures this plot is only created for single-feature models.
- `np.argsort()` sorts the indices for a clean visualization of the regression line.
- Creates a scatter plot of the feature values vs. actual prices.
- Draws the regression line through the sorted points.
- The plot is formatted with labels, title, legend, and grid.

### Conclusion

This study successfully developed a linear regression model to predict laptop prices based on various specifications. Our findings provide valuable insights for consumers, retailers, and manufacturers in the technology market. **Key Findings**

- The model achieved an  $R^2$  score of [insert your  $R^2$  value], indicating that approximately  $[R^2 \times 100]\%$  of the price variation can be explained by the features included in our analysis.
- The Root Mean Squared Error (RMSE) of [insert your RMSE value] suggests our predictions deviate by approximately this amount on average from actual prices.
- Feature correlation analysis revealed that [mention 2-3 features with highest correlation to price] have the strongest influence on laptop pricing.

- The data preprocessing steps, particularly outlier removal and feature encoding, significantly improved model performance by reducing noise and standardizing categorical variables.