# Final Project Report: Customer Churn Prediction for Lloyds Banking Group

## 1. Project Overview

This project aims to develop a machine learning model that accurately predicts customer churn for Lloyds Banking Group. By identifying at-risk customers, the business can proactively implement retention strategies, improve service quality, and reduce revenue loss. The work was completed as part of a certified data science initiative in partnership with Lloyds Banking Group.

---

## 2. Data Sources and Description

The dataset used was a structured Excel file with five sheets representing key customer information:

- **Customer_Demographics**: Age, gender, marital status, income level
- **Transaction_History**: Transaction dates, amounts, product categories
- **Customer_Service**: Interaction types and resolution statuses
- **Online_Activity**: Login frequency and service usage
- **Churn_Status**: Binary churn labels (target variable)

Each dataset was joined using a unique `CustomerID`. After merging, the final dataset had 1,041 records and included all relevant behavioral and demographic fields.

---

## 3. Data Cleaning and Preprocessing

### 3.1 Cleaning Steps

- **Null values**: Inspected and dropped non-essential columns with high null values (e.g. timestamps, ID fields)
- **Duplicate removal**: Verified no duplicate entries
- **Dropped columns**: Removed identifiers and dates to avoid data leakage

### 3.2 Feature Engineering

- Encoded categorical variables using a mix of manual mapping, `LabelEncoder`, and `OrdinalEncoder`
- Standardized `AmountSpent` using `StandardScaler`
- Ensured all features were numerical and suitable for model training

---

## 4. Exploratory Data Analysis (EDA)

- **Correlation heatmap** showed weak linear relationships, confirming need for non-linear models
- Churn distribution was imbalanced (\~20% churn)

- Certain features like service usage, login frequency, and resolution status appeared to influence churn rates

---

# 5. Model Development

## 5.1 Algorithm Selection

- Chosen model: **Random Forest Classifier**
- Rationale:
- Robust to noise and outliers
- Handles non-linear features well
- Provides feature importance for business insights
- Performs better than logistic regression or single decision trees

## 5.2 Implementation

- Training/test split: 80/20
- Cross-validation: 5-fold CV
- Hyperparameters tuned via GridSearchCV: `n_estimators`, `max_depth`, `min_samples_split`

```python
from sklearn.ensemble import RandomForestClassifier
from sklearn.model_selection import GridSearchCV

param_grid = {
  'n_estimators': [100, 150],
  'max_depth': [None, 10, 20],
  'min_samples_split': [2, 5]
}

model = GridSearchCV(RandomForestClassifier(random_state=42), param_grid,
cv=5, scoring='recall')
model.fit(X_train, y_train)
```

---

# 6. Model Performance

## 6.1 Confusion Matrix

```
[[824   2]
 [ 11 204]]
```

## 6.2 Classification Report

```
          precision    recall  f1-score    support
```

```
       0        0.99      1.00      0.99       826
       1        0.99      0.95      0.97       215

 accuracy                          0.99      1041
macro avg        0.99      0.97      0.98      1041
weighted avg     0.99      0.99      0.99      1041
```

### 6.3 ROC-AUC Score

   • ROC-AUC: **0.98** (Excellent discrimination)

**Summary:**

   • Very high recall and precision — ideal for retention use-case
   • Model generalises well based on CV scores

---

# 7. Business Recommendations

### 7.1 Using the Model

   • Assign a churn risk score to each customer
   • Prioritize retention campaigns for high-risk individuals
   • Leverage feature importance to understand behavioral patterns (e.g., low login frequency = high churn risk)

### 7.2 Deployment Suggestions

   • Integrate with a **Power BI dashboard** to track churn risk over time
   • Use **SHAP values** for interpretable predictions
   • Retrain model quarterly with new customer data

---

# 8. Certificate of Completion

The successful completion of this project was acknowledged through a certificate issued by **Lloyds Banking Group**.

   • **Certificate Status**: ✅ Awarded
   • **Recognition**: Data-driven model development for a real-world business problem

---

# 9. Conclusion

This project demonstrates the successful development and evaluation of a machine learning model for customer churn prediction. With its strong performance and business relevance, the model can be directly applied by Lloyds Banking Group to improve retention strategies and customer experience. The project was executed end-to-end, covering data acquisition, preprocessing, modeling, evaluation, and actionable insights.

◆ Tools Used: Python, Pandas, Sklearn, Seaborn, Power BI ◆ Key Outcome: 99% Accuracy, 95% Churn Recall, Lloyds Certificate

---

**Prepared by:** Manpreet Sharma\ **Date:** July 2025