

Lloyds Banking Group: Use of Retail Banking Transactional Data

Vidit Sheth

Department of Engineering Mathematics
University of Bristol
Bristol, United Kingdom
vidit.sheth.2023@bristol.ac.uk

3rd Given Name Surname

dept. name of organization (of Aff.)
name of organization (of Aff.)
City, Country
email address or ORCID

2nd Given Name Surname

dept. name of organization (of Aff.)
name of organization (of Aff.)
City, Country
email address or ORCID

4th Given Name Surname

dept. name of organization (of Aff.)
name of organization (of Aff.)
City, Country
email address or ORCID

Abstract—Bank transactions, deliver a wide range of customer activities, they serve as a foundational dataset for analyzing financial behaviors and optimizing banking operations. This project, conducted on artificially generated Transactional data given by Lloyds Bank at the University of Bristol, involved comprehensive data analysis of banking transactions to derive actionable business insights. The study involves involves 2 datasets, both with different features. The analysis was started by understanding the data and conducting RFM Analysis to evaluate customer value based on transaction recency, frequency and monetary value. This analysis helps the bank, tailor and optimize their services to each customer. Further, Clustering has been implemented to analyze spending patterns of the customers. A time series approach to forecast additional records by using ARIMA that provides a view of future customer activity and financial trend. Finally, implemented CatBoost to predict customer churn in the dataset to focus on retention strategies and customer loyalty.

I. INTRODUCTION

Facing with more complexity and competition in today's business, firms need to develop innovation activities to capture customer needs and improve customer satisfaction and retention. [1] In this regard, Customer relationship management is a broadly recognized strategy for acquisition and retention of customers. The main objective of developing customer relations is to make long-lasting and profitable relationships with customers. [2] This study focused on advanced data analytics to understand customer behavior based on their financial transactions. The goal of this study is to discover trends that can predict future behaviors, and enhance personalized services offered by the bank to keep the customers satisfied. To deal with the challenges, RFM (Recency, Frequency, Monetary) analysis is used in order to compliment customer segmentation and to plan strategies directed to each type of customer. Further, ARIMA model has been implemented to forecast future behaviors of customer transactions to have an insight of future spending of customers; these could help the bank plan offers on festive spending beforehand across the entire year. Finally, Customer churn modelled with CatBoost

predicts customer churn to understand customers drifting away from the bank or switching to competitors.

II. LITERATURE REVIEW

Before moving with any analysis, a literature survey was carried out to understand which methodologies have resulted positive giving out successful outcomes. To start with, a study investigating customer segmentation for a banking dataset was carried out by using data mining algorithms to analyze the relationships among the bank customers and their transactions. The segmentation assists the bank offer targeted strategies that focus on customer relationship, risk management, etc. The study aimed to enhance the effectiveness of customer engagement strategies and lower financial risks associated with banking operations. [3] Another study digs into deployment of machine learning techniques for fraud detection in banking. They targeted high volume and multifaceted transactions and fraud schemes. The study by Hashemi and others used class weight-tuning hyperparameters combined with Bayesian optimization to fine-tune machine learning models such as CatBoost, XGBoost, and LightGBM. Their framework enhanced the accuracy and efficiency of fraud detection systems, equipping financial institutions with more robust tools to detect fraudulent activities early and reduce potential losses through more effective and proactive fraud prevention strategies. [4] Also, another interesting paper by Zakrzewska and others focused on application of clustering algorithms for bank customer segmentation, a crucial aspect of knowledge-based marketing in the banking sector. They targeted the challenges of large and multidimensional databases by comparing results of three clustering algorithms: density-based DBSCAN, k-means, and a two-phase clustering process based on k-means. The research curated evaluation of these algorithms based on scenarios of high dimensionality with noise aiming to identify the most suitable approach for segmenting bank customers effectively. [5] Studies related to RFM (recency, frequency,

monetary) modeling and rough set theory (RS theory) combined with the K-means algorithm to enhance data mining processes in customer relationship management (CRM) by Cheng and Chen address dividing continuous attributes and using rough set theory to improve the data handling procedure and shorten training times, the study overcomes the drawbacks of conventional data mining methods. The approach clusters customer value into different classes (3, 5, and 7) to evaluate which offers the best accuracy in representing customer loyalty and assists in identifying key customer characteristics that can strengthen CRM strategies. [6]

III. METHODOLOGY

This study comprised an analytical approach to treat generated transactional data from Lloyds Bank. The aim of this study was to identify trends, spending patterns and any possible insight beneficial for the bank to enhance their services for their customers. The services could range from offering Premium Accounts with benefits for savings rate and rewards on transactions to their best customers, to retaining and establishing customer relations to the customers that have not been engaging enough with the bank. Furthermore, analyzing customer churn to check and engage with customers before they quit services from the bank. Following are the analysis done across the study:

A. RFM Analysis

Customer segmentation is a great way to distinguish the best customers that have been engaging with the company when compared to the total number of users. RFM analysis is a powerful technique used to categorize customers based on their transactional behavior. The customers are categorized by analyzing three specific aspects, Recency - How recently they have made a transaction, Frequency - How frequent they do transactions, and Monetary - How much they spend in terms of the monetary value during their transactions. These classifications are made by assigning R, F and M scores to individual accounts by setting thresholds. The metric used to quantify these scores is 'Quantiles'. The scores are allocated in such a way that accounts falling over 75% or the third quartile receive the highest score, i.e. 4, and accounts falling under 25% or the first quartile are given the score 1.

1) Assigning R, F and M scores::

- The R Score is a critical metric that evaluates the time elapsed since a customer's last transaction. In this study, the R score is calculated based on the 'Recency' feature, which measures the number of days since the last transaction was made from a particular account. The R Score depends on the Recency value; the lower the Recency, the higher the score.
- The F Score measures how often a customer made transactions over a particular period. Considering the size of the dataset, this study takes in account the total period when all the transactions were made. The total frequency of transactions for all accounts is defined by the feature,

'Frequency'. The F Score is dependent on the Frequency, the higher the frequency, the higher the F score.

- The M Score quantifies the total monetary value a customer spent over the period, providing a sum of financial contribution towards the bank. The metric has been calculated based on a feature, 'Monetary'. The data points stored in 'Monetary' represent the total sum of amount that an individual account has spent over the entire period. Therefore, the higher the monetary value, the higher the M Score.

The RFM Score is the sum of all the three metrics which is the final output for the analysis. The highest score possible is 12, which indicates accounts with best recency, frequency and monetary value. Best customers have been filtered out in the study having the perfect RFM Score.

2) *Customer Segmentation based on RFM metrics:* Customer segmentation with RFM Analysis refers to analysing the distribution of Recency, Frequency and Monetary values to categorize customers into distinct groups with respect to their transactional behavior. After careful analysis from the RFM metrics, a total of five segments were selected. Following are the customer segments with their respective threshold values:

- High-Value Customers: Accounts with the highest scores in recency, frequency, and monetary values. All metrics with a score above the 3rd Quartile.
- Loyal Customers: Accounts with high frequency and monetary scores, regardless of their recency score. F Score and M Score which is equal to or above the 3rd Quartile.
- Emerging Customers: Accounts with the highest recency values but lower frequency and monetary values. R Score is above the 3rd Quartile, but the F Score and the M Score is below the 3rd Quartile.
- Lost Customers: Accounts with low scores across recency, frequency, and monetary. All the metrics are less than or equal to the 2nd Quartile or the mean.
- Need Attention Customers: Accounts with medium scores in all categories. All metrics range from the 2nd Quartile to the 3rd Quartile.

B. Clustering with K-Means and DBSCAN

Clustering is an unsupervised machine learning statistical technique used to group data points in such a way that similar data points come close to each other to form a group. These groups are known as clusters. This technique is useful to identify patterns and create segmentations in the data. Some popular clustering techniques are K-Means, K-Medoids, Hierarchical clustering, Density based spatial clustering, etc. In this study, the dataset has been evaluated for K-Means clustering and DBSCAN.

1) *K-Means Clustering:* This clustering technique demands for number of clusters that the algorithm should look for before its execution. The value of 'K' determines the number of clusters. After preprocessing the data for clustering, selecting the optimum number of clusters is particularly important. A technique coined as Elbow method is used in the study to

find the optimum number of clusters required for the dataset. The optimum number of clusters for the dataset turned out to be four. After selecting the number of clusters, K-Means algorithm was defined and fitted to the data which gave four clusters. They have been visualized further with the help of Principal Component Analysis which helped in dimension reduction.

2) *DBSCAN*: Density-Based Spatial Clustering of Applications with Noise (DBSCAN) is a clustering technique that identifies clusters based on regions of high density separated by regions of low density. Unlike many traditional clustering techniques, DBSCAN does not require the number of clusters to be predetermined. This feature makes it particularly effective for analyzing complex datasets where the number or shape of clusters is not known beforehand. It works by exploring the dataset, classifying as noise any points that do not meet the minimum density criteria, thereby highlighting both the clusters and the gaps or noise between them. Implementing this technique gave out a more than a hundred clusters, therefore the method was not taken ahead for further analysis.

C. *ARIMA for forecasting future Volumes*

Autoregressive Integrated Moving Average is a statistical model which takes time series data to forecast future values. The model comprised of three components, p - referring to the number of autoregressive terms, d - representing the number of nonseasonal differences required for stationarity and q - which is the number of lagged forecasts errors in the prediction equation. This study utilized ARIMA to forecast 30 days' worth of new data aiming to predict future bank transactions. The optimal parameters p , d and q values were selected after trying multiple combinations that would minimize the AIC (Akaike Information Criterion) score. After the model was fitted, it was made sure that residuals were independent and normal. Finally, the model forecasts were projected where balances for the subsequent month were given out for each account present in the dataset.

D. *Churn Prediction*

Predicting churn is referred to analyzing the dataset to identify patterns and signals that could lead to potential signal that a customer could discontinue the banking service. It could be done using various features such as transactional behavior, customer interactions, demographic information, and so on. In this study, the features 'Balance' and 'Amount' has been taken into consideration for training a CatBoost model which suits perfect for a challenge like forecasting churn. The model is the good pick as it is robust enough to fight against overfitting and its ability to deal with class imbalance. In majority of cases for churn predictions, the data points for actual churns are highly low, leading to class imbalance. The forecasts successfully predicted churns with a satisfactory accuracy. The insights could be used to develop strategies to reduce churn and improve customer relationship.

IV. DATA DESCRIPTION / PREPARATION

The initial data exploration was done on dataset one, where the data was cleaned and checked for outliers. The dataset was clean therefore further steps were implemented. Feature scaling for creating a new feature was applied firstly to create a new feature `Transaction_category` which would categorise all the amounts and contain the nature of the transaction. The categories ranged from 'Very small transactions' to 'extreme transactions'. Secondly, a new feature was defined to divide the transactions into various type of transactions, for example shopping, cinema, restaurant, grocery, and so on. Furthermore, statistical metrics such as mean, variance, range of transactions for individual categories were calculated and visualised.

The second dataset contained additional features, 'Balance' and 'Timestamp'. These features were analysed to implement the methodologies used in this study to a major extent. All the steps including cleaning the dataset, feature engineering, general statistics calculation were performed on the second dataset as well. The further preparation to execute all the methodologies are mentioned below:

A. *RFM Analysis*

There were no additional data cleaning steps involved to perform RFM Analysis on the dataset. However, the feature 'Date', was converted to a datetime object as it was initially a string data type. The further steps involved in conducting the analysis were included in the Methodology section.

B. *Clustering with K-Means and DBSCAN*

During pre-processing the dataset in order to conduct clustering, it was found that the features 'Balance' and 'Amount' has missing values. These were imputed with the median rather than the mean as the latter is sensitive to outliers and is not always the best fit for financial data. The median values were taken after grouping all the transactions with their respective account numbers. Further, the dataset contained records with null values for transactions that had a single transaction, these were imputed with the overall median for the respective features. Furthermore, the features were scaled using standard scaler, and the scaled features for 'Balance' and 'Amount' were used to fit the clustering models. One-hot encoding was used to encode the categorical feature 'category' before moving ahead with the analysis.

C. *ARIMA Forecasting*

To perform forecasting of new 'Balance' data points from the existing transactions in the dataset. The pre-processing required the conversion of the 'Date' feature from string to the date format, time stamp was not required for the analysis, therefore it was dropped. The forecasted dates were calculated by adding thirty days to the last date of the transaction for a particular account. The model was then fitted and run to predict the forecasted values for the 'Balance' feature.

D. Churn Prediction

The pre-processing steps for predicting churn also involved converting the 'Date' feature to a standard date format. The data was imported every time a new methodology was applied so as to avoid any imputed parameter used in any other step. Further, the Current_date looks for the most recent date of transaction for an account and then adds a day to the same in order to calculate the number of days since the last transaction. The $mathtt{Days_since_last_transaction}$ then finally subtracts the final transaction date for each account from the Current_date to give out the total number of days since the last transaction was recorded. A threshold of more than ninety days was selected to define an account as churn and the target variable in this study. The training variables, 'Balance' and 'Amount' were further imputed for missing values and scaled before the data was splitted into training and testing. The test size of the data was 20% and the training size was 80%. The CatBoost model was further trained for a thousand iterations with a lower learning rate to achieve better performance. Evaluation metric for the model was set to AUC as it tests the model for binary classification. Finally, class weights were given so that the model pays attention to the minority class, which in this study corresponds to the target variable, churn.

V. RESULTS AND DISCUSSION

A. RFM Analysis

Following are the results of the RFM analysis as discussed above. As shown in Figure 1, the individual R, F, and M scores were calculated on the basis of features - Recency, Frequency, and Monetary.

	Recency	Frequency	Monetary	R_Score	F_Score	M_Score	RFM_Score
Account No							
101531259.0	2	123	-716.08	3	1	2	6
104832000.0	4	164	-4565.52	2	2	1	5
105375973.0	4	158	-193.77	2	2	3	7
106601471.0	1	193	10649.00	4	2	4	10
108481285.0	1	368	10038.92	4	4	4	12

Fig. 1. RFM Scores based on Recency, Frequency, and Monetary features

The figure2 is a snippet for the Best Customers who have been classified as the accounts with the highest possible RFM Score.

	Recency	Frequency	Monetary	R_Score	F_Score	M_Score	RFM_Score
Account No							
399538448.0	1	379	6965.00	4	4	4	12
719586818.0	1	336	20576.73	4	4	4	12
558119802.0	1	476	37015.18	4	4	4	12
802697323.0	1	425	17227.74	4	4	4	12
930277104.0	1	364	8672.71	4	4	4	12

Fig. 2. Best Customers with the highest RFM Score

The figure3 is the plot that demonstrates the distribution of the customer segments, as mentioned in the methodology. In order to maximise customer engagement across all segments, High-value customers should continue to get premium services and rewards for their loyalty; Personalised marketing and loyalty initiatives can help keep Loyal customers; Personalised suggestions and welcome incentives can help elevate Emerging customers and boost their transaction volume. Lost Customers should be convinced to return with customised solutions and incentives based on their feedback; and Customers who Need Attention require more personalised offers and promotions to boost their spending and loyalty.

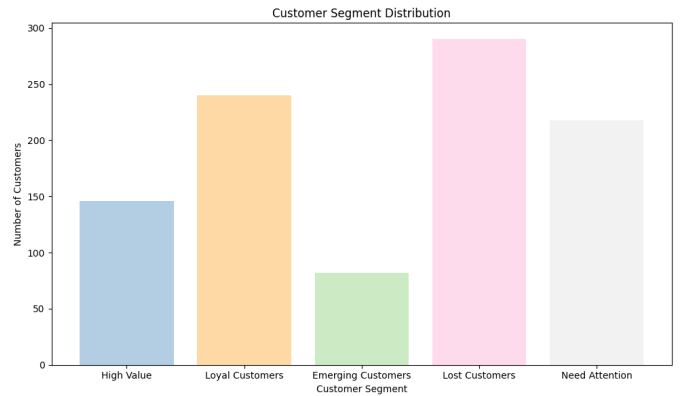


Fig. 3. Customer segmentation based on RFM Analysis with thresholds

B. Clustering

a) v: Implementation of K-Means after implementing the Elbow method gave four distinct clusters. DBSCAN on the other hand gave out a total of four hundred and fourteen clusters. Therefore K-Means is the accurate model chosen for this study. Following is the plot for the Elbow method, suggesting that the number of optimum clusters for the dataset is 4:

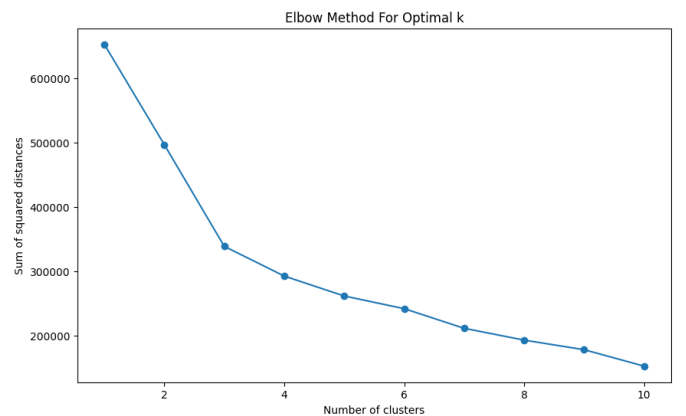


Fig. 4. Elbow plot for K-Means Clustering

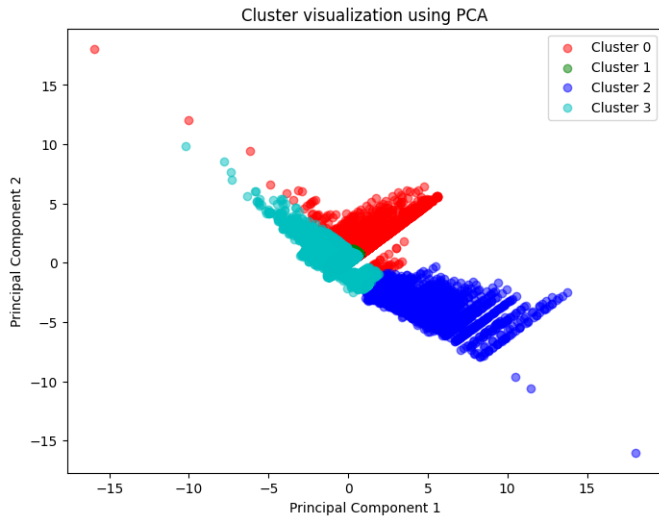


Fig. 5. K-Means with 4 Clusters

b) Further, below is the visualization plot for the clusters with reduced dimensions applying Principal Component Analysis.

C. ARIMA Forecasts

The output from the fitted ARIMA model provided forecasted balances for the next 30 days for each account, which were then saved into a CSV file. The following is depicted by the Figure6. Such forecasts can significantly aid institutions like Lloyds Bank, offering insights that facilitate strategic planning and decision-making on a large scale. Strategies like liquidity management for daily operations and loan disbursements could be managed with ease when a model gives out forecasts with good accuracy.

	Account No	Forecast Date	Forecasted Balance
0	101531259.0	2023-12-06	191.548252
1	101531259.0	2023-12-07	168.958249
2	101531259.0	2023-12-08	147.622271
3	101531259.0	2023-12-09	127.470706
4	101531259.0	2023-12-10	108.437803
5	101531259.0	2023-12-11	90.461463

Fig. 6. Head: ARIMA Forecasts with balances for all the accounts for the next 30 days

D. Churn Prediction

The rigorously trained CatBoost model was trained with scaled balances and amounts to predict churns. In the following figure, each row corresponds to an individual customer account with the forecasts and the actual churn occurrences. The (1) indicates an actual churn while (0) in the

Predicted_class indicates the model forecasting a churn as not churn. The model was trained for 1000 iterations although the AUC score of 89.59% was achieved while testing on the 149th iteration. The forecasts were further saved into a CSV file.

	Account No	Predicted_Class	Churn
89	183546640.0	1	1
210	293191074.0	0	1
276	357651163.0	1	1
433	509627909.0	1	1
488	554771381.0	0	1
500	567499591.0	1	1
525	592961759.0	1	1
623	671824917.0	0	1
649	690941877.0	1	1

Fig. 7. Churn Forecasts by the CatBoost model

VI. FURTHER WORK AND IMPROVEMENT

The focus points for further improvement of this analysis would be to increase the accuracy for the ARIMA and CatBoost model in order to achieve comprehensive results that could contribute to financial institutions in terms of monetary gains. A Churn prediction model with a near-perfect accuracy could be a game changer in the financial sector. Further, this study could be extended to predict anomalies in the transactions. Detecting fraudulent transactions can contribute to a better risk assessment service for the bank and other financial institutions. Anomaly detection can help institutions identify potential threats, improving customer safety. The study could also be extended by using deep learning models that have higher capacity of learning assisted by the hidden layers. These models can be extremely helpful when the size of data is extensively high. Recurrent neural network models such as Long-short term memory network could also be a good option to understand customer behavior and predict churn.

VII. CONCLUSION

In this work, RFM Analysis has been implemented to comprehend customer segmentation and identify the best customers with the highest Recency, Frequency and Monetary aspects. The analysis was further extended to segment the customers based on R, F and M scores. Strategies and customer relation improvement services were recommended for each segment of customers.

Next, K-Means and DBSCAN clustering algorithms were used to identify customer behavior and patterns with respect to clusters. The clusters were defined by the elbow curve plotted by analysing the sum of squared distances from the data points to the closest cluster center point.

Further, ARIMA model was implemented to forecast balances for a period of 30 days from the last date of transaction for all accounts. The results were saved in a forecasts.csv file that is available on the git-hub repository.

Finally, to detect churn in the dataset, a CatBoost model was developed to predict churns in the dataset that would fall under the category where no transactions were made for a period over 90 days. The model was trained to complete 1000 iterations and gave an AUC testing accuracy of 89.59%, after 149 iterations.

REFERENCES

- [1] Huan-Ming Chuang and Chia-Cheng Shen, "A study on the applications of data mining techniques to enhance customer lifetime value — based on the department store industry," 2008 International Conference on Machine Learning and Cybernetics, Kunming, China, 2008, pp. 168-173, doi: 10.1109/ICMLC.2008.4620398.
- [2] Ravi, Vadlamani. (2007). Advances in Banking Technology and Management: Impacts of ICT and CRM. 10.4018/978-1-59904-675-4.
- [3] Shokrgozar, Neda. (2016). Customer Segmentation of Bank Based on Discovering of Their Transactional Relation by Using Data Mining Algorithms. Modern Applied Science. 10. 283. 10.5539/mas.v10n10p283.
- [4] Hashemi Seyedeh , Mirtaheeri Seyedeh , Greco Sergio. (2022). Fraud Detection in Banking Data by Machine Learning Techniques. IEEE Access. PP. 1-1. 10.1109/ACCESS.2022.3232287.
- [5] D. Zakrzewska and J. Murlewski, "Clustering algorithms for bank customer segmentation," 5th International Conference on Intelligent Systems Design and Applications (ISDA'05), Warsaw, Poland, 2005, pp. 197-202, doi: 10.1109/ISDA.2005.33.
- [6] Cheng Ching-Hsue, Chen You-Shyang. (2009). Classifying the segmentation of customer value via RFM model and RS theory. Expert Systems with Applications. 36. 4176-4184. 10.1016/j.eswa.2008.04.003.

APPENDIX

The document up to this section should be no more than 8 pages. The appendix section is optional. You can include additional material here, but it will not be marked.