# Loan Eligibility Prediction from Mobile Device Data

## 1. Data Acquisition

**Data Source Example:**

- **Dataset:** We propose using the "Mobile Phone Specifications" dataset available on Kaggle, which aggregates GSMArena data.

- **Features:** This dataset typically includes device attributes such as device brand, model, design characteristics, launch price, battery capacity, screen size, technical specifications (e.g., resolution, processor, RAM, storage), and connectivity options (e.g., Wi-Fi, Bluetooth, GPS).

- **Proxy for Socio-Economic Status:** The underlying assumption is that users with premium devices (e.g., devices with higher launch prices, advanced technical specifications, and premium brands like Apple or Samsung) are likely to have higher disposable income and better creditworthiness.

- **Source Verification:** When implementing, confirm dataset availability and licensing terms on Kaggle.

---

## 2. Data Preprocessing & Feature Engineering

**Handling Missing Data:**

- **Imputation:**

  - For numerical attributes (e.g., battery capacity, screen size, RAM, price), apply median imputation to fill missing values.

  - Alternatively, choose models (such as XGBoost) that can handle missing values natively.

**Feature Extraction:**

- **Categorical Encoding:**

  - Encode variables such as brand, model, processor, and operating system using techniques like one-hot encoding or label encoding.

- **Deriving Premium Features:**

  - Create additional features that capture "premium-ness" (e.g., inferred device launch price, battery capacity, display quality).

  - In the absence of explicit price data, consider using a brand–oriented score (with premium brands like Apple or Samsung assigned higher scores).

- **Tools:**

  - Python libraries such as **pandas** and **numpy** are employed for data manipulation, and **scikit-learn** offers various preprocessing utilities.

---

# 3. Unsupervised Clustering for Proxy Labeling

**Rationale:**

- **Proxy Labeling:**

  - Without historical loan data, segment devices based on their attributes to create "proxy labels."

  - Devices with higher-end features or belonging to premium brands are assumed to correlate with higher income, hence a greater likelihood of loan eligibility.

  - Conversely, budget or entry-level devices indicate lower income and higher risk.

**Method:**

- **Clustering Algorithm:**

  - Use K-means clustering (from Scikit-learn) to group devices into clusters (e.g., premium, mid-range, budget).

- **Assigning Proxy Labels:**

    - Based on domain knowledge, assign provisional eligibility labels to each cluster.

        - For instance, designate the cluster with the highest average launch price as "Eligible" (proxy label = 1) and others as "Less Eligible" (proxy label = 0).

---

# 4. Supervised Classification Model

Once proxy labels are derived from the clustering stage, a supervised classification model can be trained to predict loan eligibility directly from device features.

**Model Options:**

1. **XGBoost:**

    - **Pros:**

        - Handles missing values natively.

        - Typically achieves high accuracy with structured/tabular data.

        - Offers regularization features to prevent overfitting.

    - **Cons:**

        - Requires careful parameter tuning.

        - May be viewed as a "black box," though interpretability can be enhanced with tools like SHAP.

2. **Random Forest:**

    - **Pros:**

        - Robust to noise and missing data.

        - Easier to interpret individual trees.

- Performs well with encoded categorical variables.

  - **Cons:**

    - Potentially slower with very large ensembles.

    - Sometimes slightly less accurate than boosting methods.

**Rationale for Choice:**

- Given the challenges (missing data and the absence of historical loan data), **XGBoost** is a strong candidate due to its ability to handle missing features and its performance on tabular data.

- However, if interpretability is paramount, a **Random Forest** may be preferred.

**Python Libraries:**

- **Scikit-learn:** For clustering, train-test splitting, and preprocessing.

- **XGBoost** (or RandomForestClassifier from scikit-learn) for supervised modeling.

- **pandas, numpy, matplotlib, and seaborn** for data manipulation and visualization.

---

# 5. Model Training & Evaluation

## Training Pipeline:

1. **Data Preprocessing:**

   - Clean the data, impute missing values, and encode categorical features.

   - Engineer features from raw device data (e.g., premium score, technical specifications).

2. **Clustering & Labeling:**

   - Apply K-means clustering to segment the devices.

- ○ Assign provisional loan eligibility labels based on cluster characteristics (e.g., the cluster with highest average price is labeled as "Eligible").

3. **Supervised Learning:**

   - ○ **Data Splitting:** Divide the data into training and testing sets (e.g., 80/20 split).

   - ○ **Model Training:** Train the chosen classifier (e.g., XGBoost) on the engineered features with proxy labels.

   - ○ **Model Evaluation:**

     - ■ Use cross-validation.

     - ■ Evaluate performance with metrics like accuracy, precision, recall, and a confusion matrix.

4. **Deployment:**

   - ○ Package the trained model (e.g., using Flask for an API) for real-time predictions using only the available device attributes.

---

# 6. Summary & Discussion

## Overall Workflow Recap:

- ● **Data Collection:**

  - ○ Mobile device data sourced from Kaggle (GSMArena data) acts as a proxy for user device information and, indirectly, for financial standing.

- ● **Preprocessing:**

  - ○ Missing values are imputed and categorical features are encoded.

  - ○ Features are engineered to capture aspects like brand reputation and inferred price.

- ● **Unsupervised Clustering:**

- K-means clustering segments devices into groups (e.g., premium, mid-range, budget).

- Proxy labels for loan eligibility are assigned based on the assumption that premium devices correlate with higher creditworthiness.

- **Supervised Classification:**

  - An XGBoost (or optionally Random Forest) model is trained to predict loan eligibility from the device features and proxy labels.

- **Evaluation & Deployment:**

  - Model performance is rigorously evaluated using standard classification metrics.

  - The final model can be deployed as a service for real-time predictions.


## Pros & Cons of Model Choices:

**XGBoost:**

- **Pros:**

  - Excellent performance on structured data.

  - Handles missing data well.

  - Provides regularization to reduce overfitting.

- **Cons:**

  - Requires careful parameter tuning.

  - Considered less transparent without additional interpretability measures.


**Random Forest:**

- **Pros:**

  - Robust and easier to interpret.

  - Works well out-of-the-box.

- **Cons:**

    - May be slightly less precise than boosting methods.

    - Potentially slower with very large ensembles.

## Assumptions & Limitations:

- **Assumption:** Mobile device data is a reasonable proxy for user financial status.

- **Labeling Limitations:** The proxy labels obtained via clustering rely on domain expertise and might need periodic recalibration.

- **Indirect Predictions:** Since there's no historical loan data, eligibility predictions are indirect. It is advisable to validate these predictions against external financial indicators or expert assessments.

---

# Final Notes

This report outlines a hybrid approach that combines unsupervised clustering with supervised classification to predict loan eligibility from mobile device data. The solution leverages publicly available data (from Kaggle/GSMArena), uses robust Python libraries (pandas, scikit-learn, XGBoost), and integrates domain knowledge for proxy labeling. The proposed workflow is practical for scenarios where historical loan data is missing and mobile device attributes serve as a proxy for a user's socio-economic status.