

<h2>Data Science Canvas</h2>	<b>Project:</b>	Decoding the Digital Rupee: A Predictive Analysis of India's UPI Spending					
	<b>Team:</b>	Jabin Shalom S, Suganya H					
<b>Problem Statement</b>			<b>Execution &amp; Evaluation</b>		<b>Data Collection &amp; Preparation</b>		
<b>Business Case &amp; Value Added</b>  Objective: Investigate if technology access (Device/Network) creates a "Digital Divide" limiting UPI economic participation.  Value Added: Quantifies the divide in monetary terms (Rupees).	<b>Model Selection</b>  <b>Methods:</b> Random Forest Regressor (Champion), XGBoost (Runner-up), Linear Regression (Baseline)  <b>Justification:</b> UPI spending behavior is highly non-linear. Tree-based ensemble models (Random Forest/XGBoost) effectively capture complex interactions between demographics and technology, whereas simple linear models failed ( $R^2 < 0$$ )	<b>Model Requirements</b>  Type: Supervised Regression (Predicting a continuous transaction value).  Constraints: Must be capable of handling high-cardinality categorical features (e.g., States, Merchant Categories) efficiently.  Explainability: A "White Box" approach is mandatory to fulfill the project objective of identifying key drivers; techniques like SHAP (Shapley Additive ExPlanations) are required to provide interpretable feature importance rankings.	<b>Skills</b>  Data Preparation: Proficiency in feature engineering, data cleaning, and handling missing values.  Modeling: Expertise in implementing and tuning machine learning models, particularly Random Forest, Gradient Boosting (XGBoost), and Linear Regression.  Analysis & Interpretation: Skills in interpreting model results, optimizing hyperparameters, and deriving actionable insights.  Data Visualization: Ability to create effective visualizations for EDA and final reporting.	<b>Model Evaluation</b>  Quality Indicators:  RMSE (Primary): Measures average prediction error in Rupees. Lower is better.  R <sup>2</sup> (Secondary): Validates the model's ability to capture non-linear spending patterns (Digital Divide signal).  Interpretability: SHAP values rank feature importance to confirm logical drivers (e.g., Device Type).	<b>Data Storytelling</b>  Target Group: Financial Policymakers and Bank Executives requiring actionable insights on financial inclusion and market segmentation.  Communication Strategy:  Visual Proof: Use the "Digital Privilege Index" (Staircase Box Plot) to visually demonstrate that better tech equals higher spending capacity.	<b>Data Selection &amp; Cleansing</b>  Selection: Selected relevant features: Sender_State, Device_Type, Network_Type, Merchant_Category, and Amount  Cleansing: Dropped high-cardinality identifiers (e.g., Transaction ID) and raw timestamps to prevent overfitting.	<b>Data Collection</b>  Method: Download the dataset directly from Kaggle as a CSV file.  Properties: The data must contain row-level transaction details (e.g., Amount, Device Type, Network Type) to allow for granular analysis. Since real banking data is restricted, this synthetic dataset serves as a proxy for the required properties.

			<p>Programming: Proficiency in Python and relevant libraries (Pandas, Scikit-learn).</p> <p>(one-time static analysis).</p> <p>Future: Required to detect data drift (e.g., inflation, new tech) in a live deployment.</p>	<p>based on tech profiles.</p> <p>Impact Analysis: Use the "Sector Domination" chart to highlight specific economic sectors where low-tech users are excluded.</p>	<p>Engineering: Extracted cyclic features (Hour_of_Day, Is_Weekend) and created the composite Tech_Score to quantify digital access.</p>	
<p><b>Data Landscape</b></p> <p>Required Data: Transaction-level logs containing Amount, Timestamp, and user metadata</p> <p>Available Data: "UPI Transactions 2024 Dataset" from Kaggle (Synthetic, 250,000 rows).</p> <p>.</p> <p>Data Gap: Real longitudinal user history is missing in the synthetic set but would be required for a production-grade banking model.</p>	<p><b>Software &amp; Libraries</b></p> <p>Software: Python (Jupyter Notebook/Google Colab) was used as the primary environment for analysis and model development.</p> <p>Deployment: Streamlit was used to create the interactive web application ("Insight Engine") for demonstration.</p> <p>Libraries:</p> <p>Data Handling: Pandas, NumPy.</p> <p>Visualization: Matplotlib, Seaborn, Plotly.</p> <p>Machine Learning: Scikit-learn, XGBoost.</p> <p>Explainability: SHAP (Shapley Additive exPlanations).</p>	<p>Data Integration</p> <p>System: The data is loaded into a Python environment (using Pandas within Jupyter Notebook/Google Colab) for processing and analysis.</p> <p>Migration: Since the project relies on a single CSV source file, complex data migration or integration from multiple disparate systems is not required.</p>	<p><b>Explorative Data Analysis</b></p> <p>Findings: Identified that transaction amounts are highly right-skewed with significant outliers ("Whales").</p> <p>Correlations: Initial analysis showed weak correlations due to the uniformity of the synthetic data.</p> <p>Key Action: Engineered Tech_Score and Tech_Segment features to expose the latent structure connecting digital infrastructure to spending capacity.</p>			