



TrustLens

Customer Transaction Risk & Behavior Analysis Dashboard

Complete Project Documentation

Python • Streamlit • Scikit-learn • SQLite • Pandas • Matplotlib

1. Executive Summary

TrustLens is a production-quality, end-to-end Data Science project designed to help financial institutions identify high-risk customers through transaction behavioral analysis. The platform combines Exploratory Data Analysis (EDA), SQL-powered business intelligence, and dual Machine Learning models — all surfaced through a polished, real-time interactive Streamlit dashboard.

The core business problem it solves: financial institutions process millions of transactions daily, yet lack fast, interpretable tooling to flag risky customers early. TrustLens bridges that gap by automating risk classification using behavioral signals (credit score, transaction frequency, category, account tenure, and amount) and presenting the results in a format any analyst or business stakeholder can act on immediately.

2. Project Objectives

- Analyze 5,000 synthetic customer transactions to surface behavioral patterns and risk indicators
- Build a fully automated data pipeline: generation → cleaning → SQLite storage → ML training
- Train and compare two ML classifiers (Logistic Regression and Random Forest) for high-risk customer detection
- Surface SQL-based KPIs to answer common business intelligence questions without writing code
- Enable real-time, single-customer risk scoring directly from the dashboard UI
- Produce clean, modular, well-commented code deployable on Replit with zero external setup

3. Problem Statement

Fraud and high-risk customer behavior cost financial institutions billions annually. Traditional rule-based systems are rigid and fail to capture complex behavioral patterns. Machine learning offers a more adaptive approach — but only when backed by quality data, sound feature engineering, and transparent evaluation.

This project simulates a real-world risk analytics pipeline, demonstrating how a data science team would:

1. Ingest and clean messy transactional data
2. Store it in a relational database for business queries
3. Train predictive models with proper evaluation rigour

- Deploy findings through an interactive dashboard accessible to non-technical stakeholders

4. Technology Stack

Layer	Technology	Version	Purpose
Language	Python	3.10+	Core development language
Dashboard	Streamlit	≥1.32	Interactive web UI, no frontend code needed
Data Processing	Pandas	≥2.0	DataFrames, cleaning, feature engineering
Numerics	NumPy	≥1.24	Array operations, random data generation
ML Framework	Scikit-learn	≥1.3	Model training, evaluation, preprocessing
Visualization	Matplotlib	≥3.7	Custom dark-themed chart generation
Visualization	Seaborn	≥0.12	Correlation heatmaps, statistical plots
Database	SQLite (sqlite3)	Built-in	In-memory relational DB for SQL analytics
Deployment	Replit	—	Zero-config cloud hosting for Streamlit



All dependencies are standard Python packages. No external APIs, cloud services, or database servers required. The SQLite database runs entirely in memory, making the project portable and self-contained.

5. Project File Structure

File	Lines	Responsibility
app.py	~320	Main Streamlit dashboard — layout, tabs, KPI cards, sidebar filters, CSS theming
data_generator.py	~80	Synthetic dataset creation: 5,000 transactions with realistic distributions and risk labeling
data_pipeline.py	~80	Data cleaning (imputation, outlier clipping, type casting), SQLite DB builder, SQL KPI runner

File	Lines	Responsibility
eda.py	~130	Four dark-themed EDA chart functions: risk distribution, transaction trend, correlation heatmap, category analysis
ml_models.py	~120	Model training, evaluation metrics, ROC curves, confusion matrices, live single-row prediction
requirements.txt	6	Pinned Python dependencies
README.md	~70	Setup instructions, project overview, business insights summary

6. Data Pipeline Architecture

6.1 Synthetic Data Generation (data_generator.py)

The dataset is generated deterministically using NumPy's default_rng with a fixed seed (42), ensuring perfectly reproducible results across all runs. The generator creates 500+ unique customer profiles and assigns them 5,000 transactions.

Customer-Level Attributes

- Age: uniformly distributed between 18 and 80
- Credit score: normally distributed around 650 ($\sigma=90$), clipped to [300, 850]
- Account tenure: uniformly distributed 1–240 months

Transaction-Level Attributes

- Category: 8 merchant categories with realistic occurrence weights
- Amount: log-normal distribution ($\mu=5.2$, $\sigma=1.1$) scaled by category-specific multipliers (e.g., Travel 2.5x, Grocery 0.8x)
- Date: uniformly sampled from 2023-01-01 through 2024-12-31
- Transaction frequency: uniform integer 1–150 per month

Risk Label Construction

Labels are generated via a logistic function applied to a weighted combination of features — mimicking how real risk models determine ground truth in supervised learning scenarios:

```
log_odds = -3.5
        + 0.015 * (amount / 100)      # higher amounts → more risk
        - 0.004 * (credit_score - 500) # better credit → less risk
        + 0.008 * frequency          # high frequency → more risk
```

```

- 0.003 × tenure           # longer tenure → less risk
+ 0.6 if category == 'Travel'
+ 0.4 if category == 'Electronics'
+ 0.4 if age < 25          # young customers higher risk

```

This produces a naturally imbalanced dataset (~3–8% high-risk rate) which is realistic for fraud/risk detection workloads.

Realistic Noise Injection

- 2% of amount and frequency values are set to NaN (missing values for cleaning exercise)
- 1% of amounts are multiplied by 10–30× (outliers requiring clipping)

6.2 Data Cleaning (data_pipeline.py)

Step	Method	Rationale
Missing imputation	Median fill for amount and frequency	Median is robust to outliers; safe for skewed distributions
Outlier clipping	Clip amount at 99th percentile	Preserves extreme but realistic values; eliminates data entry errors
Type enforcement	Cast all columns to correct dtypes	Prevents downstream type errors in ML pipeline
Feature derivation	year_month, day_of_week, amount_bucket	Adds temporal and categorical grouping signals

6.3 SQLite Database (data_pipeline.py)

Cleaned data is loaded into an in-memory SQLite database via the sqlite3 standard library. This enables standard SQL analytics without any external database infrastructure. The transactions table is indexed on is_high_risk and category for efficient query performance.

SQL KPIs Computed on Load

- Total transaction volume (SUM of all amounts)
- Average transaction amount
- High-risk transaction count and rate
- Count of unique customers

7. Exploratory Data Analysis

All EDA charts use a consistent dark theme (#0d1117 background) styled to match the dashboard aesthetic. Charts are generated as Matplotlib figures and rendered inline by Streamlit. Four core visualizations are included:

Chart	Type	What It Reveals
Risk Distribution	Donut + Bar chart	Overall split of high vs low risk; validates dataset imbalance level
Transaction Trend	Dual-line area chart	Monthly volume by risk class over 2 years; seasonality and drift
Correlation Heatmap	Seaborn heatmap	Pearson correlations between numeric features; identifies predictive signals
Category Analysis	Dual-axis bar + line	Average amount per category (bars) vs risk rate (line); business risk segmentation

8. Machine Learning Models

8.1 Feature Engineering

Six features are used after encoding the categorical 'category' column with scikit-learn's LabelEncoder:

Feature	Type	Preprocessing
Age	Numeric	StandardScaler (LR only)
Amount	Numeric	StandardScaler (LR only)
Credit Score	Numeric	StandardScaler (LR only)
Account Tenure (months)	Numeric	StandardScaler (LR only)
Transaction Frequency	Numeric	StandardScaler (LR only)
Category (encoded)	Ordinal int	LabelEncoder → integer

An 80/20 stratified train/test split is used to preserve the class ratio in both sets. `class_weight='balanced'` is set on both models to handle the imbalanced dataset.

8.2 Model 1 — Logistic Regression

- Algorithm: Linear classifier with L2 regularization (C=1.0)
- Preprocessing: Features are StandardScaled before fitting
- Strengths: Fast training, interpretable coefficients, strong probabilistic calibration
- Best for: Baseline benchmarking, compliance contexts requiring explainability
- max_iter=500 ensures convergence on all feature combinations

8.3 Model 2 — Random Forest

- Algorithm: Ensemble of 200 decision trees with max_depth=8
- Preprocessing: Raw features (no scaling needed for tree-based models)
- Strengths: Captures non-linear feature interactions, produces feature importances
- Best for: Maximum predictive accuracy, feature selection insights
- n_jobs=-1 enables parallel training across all CPU cores

8.4 Evaluation Metrics

Metric	Formula	Why It Matters for Risk Detection
Accuracy	Correct / Total	Overall correctness — misleading for imbalanced datasets
Precision	TP / (TP + FP)	Of predicted high-risk, how many actually are? Controls false alarms
Recall	TP / (TP + FN)	Of actual high-risk, how many caught? Critical — missed risks are costly
F1 Score	$2 \times P \times R / (P + R)$	Harmonic mean of Precision & Recall — best single imbalanced metric
ROC-AUC	Area under ROC curve	Threshold-independent discrimination power; 0.5=random, 1.0=perfect

8.5 Outputs Produced

- Metrics comparison table for both models
- Dual ROC curve plot on a single axes
- Random Forest feature importance horizontal bar chart (top 12 features)
- Confusion matrices for both models side-by-side

9. Streamlit Dashboard — Detailed Walkthrough

9.1 Sidebar Filters (Global)

All dashboard tabs respond dynamically to the sidebar filters:

- Age range slider: 18–80
- Transaction amount slider: \$0–\$10,000
- Category multiselect: include/exclude any of 8 merchant categories
- High-risk only checkbox: instantly filter to flagged transactions

9.2 Tab 1 — Exploratory Analysis

- Risk distribution (donut + bar)
- Transaction amount by category (dual-axis)
- Monthly transaction trend (area line chart)
- Feature correlation heatmap
- Scrollable raw data table (200-row sample, respects filters)

9.3 Tab 2 — SQL Query Explorer

- 4 pre-built business queries selectable from a dropdown
- Editable SQL text area — write any custom SQL against the transactions table
- Run Query button executes and renders results as an interactive Streamlit dataframe

Pre-built queries cover: top high-risk categories, monthly volume trends, risk by age bucket, and high-value flagged customers.

9.4 Tab 3 — ML Model Performance

- Side-by-side metrics table with conditional formatting (green = best per column)
- Dual ROC curve chart
- Random Forest feature importance chart
- Two confusion matrices (one per model)

9.5 Tab 4 — Live Prediction

Users enter 6 parameters (age, amount, category, monthly frequency, credit score, account tenure) via form inputs. On submission, both models score the customer in real time and display:

- Risk probability percentage (e.g., 73%)
- HIGH RISK / LOW RISK verdict with color coding
- Side-by-side comparison between Logistic Regression and Random Forest

9.6 Tab 5 — Documentation

In-app Markdown documentation covering: problem statement, data pipeline stages, ML approach summary, key business insights, how to run on Replit, and project file structure.

10. Key Business Insights



These insights are derived from the synthetic dataset's built-in risk generation logic and reflect patterns a real-world model would surface from similar data.

Insight	Signal	Recommended Action
Credit score is the #1 risk predictor	Strong negative correlation with is_high_risk	Integrate real-time credit bureau checks into transaction approval flow
Travel & Electronics are highest-risk categories	Category-specific risk rate 3–4× baseline	Apply step-up verification (e.g., OTP, manual review) for these categories
Transactions > \$5,000 show elevated risk	Amount is second strongest predictor	Trigger automated review workflow for large-value transactions
Young customers (18–24) are higher risk	Age < 25 adds +0.4 log-odds to risk	Offer credit-building products with lower limits and enhanced monitoring
Long-tenured accounts are safer	Account tenure negatively correlated	Loyalty = trust signal; use tenure as a risk-reduction modifier in scoring
High transaction frequency increases risk	Frequency positively weighted in risk function	Flag velocity anomalies — sudden bursts above customer baseline

11. How to Use the Project

11.1 Running on Replit (Recommended)

5. Create a new Replit project and select Python as the language
6. Upload all 7 project files (app.py, data_generator.py, data_pipeline.py, eda.py, ml_models.py, requirements.txt, README.md)
7. Open the Replit Shell and run the install command:

```
pip install -r requirements.txt
```

8. Start the dashboard:

```
streamlit run app.py
```

9. Click the Replit preview URL — the dashboard opens in your browser

11.2 Running Locally

```
# 1. Clone or download the project files
# 2. Create a virtual environment (recommended)
python -m venv venv
source venv/bin/activate    # Windows: venv\Scripts\activate

# 3. Install dependencies
pip install -r requirements.txt

# 4. Launch
streamlit run app.py

# Dashboard opens at http://localhost:8501
```

11.3 Dashboard Usage Guide

Goal	Where to Go	What to Do
Explore the dataset	Tab 1 — EDA	Use sidebar filters to slice by age, amount, or category; scroll through charts and raw data table
Run business queries	Tab 2 — SQL Explorer	Pick a preset query from the dropdown, or write custom SQL and click 'Run Query'
Compare model performance	Tab 3 — ML Models	Review metrics table, ROC curves, and confusion matrices to choose the best model for your use case
Score a new customer	Tab 4 — Live Prediction	Fill in the 6-field form and click 'Predict Risk' to get instant dual-model risk scores

Goal	Where to Go	What to Do
Share project context	Tab 5 — Documentation	Complete in-app reference covering methodology and business insights

12. Customization & Extension Guide

12.1 Use Real Data

Replace the generate_data() call in app.py with a CSV load:

```
import pandas as pd
raw = pd.read_csv('your_transactions.csv')
```

Ensure your CSV has columns: customer_id, age, amount, category, credit_score, account_tenure_months, transaction_frequency, transaction_date, is_high_risk.

12.2 Add More ML Models

In ml_models.py, extend the models dictionary in train_models():

```
from sklearn.ensemble import GradientBoostingClassifier
models['Gradient Boosting'] = GradientBoostingClassifier(n_estimators=100)
```

12.3 Add More Features

In ml_models.py, update the FEATURE_COLS list and adjust _prepare_features() to include your new columns after engineering them in data_pipeline.py.

12.4 Persist the SQLite DB

Change the db_path in build_db() from ':memory:' to a file path:

```
conn = build_db(df, db_path='transactions.db')
```

13. Design Decisions & Trade-offs

Decision	Choice Made	Alternative & Why Not
Dataset	Synthetic data (5,000 rows)	Real data — privacy concerns and setup complexity; synthetic is reproducible and portable

Decision	Choice Made	Alternative & Why Not
Database	SQLite in-memory	PostgreSQL — overkill for this volume; SQLite requires zero setup
ML Framework	Scikit-learn	PyTorch / TensorFlow — no deep learning needed; sklearn is simpler and faster for tabular data
Dashboard	Streamlit	Dash / Flask — Streamlit is far faster to build and requires no HTML/CSS knowledge
Scaling	StandardScaler for LR only	Scale all models — tree-based models don't benefit from scaling
Imbalanced classes	class_weight='balanced'	SMOTE oversampling — simpler and less prone to overfitting on synthetic minority examples

14. Performance & Scalability Notes

- The full pipeline (data generation → cleaning → DB → ML training) runs in under 5 seconds on a standard Replit instance
- Streamlit's `@st.cache_data` and `@st.cache_resource` decorators ensure data and models are loaded once and reused across all user interactions
- The SQLite in-memory database supports the SQL explorer with sub-millisecond query latency for the 5,000-row dataset
- For production scale (millions of rows), replace SQLite with PostgreSQL, use PySpark or Dask for distributed data processing, and deploy the Streamlit app behind a load balancer
- The Random Forest with `n_estimators=200` trains in ~2 seconds; increase to 500+ trees for marginal accuracy gains at higher latency

15. Future Roadmap

Short Term

- Add XGBoost and LightGBM models for benchmarking
- SHAP value integration for individual prediction explainability
- CSV upload support so users can score their own datasets

Medium Term

- User authentication (Streamlit-Authenticator) for multi-tenant access
- Persistent SQLite on disk with historical query saving
- Automated retraining trigger when data drift is detected
- Email/Slack alerting for high-risk prediction events

Long Term

- Deploy on AWS/GCP with a PostgreSQL backend and proper MLOps pipeline
- Real-time transaction streaming with Kafka integration
- A/B testing framework to compare model versions in production
- Regulatory compliance reporting (GDPR, Basel III) built into the dashboard

TrustLens — Customer Transaction Risk & Behavior Analysis Dashboard

Built with Python · Streamlit · Scikit-learn · SQLite · Pandas · Matplotlib · Seaborn