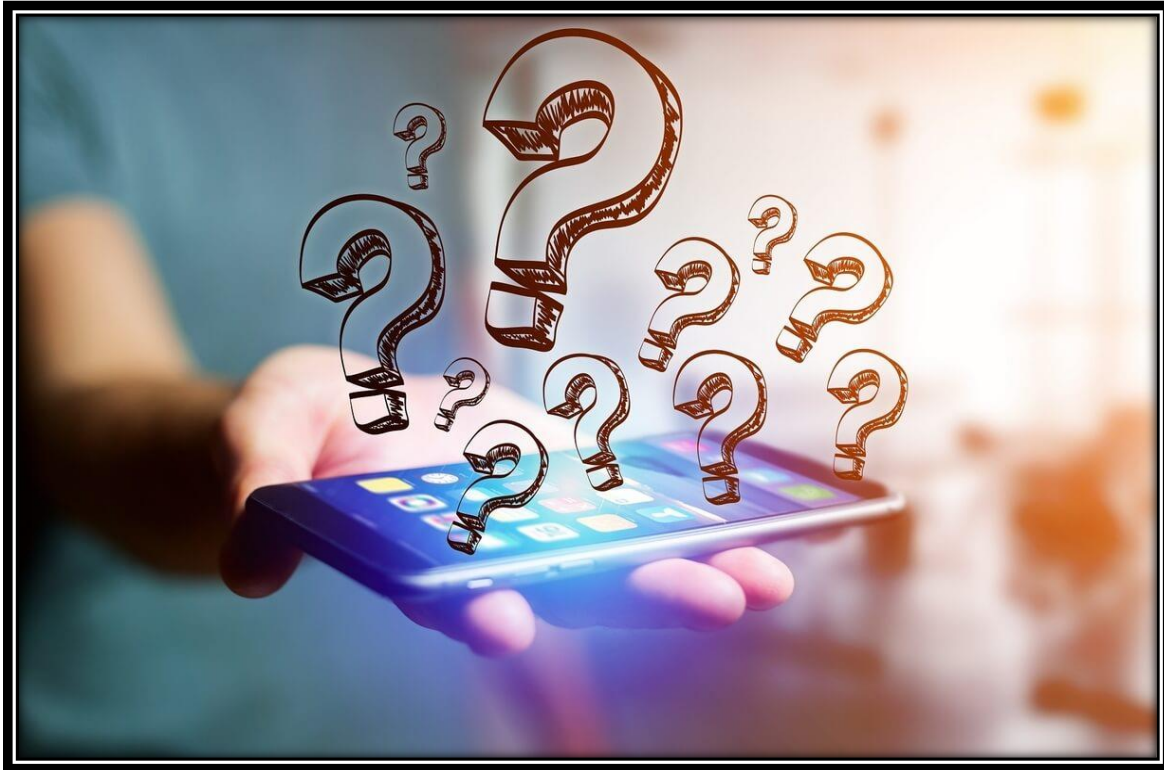University *of* New Haven

TAGLIATELA COLLEGE OF ENGINEERING

Master of Science in Data Science (MSDS)

# Cloud-Based Mobile Price Prediction



**SPRING 2024**

# Contents

# 1. Abstract:

This report presents the development and deployment of a machine learning model for predicting mobile phone prices. The model utilizes datasets containing information about mobile phone features and prices, with the goal of assisting consumers and businesses in making informed decisions about purchasing and pricing mobile devices.

# 2. Introduction:

The proliferation of mobile devices has led to a diverse range of features and price points in the market. Predicting mobile phone prices accurately can be valuable for consumers looking to purchase a device within their budget and for businesses seeking to optimize pricing strategies. In this project, we aim to develop a machine learning model capable of predicting mobile phone prices based on their features.
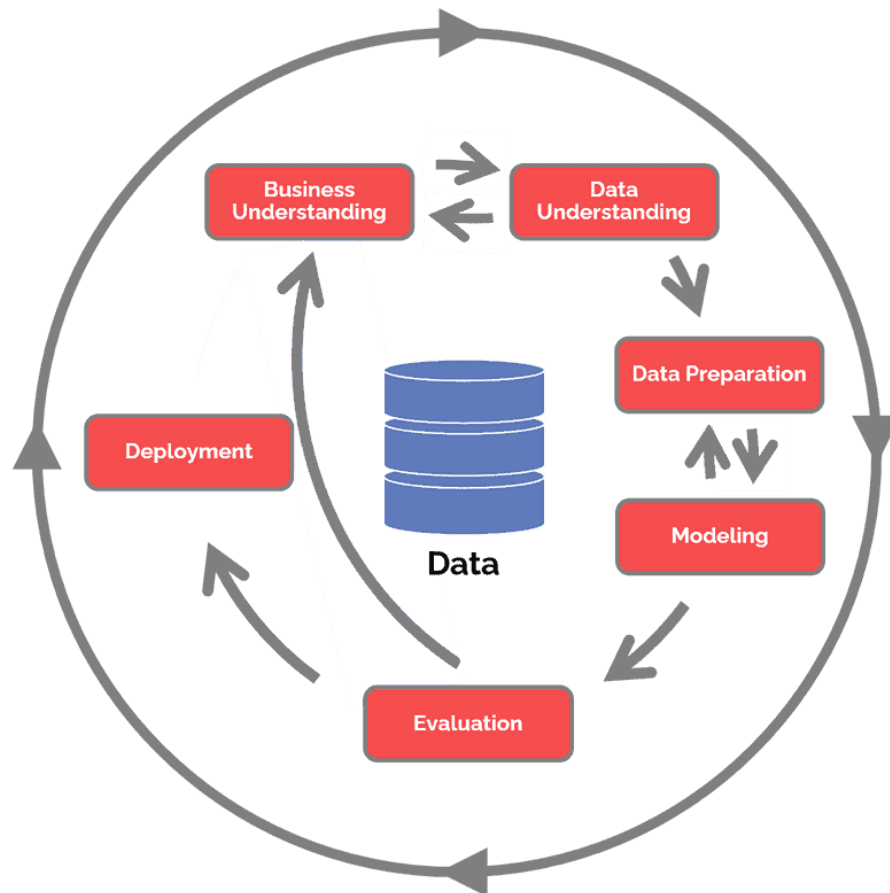
# 3. Objective:

The primary objective of this project is to develop a machine learning model that accurately predicts mobile phone prices based on their features. By leveraging datasets containing information about mobile phone specifications and prices, we aim to create a model that can assist consumers and businesses in making informed decisions about mobile device purchases and pricing strategies.

# 4. Overview:

The project follows a structured approach, encompassing data gathering, data preparation, exploratory data analysis (EDA), feature selection, model training, model evaluation, and model deployment. Key tools and services utilized include Amazon SageMaker for model development and deployment, Amazon S3 for data storage, and AWS CloudWatch for monitoring pipeline activities.

## 5. CRISP-DM:

The Cross-Industry Standard Process for Data Mining (CRISP-DM) methodology serves as the framework for this project, guiding the iterative process of data exploration, model development, and deployment.
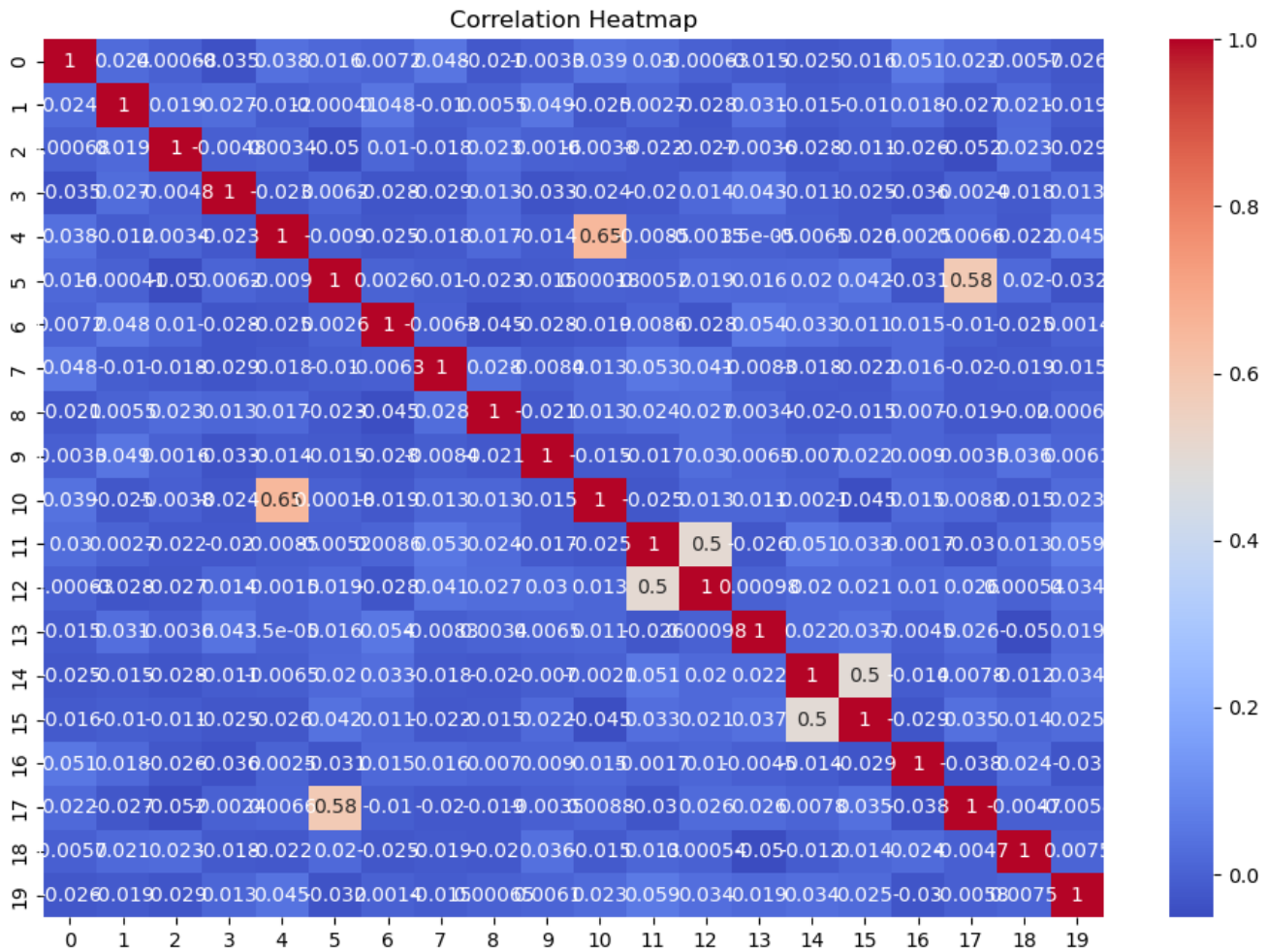


## 6. Data Gathering:

Data gathering involves collecting datasets containing information about mobile phone features and prices. These datasets are stored in Amazon S3 buckets for easy access and scalability. AWS CloudWatch is used to monitor S3 buckets for data availability, access patterns, and storage usage.

## 7. Data Preparation:

Data preparation includes cleaning and formatting the datasets to ensure they are suitable for model training. This process involves handling missing values, encoding categorical variables, and normalizing numerical features. Python libraries such as Pandas and NumPy are used for data preprocessing tasks, with AWS CloudWatch monitoring the data preparation process.

## 8. Exploratory Data Analysis (EDA):

EDA is performed to gain insights into the datasets and understand the relationships between different features and the target variable (mobile phone prices). Jupyter Notebooks and Python libraries such as Matplotlib and Seaborn are used for visualization and statistical analysis. AWS CloudWatch monitors the execution of EDA scripts, tracking metrics such as notebook execution time and memory usage.

Correlation Heatmap

## 9. Feature Selection:

Feature selection is conducted to improve model performance and efficiency by selecting the most relevant features. Features are chosen based on their importance, as determined by EDA and domain knowledge.

```python
# Calculate correlations between features and target variable
correlations = train_df.corrwith(pd.Series(y_train))

# Sort the correlations in descending order
sorted_correlations = correlations.abs().sort_values(ascending=False)

# Select the top k features (e.g., top 5)
top_k_features = sorted_correlations.head(5).index.tolist()

print("Top 5 features with highest correlation with target variable:")
print(top_k_features)
```

```
Top 5 features with highest correlation with target variable:
[13, 0, 12, 11, 6]
```

## 10.   Data Engineering Pipeline:

The data engineering pipeline leverages several AWS services, including Amazon S3 for data storage, Amazon SageMaker for model training and deployment, and AWS CloudWatch for monitoring pipeline activities.

# Project Pipeline & Tools Used:



Used Amazon S3 for storing training and testing data and model dumps

Amazon SageMaker Studio Domain

Amazon SageMaker model endpoint

Flask

Collect   Prepare   Build   Tune   Deploy

Model training and deployment using inbuilt xgboost algorithm

Amazon CloudWatch

monitoring model performance metrics.

## 11.   Model Training:

The XGBoost algorithm is chosen for model training, a popular choice for regression and classification tasks. Amazon SageMaker provides built-in support for XGBoost, making it easy to train models at scale using distributed computing resources.

### 11.1   XGBoost Classifier Overview:

XGBoost (Extreme Gradient Boosting) is a widely-used algorithm for classification problems like predicting mobile phone prices. It sequentially builds an ensemble of decision trees, correcting errors made by previous trees to improve accuracy. Key features include regularization to prevent overfitting, parallelization for scalability, and built-in cross-validation for parameter tuning.

## 11.2   Advantages:

- Efficiency: XGBoost is fast and scalable, making it suitable for large datasets like mobile phone features.

- Robustness: It handles outliers well and provides insights into feature importance, aiding in interpreting model results.

**Use Case:**

- Mobile Price Prediction: XGBoost can classify mobile phone prices based on features like battery power, camera specs, and connectivity options.

# 12.    Model Evaluation:

The trained model is evaluated using validation data to assess its performance and generalization capabilities. Evaluation metrics such as accuracy, precision, recall, and F1-score are used to measure the effectiveness of the model in predicting mobile phone prices.

# 13.    Model Deployment:

Once the model is trained and evaluated, it is deployed to a SageMaker endpoint, allowing real-time inference on new data. The endpoint is configured with appropriate settings to ensure optimal performance and cost efficiency.

Deployed Flask App: URL



# 14.  **Monitoring and Management:**

The deployed endpoint and pipeline activities are monitored using Amazon CloudWatch, providing insights into performance, resource utilization, and system health. CloudWatch enables proactive monitoring and alerting for any issues that arise during model training, evaluation, and deployment.

## 15. Conclusion:

In conclusion, the project demonstrates the end-to-end process of developing and deploying a machine learning model for predicting mobile phone prices. By leveraging AWS services and tools, we create a scalable and efficient solution that can assist consumers and businesses in making informed decisions about mobile device purchases and pricing strategies.

## 16. Future Work:

Future work may include exploring additional features and refining the model further to improve prediction accuracy. Additionally, optimization of the deployment architecture for enhanced scalability and cost efficiency could be pursued.

## 17. References:

1. [Amazon SageMaker Documentation](https://docs.aws.amazon.com/sagemaker/index.html)
2. [CRISP-DM Methodology Documentation](https://www.the-modeling-agency.com/crisp-dm.pdf)
3. [XGBoost Documentation](https://xgboost.readthedocs.io/en/latest/)
4. [AWS CloudWatch Documentation](https://docs.aws.amazon.com/cloudwatch/index.html)
5. [AWS CloudTrail Documentation](https://docs.aws.amazon.com/cloudtrail/index.html)
6. [AWS SageMaker SGBoost] https://docs.aws.amazon.com/sagemaker/latest/dg/xgboost.html