

# Mobile Price Classification

---

## **Problem Statement**

**1. Introduction:** In today's digital era, smartphones have become an essential part of our lives. With a wide range of options available in the market, it can be challenging for consumers to choose the right mobile phone that suits their needs and budget. The Mobile Price Classification problem aims to predict the price range of a mobile phone based on its features. By accurately classifying mobile phones into different price ranges, consumers can make informed decisions while purchasing a new mobile device.

**2. Dataset Information:** The dataset used for the Mobile Price Classification problem consists of various features of mobile phones, including both numerical and categorical variables. Each data instance represents a mobile phone, and the target variable is the price range, which is categorized into four classes: low cost, medium cost, high cost, and very high cost.

The features in the dataset include:

1. Battery Power: Total energy a battery can store in one time measured in mAh.
2. Bluetooth: Indicates whether the phone has Bluetooth functionality (0 for no, 1 for yes).
3. Clock Speed: Speed at which the microprocessor executes instructions.
4. Dual SIM: Indicates whether the phone has dual SIM support (0 for no, 1 for yes).
5. Front Camera: Mega pixels of the front camera.
6. 4G: Indicates whether the phone supports 4G connectivity (0 for no, 1 for yes).
7. Internal Memory: Internal memory capacity in gigabytes.
8. Mobile Depth: Depth of the mobile phone in centimeters.
9. Mobile Weight: Weight of the mobile phone.

10. Number of Cores: Number of cores in the processor.
11. Primary Camera: Mega pixels of the primary camera.
12. Pixel Resolution Height: Height of the pixel resolution.
13. Pixel Resolution Width: Width of the pixel resolution.
14. RAM: Random Access Memory capacity in mega bytes.
15. Screen Height: Screen height of the mobile phone in centimeters.
16. Screen Width: Screen width of the mobile phone in centimeters.
17. Talk Time: Longest time that a single battery charge will last.
18. 3G: Indicates whether the phone supports 3G connectivity (0 for no, 1 for yes).
19. Touch Screen: Indicates whether the phone has a touch screen (0 for no, 1 for yes).
20. Wi-Fi: Indicates whether the phone has Wi-Fi functionality (0 for no, 1 for yes).

The dataset is divided into two parts: a training set and a test set. The training set is used to train the machine learning models, while the test set is used to evaluate the performance of the models on unseen data.

**3. Background Information:** The Mobile Price Classification problem has significant real-world applications in the mobile phone industry. By accurately predicting the price range of a mobile phone based on its features, manufacturers, retailers, and consumers can benefit in the following ways:

- **Manufacturers:** Mobile phone manufacturers can use price range classification models to understand the market demand for different price segments. This information can help them design and produce mobile phones that cater to specific consumer preferences and price ranges.
- **Retailers:** Retailers can use price range classification models to recommend mobile phones to customers based on their budget and desired features. This can enhance the shopping experience for consumers and increase sales for retailers.
- **Consumers:** Consumers can use price range classification models to compare different mobile phone options and make informed decisions while purchasing a new device. They can choose a mobile phone that best meets their requirements and falls within their budget.

By building an accurate and reliable price range classification model, the Mobile Price Classification problem can contribute to improving the overall efficiency and effectiveness of the mobile phone market, benefiting both businesses and consumers.

**4. Objective:** The objective of the Mobile Price Classification problem is to develop a machine learning model that can accurately predict the price range of a mobile phone based on its features. The model should be able to classify mobile phones into the following price range categories: low cost, medium cost, high cost, and very high cost.

**To achieve this objective, the following steps will be performed:**

1. Data Preprocessing: The dataset will be preprocessed to handle missing values, handle categorical variables, and scale numerical features if required.
2. Exploratory Data Analysis (EDA): The dataset will be analyzed to gain insights into the distribution of features and their relationships with the target variable. This step will help in understanding the data and identifying any patterns or trends.
3. Model Training: Several machine learning algorithms will be trained on the preprocessed dataset. The models will be trained using the training set, and their performance will be evaluated using suitable evaluation metrics.
4. Model Evaluation: The trained models will be evaluated on the test set to measure their performance in predicting the price range of mobile phones. The evaluation metrics used may include accuracy, precision, recall, and F1 score.
5. Model Selection and Deployment: Based on the model evaluation results, the best-performing model will be selected and deployed for predicting the price range of new, unseen mobile phones.

By successfully building a mobile price classification model, this project aims to provide a useful tool for consumers, manufacturers, and retailers in making informed decisions regarding mobile phone purchases and market strategies.

## **Framework**

**Importing Libraries:** Begin by importing the necessary libraries for data manipulation, visualization, and machine learning. Commonly used libraries include Pandas, NumPy, Matplotlib, and scikit-learn.

**Loading the Dataset:** Load the dataset into your code. You can use the Pandas library to read the dataset from a CSV file or any other appropriate format. Store the dataset in a Pandas DataFrame for further processing.

**Data Preprocessing:** Perform data preprocessing steps to prepare the dataset for model training. This may include handling missing values, encoding categorical variables (if any), and scaling numerical features (if required). Use appropriate techniques such as imputation, one-hot encoding, and feature scaling.

**Exploratory Data Analysis (EDA):** Conduct exploratory data analysis to gain insights into the dataset. Visualize the distributions of features and their relationships with the target variable using plots, histograms, scatter plots, and other visualization techniques. This step will help you understand the data better and identify any patterns or trends.

**Feature Selection/Engineering:** Analyze the importance of features and perform feature selection or engineering if necessary. This step involves identifying the most relevant features for the classification task and selecting or creating new features that might improve the model's performance. Techniques such as correlation analysis, feature importance, and domain knowledge can be used.

**Data Splitting:** Split the dataset into training and testing sets. The training set will be used to train the machine learning model, while the testing set will be used to evaluate its performance on unseen data. Use scikit-learn's `train_test_split` function to split the dataset into the desired proportions (e.g., 80% training and 20% testing).

**Model Training:** Choose suitable machine learning algorithms for the Mobile Price Classification task and train them using the training dataset. Commonly used algorithms for classification include logistic regression, decision trees, random forests, support vector machines (SVM), and neural networks. Instantiate the models and fit them to the training data.

**Model Evaluation:** Evaluate the trained models using appropriate evaluation metrics. Calculate metrics such as accuracy, precision, recall, and F1 score to measure the performance of each model on the testing set. Compare the results to select the best-performing model.

**Hyperparameter Tuning:** Fine-tune the hyperparameters of the selected model to optimize its performance. Use techniques like grid search or random search to explore different combinations of hyperparameters and select the best ones. This step helps in improving the model's accuracy and generalization ability.

**Final Model Training:** Retrain the selected model using the entire dataset (both training and testing sets) with the optimized hyperparameters. This step allows the model to utilize the maximum available data for better performance.

**Model Deployment:** Once the final model is trained, save it for future use. You can serialize the model using libraries like joblib or pickle. This step prepares the model for deployment in a production environment or for making predictions on new, unseen data.

**Prediction and Results:** Use the deployed model to make predictions on new mobile phone instances or the entire dataset. Evaluate the predictions using appropriate metrics and analyze the results. This step helps in understanding the model's performance on real-world data and validating its effectiveness.

By following this detailed outline, you can write the code for the Mobile Price Classification project, from data preprocessing and model training to evaluation and deployment. Remember to document your code properly, including comments and explanations, for better understanding and maintainability.

## **Code Explanation**

The code you provided for the Mobile Price Classification project consists of several functions that perform different tasks. Let's understand each function and its purpose:

**load\_data:** This function is responsible for loading the dataset into the code. It takes the file path as an argument and uses the Pandas library to read the dataset from a CSV file. It returns a Pandas DataFrame containing the dataset.

**preprocess\_data:** The preprocess\_data function is responsible for data preprocessing tasks. It takes the dataset as input and performs several operations to prepare the data for model training. These operations include handling missing values, encoding categorical variables, and scaling numerical features. It returns the preprocessed data.

**perform\_eda:** The perform\_eda function conducts Exploratory Data Analysis (EDA) on the dataset. EDA helps us understand the data better and uncover any patterns or relationships. In this function, you can perform various analysis and visualization tasks such as plotting histograms, scatter plots, and bar charts to explore the distributions and relationships between features.

**feature\_selection:** The feature\_selection function is responsible for selecting the most relevant features for the classification task. Feature selection helps in reducing the dimensionality of the dataset and improving the model's performance. This function can use techniques such as correlation analysis or feature importance to identify the important features.

**split\_data:** The split\_data function is used to split the dataset into training and testing sets. This step is important to evaluate the model's performance on unseen data. The function takes the preprocessed dataset and the desired split ratio (e.g., 0.8 for an 80-20 split) as inputs and returns the training and testing sets.

**train\_model:** The train\_model function is responsible for training the machine learning model on the training dataset. It takes the training set and the chosen machine learning algorithm as inputs and fits the model to the data. This step involves finding patterns and relationships in the data to make predictions.

**evaluate\_model:** The evaluate\_model function evaluates the trained model's performance on the testing set. It takes the trained model and the testing set as inputs

and calculates various evaluation metrics such as accuracy, precision, recall, and F1 score. These metrics help us assess how well the model is performing.

**predict:** The predict function is used to make predictions on new, unseen data using the trained model. It takes the trained model and new data as inputs and returns the predicted mobile price labels. This function can be used to make predictions on individual instances or on the entire dataset.

**main:** The main function is the entry point of the code. It executes the main workflow of the project by calling the above functions in a specific order. It starts by loading the dataset, performs data preprocessing, conducts EDA, selects relevant features, splits the data into training and testing sets, trains the model, evaluates its performance, and finally makes predictions.

The workflow of the code follows a typical machine learning pipeline:

1. Load the dataset.
2. Preprocess the data by handling missing values, encoding categorical variables, and scaling numerical features.
3. Perform EDA to gain insights into the dataset.
4. Select relevant features to improve the model's performance.
5. Split the data into training and testing sets.
6. Train the chosen machine learning model using the training set.
7. Evaluate the model's performance on the testing set.
8. Use the trained model to make predictions on new, unseen data.
9. Analyze the results and assess the model's effectiveness.

By following this code and understanding its functions, you can build a Mobile Price Classification model from scratch. Remember to adapt the code to your specific requirements, such as selecting appropriate algorithms and evaluating metrics based on the problem at hand.

## **Future Work**

After successfully implementing the Mobile Price Classification project, there are several avenues for future work and improvements. Let's explore each step in detail and provide a step-by-step guide on how to implement these enhancements:

**1. Collecting More Data:** To improve the model's accuracy and generalization, it is beneficial to collect a larger and more diverse dataset. Gathering data from additional sources, such as different regions or mobile manufacturers, can provide a broader perspective and capture a wider range of mobile price variations.

**2. Exploring Advanced Feature Engineering:** Consider exploring advanced feature engineering techniques to enhance the model's predictive power. This can involve creating new features based on domain knowledge or using techniques such as feature interaction, polynomial features, or dimensionality reduction methods like Principal Component Analysis (PCA) to capture complex relationships in the data.

**3. Trying Different Machine Learning Algorithms:** Experiment with different machine learning algorithms to find the best fit for the Mobile Price Classification task. Start by exploring algorithms such as Random Forest, Gradient Boosting, Support Vector Machines (SVM), or even deep learning approaches like Neural Networks. Compare their performance and choose the algorithm that provides the highest accuracy or desired trade-off between accuracy and computational complexity.

**4. Hyperparameter Tuning:** Optimize the hyperparameters of the chosen machine learning algorithm to improve the model's performance. Hyperparameters are settings that control the learning process of the algorithm. Utilize techniques such as grid search, random search, or Bayesian optimization to find the best combination of hyperparameters that yields the highest accuracy on the validation set.

**5. Implementing Cross-Validation:** To obtain a more robust estimation of the model's performance, implement cross-validation. Cross-validation involves splitting the dataset into multiple subsets and iteratively training and evaluating the model on different combinations of these subsets. This approach provides a better understanding of the model's performance across different data samples and reduces the risk of overfitting.

**6. Handling Class Imbalance:** If the dataset exhibits class imbalance, where some classes have significantly fewer samples than others, consider employing techniques to



handle this issue. Techniques like oversampling the minority class, undersampling the majority class, or using advanced methods like SMOTE (Synthetic Minority Over-sampling Technique) can help balance the class distribution and prevent biased predictions.

**7. Model Ensemble and Stacking:** Explore the concept of model ensemble and stacking to improve the overall prediction performance. Ensemble methods combine predictions from multiple models to make a final prediction. Consider using techniques such as bagging, boosting, or stacking to leverage the strengths of different models and enhance the overall predictive power.

**8. Deployment and Real-Time Prediction:** Once the model is finalized, consider deploying it to make real-time predictions on new, unseen mobile data. This can involve building an application or API that takes mobile specifications as input and provides the predicted price range as output. Consider using frameworks like Flask or Django for building a web application or leveraging cloud services for scalable deployment.

**Step-by-Step Guide for Implementation: To implement the future work for the Mobile Price Classification project, follow these steps:**

1. Data Collection: Gather a larger and more diverse dataset containing mobile specifications and their corresponding prices.
2. Feature Engineering: Explore advanced feature engineering techniques to extract meaningful features from the dataset.
3. Algorithm Exploration: Experiment with different machine learning algorithms such as Random Forest, Gradient Boosting, SVM, or Neural Networks to find the best performing one.
4. Hyperparameter Tuning: Optimize the hyperparameters of the chosen algorithm using techniques like grid search or random search.
5. Cross-Validation: Implement cross-validation to assess the model's performance across different data subsets and reduce overfitting.
6. Class Imbalance Handling: Address class imbalance issues using techniques like oversampling, undersampling, or SMOTE.
7. Model Ensemble and Stacking: Explore ensemble methods like bagging, boosting, or stacking to improve prediction performance.
8. Deployment and Real-Time Prediction: Build an application or API to deploy the final model for real-time predictions on new mobile data.

By following these steps, you can enhance the Mobile Price Classification project and further improve the accuracy and usability of the model. Remember to adapt the steps to your specific requirements and explore additional techniques and algorithms as necessary. Happy exploring and building!

## **Concept Explanation**

Imagine you're on a treasure hunt in a magical forest. You have a map with clues, but it's written in a language only the forest creatures understand. Your goal is to reach the hidden treasure chest, but how do you decipher the clues and find the right path?

Well, let me introduce you to the fantastic Decision Tree, your loyal guide in this enchanted adventure! The Decision Tree is like a magical creature that can understand the clues and guide you to the treasure.

So, what exactly is a Decision Tree? Imagine it as a flowchart with branches and leaves. Each branch represents a decision or a clue, and each leaf represents an outcome or a final destination. It's like a giant "Choose Your Own Adventure" book!

To understand how a Decision Tree works, let's take a fun example. Suppose you want to decide which mobile phone to buy, and you've collected some data about different phones. You want to classify them into price ranges: "Low," "Medium," "High," or "Very High." Here's how the Decision Tree helps you make the right choice:

**Root Node** - The Starting Point: At the very top of the tree, we have the root node. It's like the starting point of your adventure. In our case, the root node might ask, "Does the phone have a dual camera?" This clue helps us make the first decision.

**Branches** - Decisions, Decisions: From the root node, branches grow like magical vines. Each branch represents a decision based on a clue. For example, if the phone has a dual camera, we might follow the branch that says "Yes," and if not, we take the branch that says "No."

**Intermediate Nodes** - More Clues: As we follow the branches, we encounter more clues in the form of intermediate nodes. These nodes help us make more decisions to get closer to the treasure. For instance, the intermediate node might ask, "Does the phone have a large screen size?"

**Leaves** - Finding the Treasure: Finally, at the end of each branch, we reach the leaves of the tree. These are like the treasure chests containing our classification outcomes. In our example, each leaf represents a price range: "Low," "Medium," "High," or "Very High." Based on the clues we followed, we discover the most suitable price range for the mobile phone.

The beauty of the Decision Tree is that it learns from the data you provide. It examines different clues and their outcomes, and it magically creates the branches and leaves to guide you. It tries to find the best clues that lead to accurate predictions.

But how does it learn? It's like the Decision Tree attends a magic school and takes classes on data analysis. It looks at the features of the phones, such as the camera, screen size, battery life, and more, to make the right decisions. It learns from examples you provide, where each example tells it about a phone's features and the corresponding price range.

Once the Decision Tree is trained, it becomes a wise and helpful companion. You can ask it about any phone's features, and it will use its knowledge to predict the price range. It's like having a magical oracle at your fingertips!

Remember, the Decision Tree is just one of the many enchanting algorithms in the magical forest of machine learning. It has its own strengths and weaknesses, and there are other creatures like Random Forests, Gradient Boosting, and Neural Networks that offer different powers.

So, next time you're on a treasure hunt for answers in a sea of data, remember the fantastic Decision Tree! It will guide you through the twists and turns of the clues, helping you find the hidden gems of knowledge. Happy adventuring, my friend!

## Exercise Questions

**Question: How does the Decision Tree algorithm handle missing values in the dataset?**

**Answer:** The Decision Tree algorithm handles missing values by employing a technique called "missing value imputation." There are several common strategies for imputing missing values in Decision Trees. One approach is to assign the most frequent value of the attribute to the missing values. Another approach is to use the average value of the attribute. Alternatively, a more sophisticated method is to use the values from similar instances based on other attributes. By imputing missing values, the Decision Tree algorithm ensures that all instances in the dataset can be properly classified.

**Question: What are some advantages of using the Decision Tree algorithm?**

**Answer:** The Decision Tree algorithm offers several advantages, including:

- Easy to understand and interpret: Decision Trees provide a clear and intuitive representation of the decision-making process, making it easy for humans to understand and interpret the results.
- Handles both categorical and numerical data: Decision Trees can handle both categorical and numerical attributes without requiring additional preprocessing steps.
- Can handle missing values: Decision Trees have built-in mechanisms to handle missing values by imputing them during the training process.
- Nonlinear relationships: Decision Trees can capture nonlinear relationships between features, making them suitable for complex datasets.
- Feature importance: Decision Trees can measure the importance of each feature in the classification process, providing valuable insights into the dataset.

**Question: How can you prevent overfitting when using Decision Trees?**

**Answer:** Overfitting occurs when a Decision Tree becomes overly complex and starts to memorize the training data instead of generalizing well to unseen data. To prevent overfitting, you can apply the following techniques:

- Pruning: Pruning involves trimming the Decision Tree by removing branches that do not contribute significantly to the overall performance. This helps simplify the model and reduce overfitting.
- Setting maximum depth: You can limit the maximum depth of the Decision Tree, controlling the number of levels it can grow. This prevents the tree from becoming too deep and overfitting the training data.
- Minimum sample split: By specifying a minimum number of samples required to split a node, you can prevent the Decision Tree from creating branches with too few instances, which may lead to overfitting.

**Question: Can Decision Trees handle regression problems in addition to classification?**

**Answer:** Yes, Decision Trees can handle both classification and regression problems. When used for regression, Decision Trees predict a continuous value instead of a class label. The algorithm splits the data based on different attributes and aims to minimize the variance within each split. The predicted value for a new instance is then calculated based on the average value of the instances falling into the leaf node.

**Question: What are some limitations of the Decision Tree algorithm?**

**Answer:** While Decision Trees offer many benefits, they also have certain limitations, including:

- Overfitting: Decision Trees can be prone to overfitting, especially if they are allowed to grow too deep or if the dataset is noisy.
- Instability: Decision Trees can be sensitive to small changes in the training data, which may lead to different trees being generated.
- Difficulty in handling continuous variables: Decision Trees work best with categorical or binary features. Handling continuous variables requires appropriate binning or other preprocessing techniques.
- Bias towards features with more levels: Decision Trees tend to favor attributes with more levels since they can provide more information gain, potentially ignoring other relevant features.