

Predictive Modeling and Privacy-Preserving Identity for Small Retail Analytics

Author

Tejas Shinde

Independent Research Project (Undergraduate Level)

Abstract

Small retail businesses often lack access to advanced analytical tools to understand and predict customer return behavior, while simultaneously facing privacy, cost, and data-collection constraints. This research proposes a privacy-preserving retail analytics framework that models customer revisit probability using anonymized identifiers and event-based behavioral data. The system avoids storing personally identifiable information and instead relies on interpretable machine learning models suitable for low-resource environments. Due to infrastructure and deployment limitations, the project focuses on system design, modeling logic, feature engineering, and evaluation methodology rather than large-scale empirical training. The study demonstrates how meaningful predictive insights can be generated ethically and practically for small retail settings.

1. Problem Statement

Customer retention is critical for the sustainability of small retail businesses, yet most small retailers operate without data-driven insights into customer behavior. Traditional loyalty systems either provide simple reward tracking or require the collection of invasive personal information such as phone numbers or names, creating privacy risks and operational burdens.

The core problem addressed in this research is:

How can an AI-driven system predict customer return probability without storing or processing personally identifiable information, while remaining interpretable, ethical, and feasible for small-scale retail deployment?

2. Motivation

Advanced customer analytics are largely accessible only to large enterprises with significant infrastructure and data resources. Small retailers are often excluded from AI-driven decision-making tools due to cost, technical complexity, and privacy concerns.

The goal of this research is to explore how machine learning techniques can be adapted to low-resource environments, prioritizing ethical data handling and model interpretability. By focusing on anonymized behavioral patterns rather than personal identity, the system aims to empower small businesses while respecting user privacy and responsible AI principles.

2.5 Practical Context and Prototype Motivation

This research originated from the author's attempt to design a digital loyalty pass system for small retail businesses. The proposed system allows customers to earn loyalty points through visit-based interactions such as QR scans or manual check-ins, with rewards unlocked after a fixed number of visits. During the design phase, key challenges emerged related to customer identity management, privacy preservation, and understanding return behavior. This research formalizes those challenges into an AI-driven, privacy-preserving analytics framework that could support such a loyalty system.

3. System Design Overview

The proposed system consists of four conceptual components:

- 1. Data Collection Layer**

Captures timestamped customer visit events without collecting personal identifiers.

- 2. Anonymized Identity Layer**

Generates hashed identifiers for each customer, ensuring that no reversible personal identity is stored.

- 3. Feature Engineering Module**

Transforms raw visit events into meaningful behavioral features.

- 4. Prediction and Analytics Layer**

Uses interpretable machine learning models to estimate customer revisit probability

and provide insights.

The system is designed for simplicity, transparency, and deployability in small retail environments. While not fully deployed, the architecture is feasible to implement using lightweight, mobile-friendly tools suitable for low-resource contexts.

4. Data Representation and Features

Each customer is represented using an anonymized hashed identifier. The system relies exclusively on event-based behavioral data, avoiding personal attributes.

Key features include:

- **Visit Frequency:** Number of visits within a defined time window
- **Recency:** Time elapsed since the last visit
- **Inter-Visit Gaps:** Average time between visits
- **Visit Consistency:** Regularity of visits over time
- **Temporal Patterns:** Day-of-week or time-based visit trends

These features are chosen for their predictive relevance and ethical feasibility.

5. Machine Learning Models

Customer return prediction is formulated as a supervised learning problem where the objective is to estimate the probability of a customer revisiting within a future time window.

Proposed models include:

- **Logistic Regression:** For probabilistic interpretation and feature transparency
- **Decision Trees:** For rule-based explanations understandable by non-technical users

Model selection prioritizes interpretability and computational efficiency over complex architectures such as deep neural networks.

6. Privacy-Preserving Approach

Privacy preservation is a foundational principle of the system. The design ensures that:

- No personally identifiable information is stored
- Customer identities are represented only through hashed identifiers
- Analytics operate on aggregated behavioral data
- Outputs focus on probabilities and trends rather than individual profiling

This approach aligns with responsible AI practices and reduces ethical risks in small-scale deployments.

7. Experiments and Evaluation Methodology

Due to the absence of real-world deployment and live datasets, evaluation is conducted conceptually using simulated behavioral patterns.

Evaluation criteria include:

- Logical consistency of predictions
- Interpretability of model outputs
- Actionability of insights for small business owners

Standard machine learning metrics such as accuracy, precision, recall, and probability calibration are proposed for future empirical evaluation.

8. Expected Results and Reasoning

The expected outcome of the system is the ability to:

- Identify customers with a higher likelihood of return
- Provide interpretable explanations for predictions
- Enable data-driven decisions without compromising privacy

The system prioritizes ethical utility and real-world feasibility over maximum predictive accuracy.

9. Limitations

Key limitations include:

- Lack of real-world deployment
- Absence of large-scale empirical training
- Simplified modeling assumptions

These limitations reflect scope and resource constraints rather than conceptual weaknesses.

10. Future Work

Future extensions include:

- Real-world deployment with explicit customer consent
 - Collection of larger anonymized datasets
 - Exploration of temporal and probabilistic sequence models
 - Integration of advanced privacy techniques such as differential privacy
-

11. Conclusion

This research demonstrates that meaningful customer behavior analytics can be designed without invasive data collection or advanced infrastructure. By combining privacy-preserving identity mechanisms with interpretable machine learning models, small retail businesses can gain predictive insights while respecting ethical constraints. The project highlights the importance of responsible, accessible AI design for real-world, low-resource environments.

Declaration of Originality

This project is an independent undergraduate-level research work conducted by the author. All system design decisions, modeling logic, and analysis were developed for academic exploration purposes. No real customer data was collected or used.