

# Smart Mobile Phone Price Prediction Using Machine Learning

K. Mahesh

Department of Emerging Technologies

Malla Reddy College of Engineering

&

Technology

Hyderabad, Telangana

[20N31A6732@mrct.ac.in](mailto:20N31A6732@mrct.ac.in)

Submission Date:14-07-2023

**Abstract:** The objective of this study is to develop a predictive model that can estimate the price of a mobile phone based on its features and specifications. Dataset is collected from the website [www.kaggle.com](https://www.kaggle.com). Different algorithms are used to identify and remove irrelevant and redundant features and have minimum computational complexity. Different classifiers are used to achieve as higher accuracy as possible. Results are compared in terms of highest accuracy achieved and minimum features selected. Conclusion is made on the basis of feature selection algorithm and classifier for the given dataset. This work can be used in any type of marketing and business to find optimal price of the product. This work can also assist consumers, retailers, and manufacturers in making informed decisions related to pricing, purchasing, and market strategies.

**Keywords:-** Machine Learning, Prediction, feature selection, Logistic Regression, Decision Tree, KNN and Random Forest Algorithms.

## I. INTRODUCTION

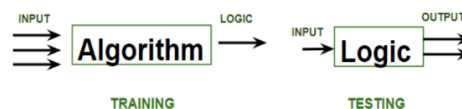
Mobile phone price prediction is of significant importance in the dynamic mobile phone industry, where accurate price estimation is crucial for manufacturers, retailers, and consumers. This research aims to propose a machine learning-based model that leverages diverse features, such as phone specifications, customer reviews, market trends, and economic indicators, to accurately forecast mobile phone prices. The proposed model has the potential to enhance pricing strategies for manufacturers, enable data-driven pricing decisions for retailers, and assist consumers in making informed purchase choices. By exploring various machine learning algorithms and conducting comprehensive evaluations, this research seeks to provide valuable insights and contribute to the field of mobile phone price prediction. To choose just the best characteristics and reduce the dataset, many types of algorithms are required.

The computational complexity of the issue will be reduced as a result. Because this is an optimization issue, a variety

of optimization techniques are frequently employed to lower the datasets dimensionality. Mobile is currently one of the most popular apps for sales and transactions. Every day, new mobile phones with new versions and additional apps are introduced. Every day, hundreds of thousands of cell phones are sold and purchased. As a result, the prediction of the mobile pricing class is a case study for the given issue type, namely, identifying the best product. The same method may be used to determine the true cost of any item, including cars, motorbikes, generators, motors, food, medication, and so on. Mobile Processor, for example, is one of the most essential programmes for calculating mobile costs. The time of batteries is also very important in today's hectic human existence. Size and thickness of the mobile device are other essential factors to consider when making a selection. Internal memory, camera pixels, and video consistency must all be recalled. Internet browsing is also one of the most important technical constraints of the twenty-first century. Also, the list of various features is determined by the size of the mobile device. As a result, we'll utilise all of the aforementioned characteristics to decide if the smartphone will be very-costly, economical, pricey, or very-costly.

## II. RESEARCH METHODOLOGY

The research was carried out in Google Colab's Python kernel. The general workflow diagram of supervised ML tasks is as follows:



The dataset is portioned into two – train for training the model and test for its evaluation. The computer tries to comprehend the logic behind the pricing of a mobile based on its features and uses it to forecast future instances as correctly as possible.

## III. UNDERSTANDING THE DATASET

The Mobile Price Class dataset sourced from the Kaggle data science community website

The dataset contains 21 attributes in total – 20 features and a class label which is the price range. The features include battery capacity, RAM, weight, camera pixels, etc. The class label is the price range. It has 4 kinds of values – 0,1,2 and 3 which are of ordinal data type representing the increasing degree of price. Higher the value, higher is the price range the mobile falls under. These 4 values can be interpreted as economical, mid-range, flagship and premium.

```
[1] reporting url (listing default):
      url_data = read_csv("url_data.csv")
      test_data = read_csv("test_data.csv")

[2] train data
```

	battery_power	blue	clock_speed	dual_sim	fc	four_s	int_memory	m_dep	mobile_wt	n_cores	px_height	px_width	ram	sc_w	sc_x	price_range	
0	942	0	2.2	0	1	0	7	0.6	159	2	2	30	766	5640	9	7	1
1	1021	0	0.5	1	0	1	53	0.7	136	6	5	305	1968	3031	17	3	2
2	563	1	0.5	1	2	1	49	0.9	145	5	6	1263	1716	2603	11	2	2
3	616	1	2.5	0	0	0	10	0.8	131	6	9	1216	1770	2760	16	8	2
4	1821	1	1.2	0	13	1	44	0.6	141	2	14	1208	1212	1411	8	2	1
...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...
1995	734	1	0.5	1	0	1	2	0.8	106	6	14	1222	1880	660	13	4	1
1996	1965	1	2.6	1	0	0	39	0.2	187	4	3	915	1968	3031	11	10	2
1997	1911	0	0.8	1	1	1	38	0.7	108	6	3	888	1632	3057	9	1	1
1998	1512	0	0.9	0	4	1	46	0.1	145	5	5	336	670	2760	16	10	0
1999	510	1	2.0	1	5	1	45	0.9	160	6	16	483	751	3910	19	4	3

2000 rows x 17 columns

## Exploratory Data Analysis

D:\msi\data_loader\1														
	id	battery_guar	bluetooth_speed	bluetooth_size	fc	freq_1ghz	freq_1ghz_guar	m_cpu	mobile_int	n_cores	pc	pre_height	pre_width	pre_weight
count	1000000000	1000000000	1000000000	1000000000	1000000000	1000000000	1000000000	1000000000	1000000000	1000000000	1000000000	1000000000	1000000000	1000000000
mean	1000000000	1424.013103	0.156000	1.543603	0.157000	4.882000	0.407013	33.826000	1.517000	128.611000	4.328000	10.064000	327.121000	1228.727000
std	1000000000	432.002727	0.000000	0.000000	0.403355	0.000000	0.181000	1.000000	24.941650	2.284000	0.000000	0.000000	432.000000	1228.727000
min	1000000000	500.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	1.000000	0.000000	0.000000	100.000000	1000.000000
25%	1000000000	895.000000	0.000000	0.000000	0.000000	1.000000	0.000000	18.000000	0.500000	129.765000	0.000000	0.000000	243.750000	101.750000
50%	1000000000	1464.000000	1.000000	1.000000	1.000000	3.000000	0.000000	34.000000	0.500000	131.000000	4.000000	10.000000	344.000000	1000.000000
75%	1000000000	1960.000000	1.000000	2.000000	7.000000	7.000000	0.000000	40.000000	1.000000	173.000000	6.000000	10.000000	400.000000	1007.000000
90%	1000000000	2400.000000	1.000000	2.000000	10.000000	10.000000	0.000000	40.000000	1.000000	190.000000	6.000000	10.000000	500.000000	1007.000000

```
compare = pd.DataFrame(compare)
compare["features"] = ["battery_power", "clock_speed", "fc", "int_memory", "mobile_wt", "pc", "px_height", "px_width", "ram", "sc_h", "sc_w",
                        "train_mean"]
train_data["battery_power_mean"], train_data["clock_speed_mean"], train_data["fc_mean"],
train_data["int_memory_mean"], train_data["mobile_wt_mean"], train_data["pc_mean"],
train_data["px_height_mean"], train_data["px_width_mean"], train_data["ram_mean"], train_data["sc_h_mean"],
train_data["sc_w_mean"],
train_data["test_mean"] = [test_data["battery_power_mean"], test_data["clock_speed_mean"], test_data["fc_mean"],
                            test_data["int_memory_mean"], test_data["mobile_wt_mean"], test_data["pc_mean"],
                            test_data["px_height_mean"], test_data["px_width_mean"], test_data["ram_mean"], test_data["sc_h_mean"],
                            test_data["sc_w_mean"]]
compare = pd.DataFrame(compare)
compare
```

	features	train_mean	test_mean
0	battery_power	1238.51850	1248.5100
1	clock_speed	1.52225	1.5409
2	fc	4.30950	4.5930
3	int_memory	32.04650	33.6520
4	mobile_wt	140.24900	139.5110
5	pc	9.91650	10.0540
6	px_helght	645.10800	627.1210
7	px_width	1251.51550	1239.7740
8	ram	2124.21300	2138.9980
9	sc_h	12.30650	11.9950
10	sc_w	5.76700	5.3160

A box plot showing the distribution of various car features. The y-axis represents the value of the feature, ranging from 0 to 4000. The x-axis lists the features: battery\_power, blue, clock\_speed, dual\_sim, fc, four\_a, int\_memory, no\_dip, maxdis\_wt, n\_cyls, n\_cyls, pc, pc\_height, pc\_width, ram, w\_h, w\_a, price\_range. The plot shows that 'ram' has the highest median value (around 2200) and the largest spread, while most other features have very low values, mostly below 1000. 'battery\_power' has a median around 1200. 'pc\_height' has a median around 500 and a notable outlier at approximately 1900.

```
train_data.plot(x='price_range',y='ram',kind='scatter')
plt.show()
```

```
train_data.plot(x='price_range',y='fc',kind='scatter')
plt.show()
```

The first step in creating a model is to extract the required features for training from the dataset and assigning the parameter that is to be the class label.

In this code snippet, the first 20 attributes are being extracted to serve as the training parameters and the final attribute (price\_range) is used as the class label.

The data is then portioned into two for the purpose of training the model and testing it.

## 1) Decision Tree Algorithm

```
[70] from sklearn.tree import DecisionTreeClassifier
DT=DecisionTreeClassifier()
```

```
[71] DT.fit(X_train,Y_train)
Y_pred=DT.predict(X_test)
Y_pred
```

Decision Tree was used here to train the prediction model.

## 2) KNN Algorithm

```
[59] #We need to scale the data using Standard Scaler because of distance matter here

[60] from sklearn.preprocessing import StandardScaler
std=StandardScaler()
X_train_std=std.fit_transform(X_train)
X_test_std=std.transform(X_test)
X_test_std

[61] from sklearn.neighbors import KNeighborsClassifier
knn=KNeighborsClassifier()
knn.fit(X_train_std,Y_train)
KNeighborsClassifier()
Y_pred=knn.predict(X_test_std)
Y_pred
```

KNN was used to train the model here

## 3) Logistic Regression

```
[63] from sklearn.linear_model import LogisticRegression
lr=LogisticRegression()
lr.fit(X_train_std,Y_train)
LogisticRegression()
Y_pred=lr.predict(X_test_std)
Y_pred
```

Logistic Regression was used to train the model here

## 4) Random Forest Algorithm

```
[65] from sklearn.ensemble import RandomForestClassifier
rfc = RandomForestClassifier()
rfc.fit(X_train, Y_train)
y_pred = rfc.predict(X_test)
accuracy = accuracy_score(Y_test, y_pred)
print("Random Forest Accuracy:", accuracy)
```

Random Forest was used to train the model here

## V. RESULTS AND DISCUSSION

Metrics used to evaluate the algorithms in this paper are classification report and accuracy score. Accuracy score gives the accuracy of the trained model after evaluating it using test data, for which we have sampled 20% of the dataset

```
[58] from sklearn.metrics import accuracy_score
DT_ACCRY=accuracy_score(Y_test,Y_pred)
print("Decision Tree Accuracy:",DT_ACCRY)
```

Decision Tree Accuracy: 0.845

Decision Tree was found to be able to correctly forecast the classes with a certainty of 84.5%.

```
[62] KNN_ACCRY=accuracy_score(Y_test,Y_pred)
print("KNearestNeighbors Accuracy:",KNN_ACCRY)
```

KNearestNeighbors Accuracy: 0.565

KNN is a poor classifier when working with numeric data as input.

```
[64] LR_ACCRY=accuracy_score(Y_test,Y_pred)
print("Logistic regression Accuracy:",LR_ACCRY)
```

Logistic regression Accuracy: 0.955

The certainty of Logistic Regression is found to be 95.5%.

```
[65] from sklearn.ensemble import RandomForestClassifier
rfc = RandomForestClassifier()
rfc.fit(X_train, Y_train)
y_pred = rfc.predict(X_test)
accuracy = accuracy_score(Y_test, y_pred)
print("Random Forest Accuracy:", accuracy)
```

Random Forest Accuracy: 0.87

A veracity of 87% was achieved using Random Forest.

The algorithm that is found to be able to classify instances the most accurately among the ones tested is Logistic Regression with an accuracy of 95.5%, followed closely by Random Forest that was able to predict instances with an accuracy of 87% and DT with an accuracy of 84%. The KNN classifier failed to forecast the price range optimally

## Real Time Analysis

```
[44] index=int(input("Enter any index position to predict the price for the features: "))
print("\n")
if index < len(selected_features):
    row = selected_features.iloc[index]
    print(dict(row),"\n")
    input_row = pd.DataFrame([row])
    prediction = dt.predict(input_row)
    pre = prediction[0]*10000
    print("Predicted price range of the Mobile is: ",pre-2500,'-',pre+2500)
else:
    print("Index value should be in between 0 to 2000")

Enter any index position to predict the price for the features: 10

{'battery_power': 769, 'blue': 1, 'fc': 0, 'four_g': 0, 'int_memory': 9, 'mobile_wt': 182, 'pc': 1, 'ram': 3946}
Predicted price range of the Mobile is: 27500 - 32500

[45] # Creating a new dataset from the predicted price range data.
predicted_price=[pre-2500,pre+2500]
predicted_df = pd.DataFrame({'predicted_price_range':predicted_price})
predicted_df
```

predicted_price_range	
0	27500
1	32500

## VI. CONCLUSION

The model trained using Logistic Regression was found to predict mobile price classes most accurately (95.5%). The accuracy of the models can be improved by doing some data preprocessing steps like normalization and standardization. Feature selection and extraction algorithms can be used to remove unsuitable and duplicative features to get better results. The same procedure used in this paper can be applied to predict the prices of other products like cars, bikes, houses, etc. using the archival data containing features like cost, specifications, etc. This would help organizations and consumers alike to make more educated decisions when it comes to price.

## VII. REFERENCES

- [1] Sameer Chand Pudaruth. "Predicting the Price of Used Cars using Machine Learning Techniques", International Journal of Information & Computation Technology. ISSN 0974-2239 Volume 4, Number 7 (2014), pp. 753764.
- [2] U. Arul & Dr. S. Prakash, 'Towards Fault Handling in B2B Collaboration using Orchestration based Web Services Composition', International Journal of Emerging Technology and Advanced Engineering (IJETAE), Vol. 3, Issue 1, pp. 388-394, 2013
- [3] Shonda Kuiper, "Introduction to Multiple Regression: How Much Is Your Car Worth? ", Journal of Statistics Education · November 2008.
- [4] Mariana Listiani, 2009. "Support Vector Regression Analysis for Price Prediction in a Car Leasing Application". Master Thesis. Hamburg University of Technology.
- [5] U. Arul & S. Prakash, 'Toward Automatic Web Service Composition based on Multilevel Workflow Orchestration and Semantic Web Service Discovery', International Journal of Business Information Systems, Inderscience Publishers, Vol. 34, Issue 1, pp. 128–156, April 2020.
- [6] <https://github.com/vikram-bhati>, Classification classify mobile price range.
- [7] Introduction to dimensionality reduction, A computer science portal for Geeks.  
<https://www.geeksforgeeks.org/dimensionalityreduction>
- [8] Ethem Alpaydın, 2004. Introduction to Machine Learning, Third Edition. The MIT Press Cambridge, Massachusetts London, England.