# Churn Prediction Model Documentation

## Data Preprocessing

## Load Dataset and Initial Exploration

- The provided dataset was loaded using the '**pd.read_excel()**' function to read the Excel file 'churn_dataset.xlsx'.
- Initial exploration was performed using '**df.head()**' to examine the first five rows of the dataset.

## Handling Missing Data and Outliers

- The dataset's structure was checked using '***df.info()***' to identify missing data in columns.
- Visualization techniques such as box plots were employed to identify potential outliers.
- Based on the extent of missing data, a decision was made to impute missing values using mean imputation.
- Outliers were managed using the Z-score thresholding technique to identify and potentially correct extreme values.

## Encoding Categorical Variables and Train-Test Split

- Categorical variables ('Gender' and 'Location') were encoded using '***pd.get_dummies()***' to create one-hot encoded columns.
- The features (X) and target variable (y) were defined as follows:
- Features (X): All columns except '*Churn*'
- Target (y): '*Churn*'
- Data was split into training and testing sets using ***train_test_split()*** with a test size of 0.2.

## Feature Engineering

## Generating Relevant Features

- A new feature *'Usage_Bill_Ratio'* was created by calculating the ratio of *'Total_Usage_GB'* to *'Monthly_Bill'*.

## Feature Scaling

- Numerical features were scaled using the ***StandardScaler*** to ensure consistent scaling across different variables.
- Scaling was applied separately to both training and testing sets.

# Model Building

## Choosing Machine Learning Algorithms

- *Logistic Regression and Random Forest Classifier* were selected as potential models for churn prediction.

```
Logistic Regression:
Accuracy: 0.5037
Precision: 0.49966777408637875
Recall: 0.3789940530188489
F1 Score: 0.431044365470595

Random Forest Classifier:
Accuracy: 0.497
Precision: 0.49268189954722547
Recall: 0.4716258441689346
F1 Score: 0.48192398805232267

Support Vector Machine (SVM) Classifier:
Accuracy: 0.50185
Precision: 0.4975756176402678
Recall: 0.43443201290192524
F1 Score: 0.4638648226874025

Neural Network model:
Accuracy: 0.50185
Precision: 0.4975756176402678
Recall: 0.43443201290192524
F1 Score: 0.4638648226874025

Gradient Boosting Classifier:
Accuracy: 0.50185
Precision: 0.4975756176402678
Recall: 0.43443201290192524
F1 Score: 0.4638648226874025
```

# Training and Validating the Model

- The *Logistic Regression* model was fitted to the training data using the ***model.fit()*** function.
- Cross-validation with 5 folds was applied to assess the model's performance.
- Model metrics ***(accuracy, precision, recall, F1-score)*** were calculated using cross-validation results.

# Model Optimization

## Fine-tuning Model Parameters

- Hyperparameters of the Random Forest Classifier were fine-tuned using ***GridSearchCV***.
- Parameters considered for tuning include ***'n_estimators'***, ***'max_depth'***, and ***'min_samples_split'***.
- Cross-Validation and Hyperparameter Tuning

- Cross-validation was performed with 5 folds to evaluate model performance robustly.
- *GridSearchCV* was used to search for optimal hyperparameter values for the Random Forest Classifier.

# Model Deployment

## Deployment in Development Environment

- A **Flask web application** was developed to serve the trained churn prediction model.
- An API *endpoint '/predict'* was created to accept input data and return churn predictions.
- The trained **Random Forest Classifier** model was serialized and saved as *'trained_model.pkl'*.
- The **StandardScaler** used for preprocessing was serialized and saved as *'scaler.pkl'*.

```python
from flask import Flask, request, jsonify
import pickle
import pandas as pd
from sklearn.preprocessing import StandardScaler

app = Flask(__name__)

# Load the trained model
with open('trained_model.pkl', 'rb') as model_file:
    model = pickle.load(model_file)

# Load the scaler
with open('scaler.pkl', 'rb') as scaler_file:
    scaler = pickle.load(scaler_file)

# Endpoint to receive input data and return predictions
@app.route('/predict', methods=['POST'])
def predict():
    try:
        # Get input data from the request
        input_data = request.json

        # Preprocess input data
        input_df = pd.DataFrame(input_data, index=[0])
        input_scaled = scaler.transform(input_df)

        # Make predictions
        predictions = model.predict(input_scaled)

        # Format predictions
        churn_labels = ['Not Churn', 'Churn']
        prediction_labels = [churn_labels[pred] for pred in predictions]

        return jsonify({'predictions': prediction_labels})
    except Exception as e:
        return jsonify({'error': str(e)})

if __name__ == '__main__':
    app.run(host='127.0.0.1', port=5000)
```

# Testing Model Deployment

- The **Flask** application was run, and the API was tested using sample input data.
- The model's ability to handle new customer data and provide accurate **churn predictions** was verified.

```
 * Serving Flask app '__main__'
 * Debug mode: off
INFO:werkzeug:WARNING: This is a development server. Do not use it in a production deployment. Use a production WSGI server instead.
 * Running on http://127.0.0.1:5000
INFO:werkzeug:Press CTRL+C to quit
```

# Deliverables

## Jupyter Notebook or Python Script:

A comprehensive **Jupyter Notebook** or **Python** script named *'Churn_Prediction_Model.ipynb'* containing the entire process.

## Report on Approach:

- A detailed summary of the approach taken, including explanations of *data preprocessing, feature engineering, and model selection decisions*.
- This report is saved as *'Approach_Report.pdf'*.

## Model Performance Metrics and Visualizations:

*Tables and visualizations* showcasing model performance metrics such as *accuracy, precision, recall, and F1-score*.

## Additional Tips

- Each step of the code process is well-commented to provide clarity and understanding.
- Markdown cells are used in the Jupyter Notebook to provide detailed explanations and documentation.
- The code follows best practices, is organized, and includes relevant libraries, imports, and dependencies.