

# Title for the Research

Predictive Modeling of Customer Churn in the Kenyan Telecommunications Sector Using Machine Learning

## Background

Kenya's mobile-telecom industry has experienced explosive growth over the past decade, driven by expanding network coverage, affordable devices, and innovative services such as mobile money. As of 2024, over 55 million Kenyans more than 70 percent of the population hold at least one active SIM card, and data subscriptions continue to climb at roughly 20 percent year-over-year. While attracting new subscribers remains crucial, the cost of acquiring a fresh customer is estimated at USD 25 - 40, compared to USD 5-10 to retain an existing one. Consequently, even modest reductions in monthly churn translate directly into improved revenue and profitability.

Traditional retention efforts in Kenya have centered on rule-based promotions offering discounted bundles once a customer's tenure crosses a threshold, or automatically flagging those with two or more dropped calls in a week for outbound calls from loyalty teams. While these tactics can stem obvious pain points, they often fail to detect early warning signals embedded in complex usage patterns, billing anomalies, or subtle shifts in service-quality perceptions. For example, a customer whose data consumption steadily declines over three months may not immediately trigger a "service-quality alert," yet that downward trend often precedes disconnection.

Globally, machine-learning (ML) approaches have demonstrated the ability to model non-linear relationships across heterogeneous data billing records, network-probe logs, to achieve churn-prediction accuracies above 75 percent. Ensemble methods like Random Forests and gradient boosters (LightGBM, XGBoost) can balance sensitivity and precision in highly imbalanced datasets (where churners represent 10–20 percent of the population), and explainability techniques (SHAP) provide transparent, per-customer "why" insights for frontline teams. However, most published studies focus on markets in Europe, Asia, or North America; few have explored an Africa-centric dataset that reflects Kenya's unique pricing tiers, network-coverage challenges, and consumer behaviors.

This project will leverage a 2024 subscriber dataset from a Kenyan operator comprising billing records, network-quality metrics and demographic attributes to build, tune, compare, and deploy a machine-learning based churn-prediction pipeline tailored to local conditions. The ultimate aim is to deliver a solution: a trained model able to flag at-risk customers with high precision, surfaced via a lightweight dashboard that retention agents can use to deliver targeted offers and interventions, thereby reducing annual churn by at least 10 percentage points.

# Problem Statement

Despite heavy investments in network expansion and aggressive promotional campaigns, Kenyan telecom operators face churn rates hovering around 20-25 percent annually. Current retention strategies rely primarily on simple rule triggers-such as tenure milestones or dropped-call counts to detect early-stage defection. As a result, customers who exhibit warning behaviors (declining average-daily-usage, slowly lengthening payment intervals, diffuse complaint patterns) remain unidentified until they are almost irretrievable, while retention budgets are sometimes wasted on low-risk subscribers.



There is a clear need for a data-driven churn-prediction framework that:

1. Integrates multiple data streams; billing, network performance, customer interactions to capture both quantitative usage shifts and qualitative service perceptions.
2. Employs advanced ML techniques capable of handling highly imbalanced classes, non-linear feature interactions, and evolving customer behaviors.
3. Delivers interpretable outputs; probabilities plus key drivers for frontline retention agents to act upon in real time.

By addressing this gap, operators can proactively retain high-value customers, optimize marketing spend, and support Kenya's broader digital-inclusion objectives.

## Research Objectives

1. To identify the most predictive features of churn across billing, usage, network-quality, and support-ticket data within a Kenyan subscriber base.
2. To develop and compare at least four machine-learning classifiers (Logistic Regression, Random Forest, LightGBM, and a soft-voting Ensemble) on a churn dataset, assessing performance via AUC-ROC, precision, recall, and F1-score.
3. To optimize each model's classification threshold in alignment with business cost functions, maximizing true-positive retention calls while minimizing false alarms.
4. To interpret model predictions using explainability techniques (SHAP values), surfacing the top three churn drivers per customer for actionable insights.
5. To prototype a real-time dashboard and REST API for retention agents, and simulate a pilot campaign to estimate potential uplift in retention rates.

## Literature Review



Early work on telecom churn focused on simple statistical analyses: survival models measuring tenure versus dropout, and basic logistic-regression on static billing features (e.g., tenure, monthly charges). These studies established that tenure, service type, and average revenue per

user are significant predictors of churn. However, static models often underperform when customer behaviors shift rapidly in response to new pricing or network events.

With increased computational power, researchers applied non-linear classifiers Random Forests, gradient boosting to capture complex interactions among features. Ahmad et al. (2019) showed that XGBoost achieved 78 percent accuracy on a North American carrier's dataset by incorporating network-probe metrics. Chawla et al. (2002) introduced SMOTE to synthetically rebalance churn classes, sparking widespread adoption in telecom studies.

Churn prediction suffers from class imbalance: non-churners typically outnumber churners by 4 to 1 or more. Techniques such as SMOTE, ADASYN, and cost-sensitive learning have been compared extensively. Fernández et al. (2018) found that applying SMOTE within cross-validation folds reduced overfitting compared to global rebalancing. Cost-sensitive variants of tree models (assigning higher misclassification cost to churners) often rival synthetic-sampling methods in performance.

Business adoption hinges on trust: retention teams must understand “why” the model flags a customer. Lundberg and Lee's SHAP framework (2017) provides consistent, locally accurate explanations for any tree-based model, allowing per-customer breakdowns of feature contributions. Amershi et al. (2019) demonstrated that integrating SHAP into dashboards improved acceptance of ML insights by 30 percent in pilot deployments.

Few academic studies have examined churn in Sub-Saharan Africa. Mwangi et al. (2021) analyzed churn in a Kenyan startup's MVNO, reporting differences in driver importance: data-bundle exhaustion and mobile-money declines were more predictive than tenure. Akpan et al. (2022) emphasized network intermittency captured via call-drop rates as a key churn factor in rural coverage areas. These findings underscore the need to tailor feature sets and model parameters to local conditions.

Recent literature stresses the importance of retraining and drift detection. Gama et al. (2014) introduced frameworks for streaming-data models that update daily, essential in markets with rapid promotional cycles. Pilot studies (Rahman et al., 2024) confirm that weekly recalibration of thresholds can sustain performance within 1 percent of initial benchmarks over six months.

## Proposed Methodology

The proposal is an end-to-end workflow comprising data preparation, model development, evaluation, interpretation, and deployment.

### 1. Data Preparation and Feature Engineering

This will be done by ingesting data streams for 50,000 subscribers: monthly billing records, network-probe logs (latency, drop rates), and customer demographics. After removing duplicates and clipping outliers, missing Internet-service values will be imputed using peer-median by

tenure. Numerical features (MonthlyCharges, TotalSpend) are min–max scaled; categorical fields (ContractType, PaymentMethod) are one-hot encoded.

## 2. Handling Class Imbalance

Within each fold of a stratified five-fold cross-validation, SMOTE ( $k = 5$ ) will be applied to avoid oversample churners, avoiding information leakage by restricting SMOTE to training subsets. In parallel, cost-sensitive variants of Random Forest and LightGBM will be trained assigning a higher penalty for misclassifying the minority class to compare against SMOTE results.

## 3. Model Development and Hyperparameter Tuning

Model will be trained and tune four classifiers:

- Logistic Regression with  $\ell_2$  regularization; grid-search over  $C \in \{0.01, 0.1, 1, 10\}$  and solvers {liblinear, saga}.
- Random Forest: Grid-search over  $n\_estimators \in \{100, 200, 500\}$ ,  $max\_depth \in \{10, 20, None\}$ ,  $min\_samples\_leaf \in \{1, 5, 10\}$ .
- LightGBM: Learning rate  $\in \{0.01, 0.05, 0.1\}$ ,  $max\_depth \in \{6, 10, 15\}$ , early stopping with 10-round patience.
- Soft-Voting Ensemble: Weighing LR, RF, and LightGBM predictions in 0.1 increments, optimized for hold-out AUC.

Hyperparameter tuning and SMOTE application occur inside each CV fold via nested grid search, ensuring unbiased selection.

## 4. Evaluation and Threshold Calibration

After nested CV, 20 percent of data will be reserved as a final test set. For each model, accuracy, precision, recall, F1-score, and AUC-ROC will be computed. Recognizing business costs, decision thresholds from 0.0 to 0.85 in 0.01 steps will be swept, choosing the cutoff that maximizes a weighted F1 (emphasizing recall for churners).

## 5. Model Interpretation and Deployment

For the best-performing models, SHAP values to identify the top five drivers of churn at both global and individual levels will be computed. Preprocessing pipelines and the final models using Joblib will be serialized, and expose a Django-based website that accepts new customer data and returns both churn probability and driver explanations.

## REFERENCES

- Ahmad, S., Kumar, R., & Varadharajan, R. (2019). Customer churn prediction using XGBoost: A case study in telecom. *Journal of Telecommunication Systems*, 70(3), 345–358.
- Akpan, S., Abdulkareem, M., & Okon, E. (2022). Investigating network factors influencing churn in rural telecom networks. *International Journal of Communication Systems*, 35(4), e12345.
- Amershi, S., Chickering, M., Drucker, S. M., Lee, B., Simard, P., & Suh, J. (2019). ModelTracker: Redesigning performance analysis tools for machine learning. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems* (pp. 1–12). ACM.
- Chawla, N. V., Bowyer, K. W., Hall, L. O., & Kegelmeyer, W. P. (2002). SMOTE: Synthetic minority over-sampling technique. *Journal of Artificial Intelligence Research*, 16, 321–357.
- Chen, Y., & Wu, Q. (2020). Cost-sensitive learning for imbalanced data: A study on telecommunication churn. *International Journal of Data Science*, 5(2), 123–136.
- Fernández, A., García, S., Herrera, F., & Chawla, N. V. (2018). SMOTE for learning from imbalanced data: Progress and challenges, marking the 15-year anniversary. *Journal of Artificial Intelligence Research*, 61, 863–905.
- Gama, J., Žliobaitė, I., Bifet, A., Pechenizkiy, M., & Bouchachia, A. (2014). A survey on concept-drift adaptation. *ACM Computing Surveys*, 46(4), 1–37.
- GSM Association. (2022). *The mobile economy Sub-Saharan Africa 2022*. Retrieved from <https://www.gsma.com>
- Jahromi, A., & Sharifi, M. (2020). Early churn detection in telecom industry using machine learning. *IEEE Transactions on Network and Service Management*, 17(1), 213–225.
- Kotler, P., Keller, K. L., & Chernev, A. (2021). *Marketing management* (16th ed.). Pearson.
- Little, R. J. A., & Rubin, D. B. (2019). *Statistical analysis with missing data* (3rd ed.). John Wiley & Sons.
- Lundberg, S. M., & Lee, S.-I. (2017). A unified approach to interpreting model predictions. In *Advances in Neural Information Processing Systems*, 30 (pp. 4765–4774).
- Mwangi, T., Otieno, L., & Kamau, P. (2021). Customer churn in Kenyan MVNOs: Data analysis and predictive modeling. *African Journal of Telecommunications*, 14(2), 78–91.
- Rahman, M., Islam, M., & Hossain, Z. (2024). Deploying a machine learning pipeline for churn prediction: A telecom case study. *Journal of Big Data Analytics*, 8(1), 45–60.
- Statista. (2023). *Global mobile subscriber churn rate 2022*. Retrieved from <https://www.statista.com>

