

Data Science Canvas		Project:	Nigerian Banking Fraud Detection System				
		Team:	Data Conduits				
Problem Statement				Execution & Evaluation		Data Collection & Preparation	
Business Case & Value Added <p>Nigerian banks lost ₦17.67B (~\$216M) in 2023, with fraud cases rising 23% YoY. Traditional rule-based systems struggle with evolving fraud patterns and generate high false positives. This project delivers a machine-learning-driven fraud detection system that:</p> <ul style="list-style-type: none">i. Reduces fraud losses through early detectionii. Minimizes false positivesiii. Automates low-risk decisions to reduce manual review workloadiv. Enables fraud teams to focus on high-risk, high-value casesv. Improves compliance and customer trust with explainable predictions	Model Selection <p>Given the extreme class imbalance (0.3% fraud cases), the solution uses:</p> <ul style="list-style-type: none">Gradient boosting models (LightGBM, XGBoost) for non-linear patternsEnsemble stacking to leverage complementary feature importanceLSTM embeddings to capture short-term temporal transaction patternsOptuna for hyperparameter and ensemble weight optimizationPrecision-Recall threshold tuning to maintain extremely low FPR	Model Requirements <p>Performance: $F1 \geq 0.90$ (achieved 0.8847), $AUC-ROC \geq 0.95$, $FPR < 0.1\%$ (achieved 0.001%)</p> <p>Explainability: SHAP-based interpretability for audit and compliance</p>	Skills <p>Data Engineering: Pipeline design for high-volume transaction data</p> <p>Data Science: Feature engineering (temporal, velocity, interaction), ensemble modeling</p> <p>ML Engineering: Optuna tuning, model serving, real-time scoring pipelines</p> <p>Risk/Fraud Expertise: Interpret patterns, label verification, regulatory compliance</p>	Model Evaluation <p>Performance is assessed using $F1$, $AUC-ROC$, PR AUC, and precision-recall balance due to extreme imbalance. The ensemble achieved:</p> <p>AUC: 0.9638 $F1$: 0.8847</p> <p>FPR: 0.001% (100x better than requirement)</p> <p>Confusion matrix analysis guides threshold adjustment. Real-time monitoring includes drift detection, data quality checks, and alerting when fraud distribution shifts.</p>	Data Storytelling <p>Target Group Reqs Clear, business-focused insights on fraud patterns affecting Nigerian payment channels.</p> <p>Operationally relevant metrics (fraud hotspots, transaction patterns, risk levels by customer/merchant segment).</p> <p>Actionable recommendations that regulators, banks, and fintech operators can implement immediately.</p> <p>Simple, interpretable visuals that work for mixed stakeholders (analysts, compliance teams, executives).</p> <p>Traceability and transparency in how fraud was detected, especially for regulatory review.</p> <p>How? Lead with key fraud insights (where, how, and why fraud occurs in the Nigerian context). Use clean visuals (heatmaps, trend lines, risk scores) instead of technical plots. Translate findings into operational actions: improved KYC checks, velocity limits, merchant risk tiers, model cutoffs. Highlight local relevance: mobile money patterns, high-risk time windows, account takeover behaviors typical in Nigeria.</p>	Data Selection & Cleansing <p>Relevant features like: Temporal: hour, day, rolling time gaps Velocity: tx_count_24h, amount_sum_24h Behavioral: ratio_to_customer_max, deviation features Interaction: composite risk scores</p> <p>Data cleansing included: Handling missing timestamps Outlier detection for extreme transaction amounts Encoding categorical fields (channel, location, merchant)</p>	Data Collection <p>Real-time data from mobile banking, transfers, POS, ATM, web</p> <p>Fraud labels from investigation teams</p> <p>Need for sequence data per customer (for LSTM/transformer models)</p>
Data Landscape <p>NIBSS Fraud Dataset (Kaggle) — synthetic data that is representative of transaction data reflecting Nigerian banking characteristics.</p> <p>Synthetic data is chosen as banking data is sensitive, and personal.</p>		Software & Libraries <p>Python pandas, NumPy scikit-learn LightGBM, XGBoost, CatBoost Optuna PyTorch/Keras for LSTM SHAP Matplotlib/Seaborn</p>			Data Integration <p>Different data sources (transaction logs, customer profiles, merchant metadata, fraud investigation labels) are unified into a centralized feature store.</p>	Explorative Data Analysis <p>Fraud prevalence: 0.3% Fraud concentrated 1 AM – 4 AM High-risk channels: mobile, web Fraud values skew higher than legitimate ones Lagos & Abuja show highest fraud density Temporal and velocity patterns drive strong predictive power</p>	