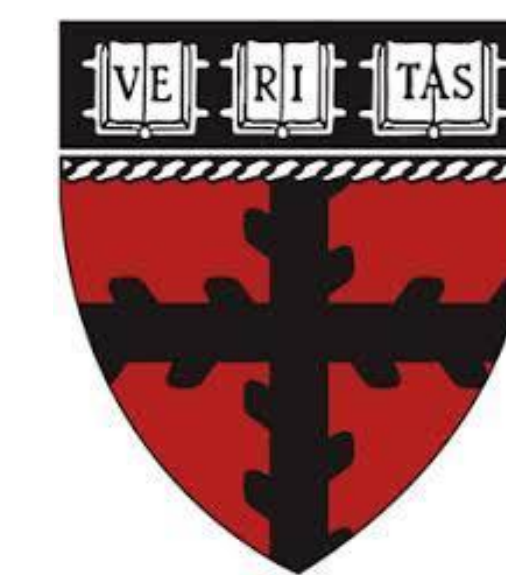# Differential Privacy on Menstruation Data

Nathan Dennis, Raima Islam, Lia Zheng, Shibani Rana

## Introduction

- Period tracking apps collect highly sensitive personal data such as age, BMI, ethnicity, cycle length.
- Vulnerable to privacy breaches.
- Our goal: Explore how Differential Privacy can protect individuals in menstruation datasets while preserving data utility.
- Uses cases:
  - Users
  - Third Parties (researchers, advertisers)
- Methods:
  - DP aggregates (mean, histogram) with Laplace noise.
  - Predictive Deep Learning models with DP-SGD (Adam variant).
- Test multiply privacy budgets, ε, to assess the privacy-accuracy trade off.
- Open doors for further advancements in privacy preserving techniques in the healthcare industry
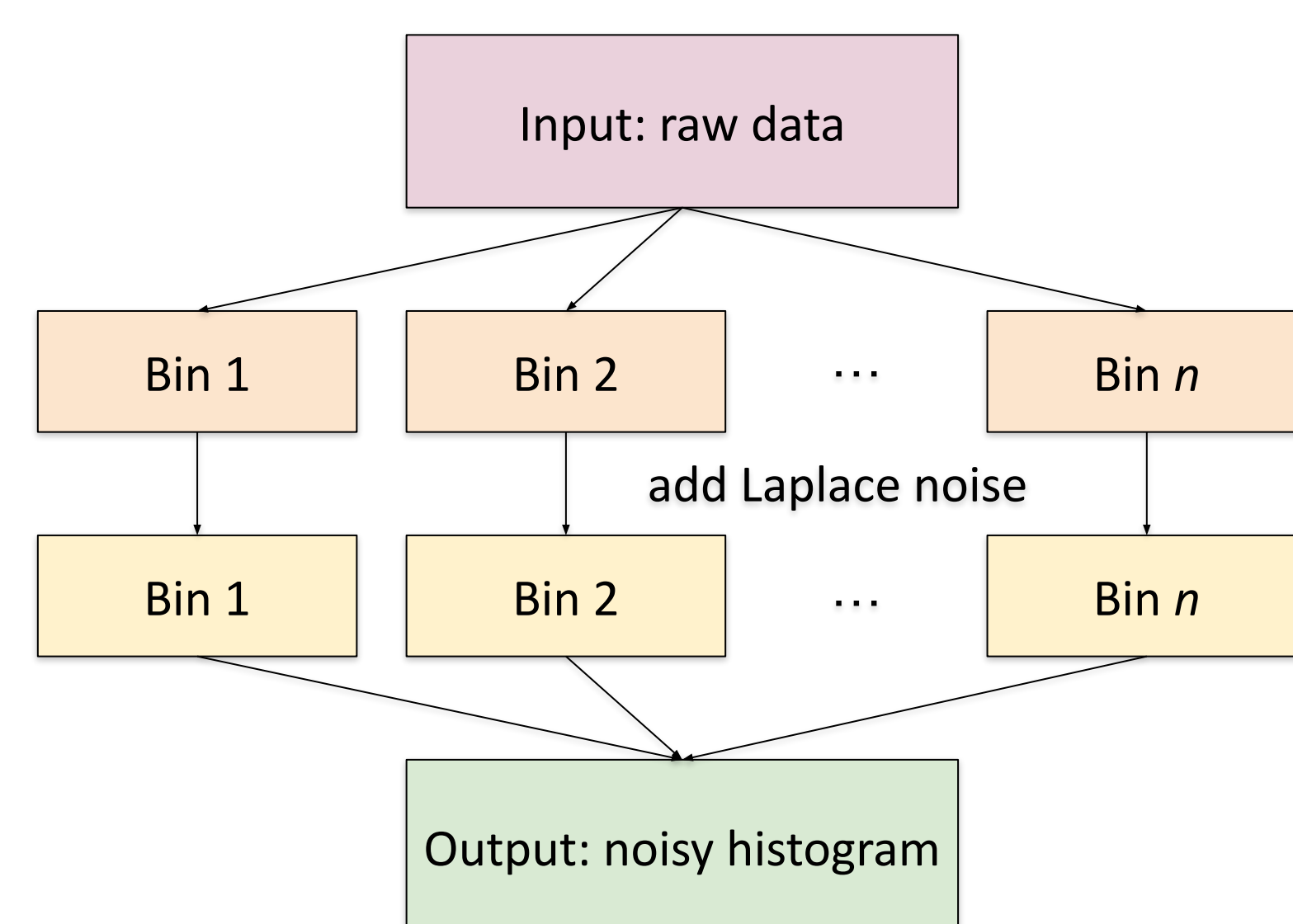
## Background

- **Differential Privacy (DP)** adds noise to data to prevent identification of individuals.
- **2020 U.S Census:** Used DP to protect responses while maintaining data utility.
- **Healthcare:** DP used to aggregate statistics (genomics, wearables)
- **DP-SGD:** Enables private model training with good accuracy dependent on ε.
  - **Dopamine:** Combines DP-SGD + federated learning for private medical ML.
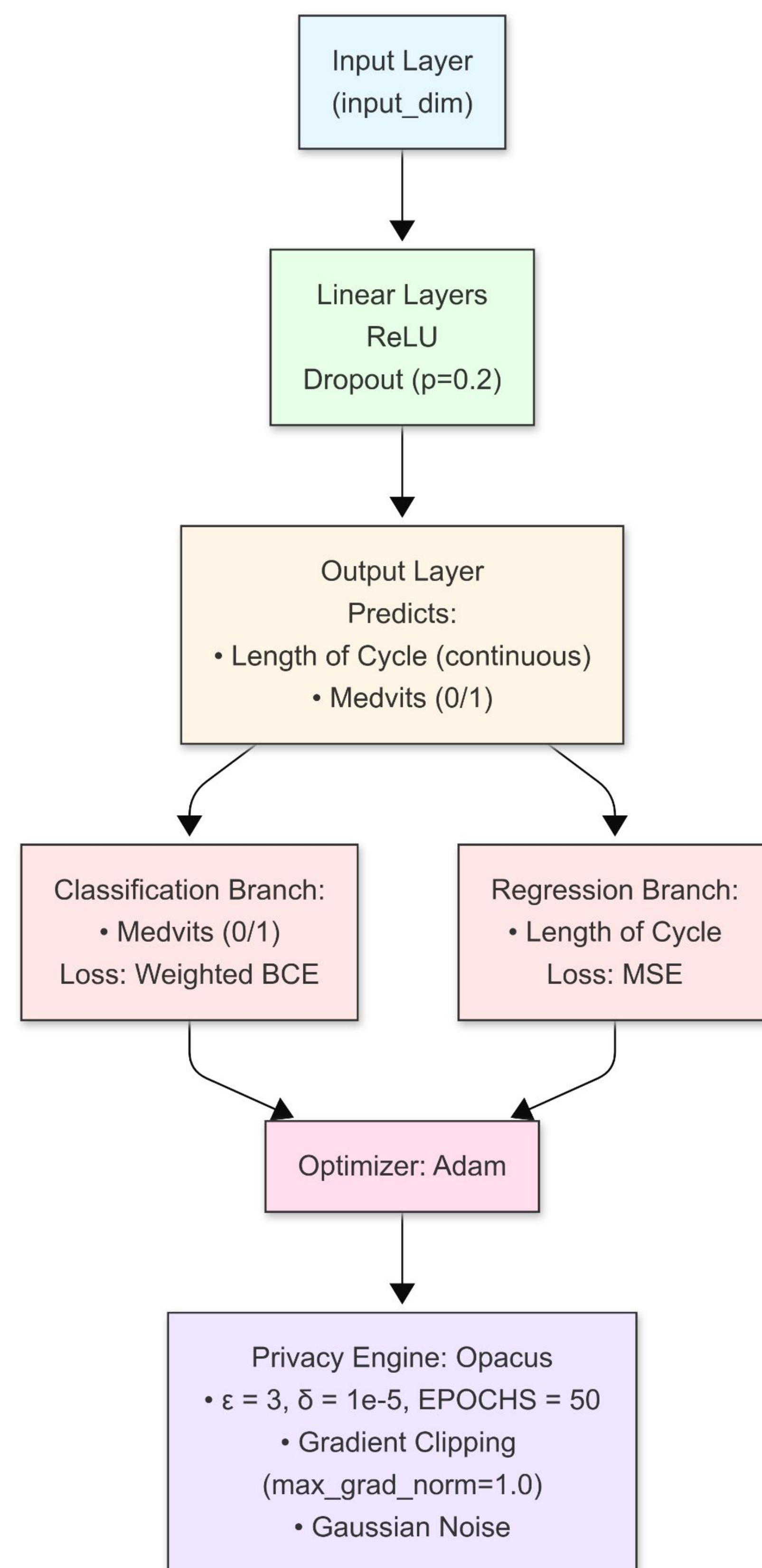- **Gap:** Limited DP research on menstruation data.

## Data

- "**Menstrual Cycle Data**" from 2012 randomized clinical trial conducted by Richard J. Fehring at Marquette University.
- Women tracking cycles in fertility study.
- Key variables: Age, ethnicity, BMI, fertility indicators, cycle length.
- Underwent extensive feature selection process using domain expertise, literature review, exploratory data and statistical analysis
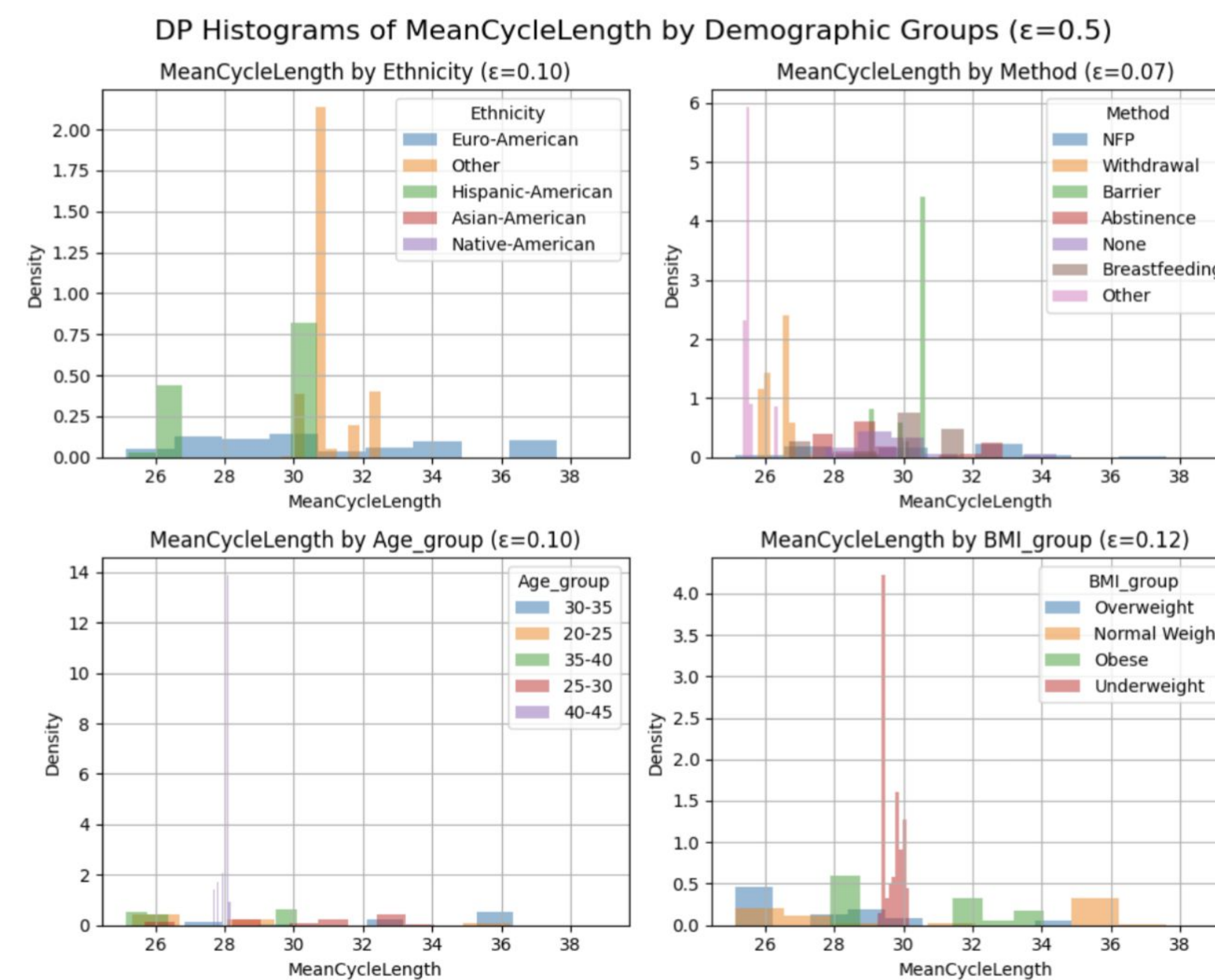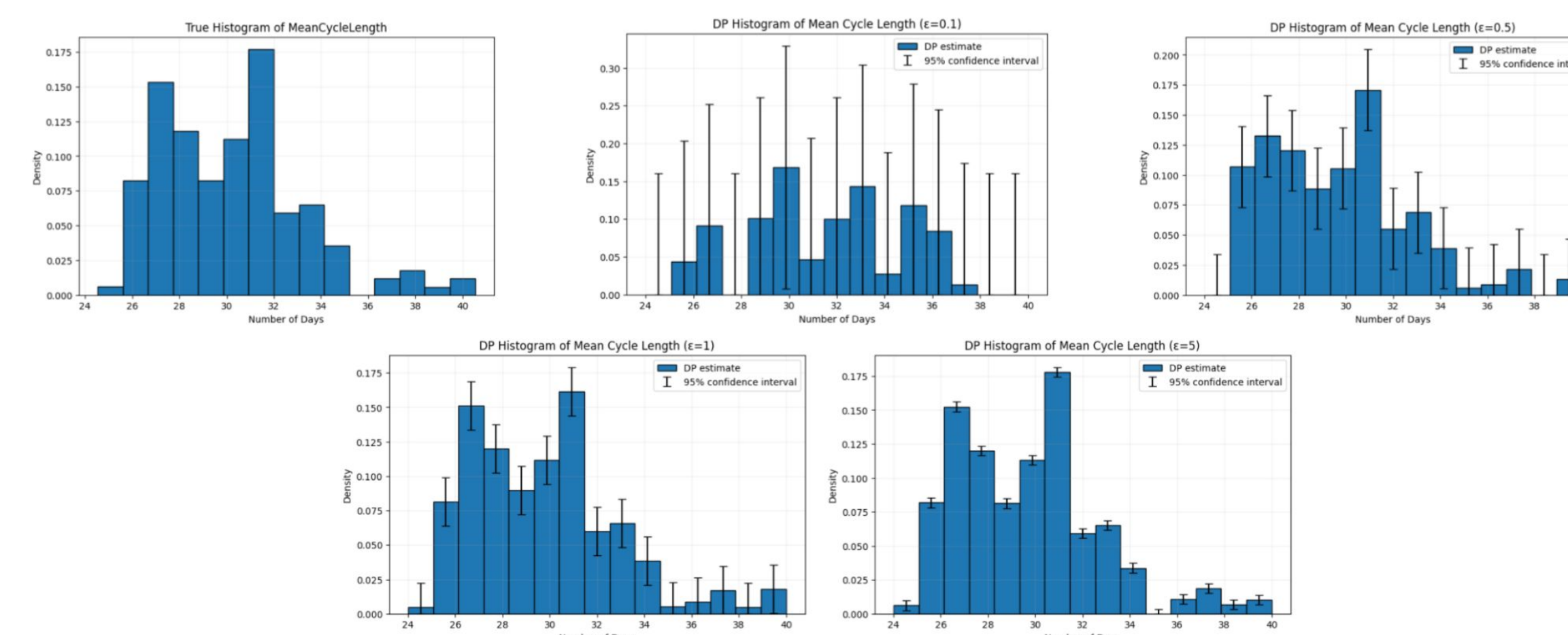
## Methodology

### Private Aggregates



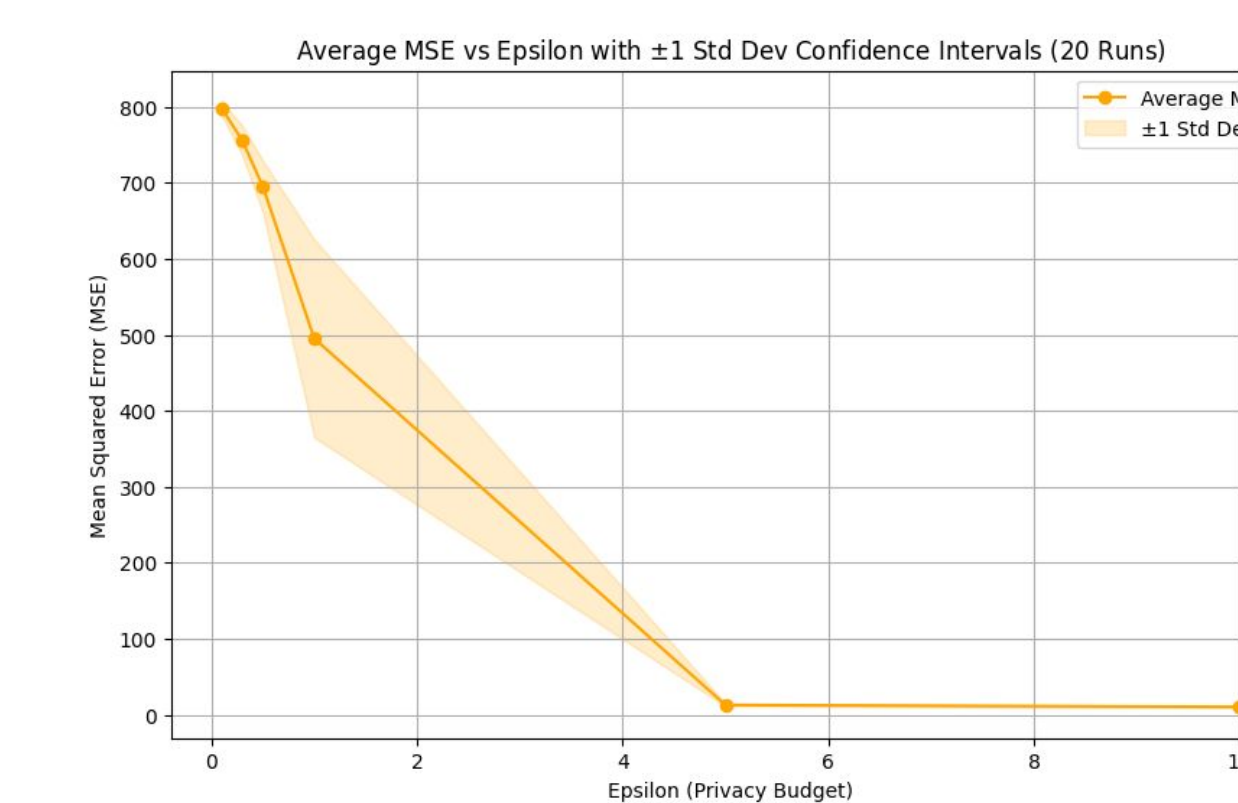Input: raw data → Bin 1, Bin 2, … Bin n → add Laplace noise → Bin 1, Bin 2, … Bin n → Output: noisy histogram

### Predictive Modeling: DP-SGD



Input Layer (input_dim)
Linear Layers ReLU Dropout (p=0.2)
Output Layer Predicts:
• Length of Cycle (continuous)
• Medvits (0/1)

Classification Branch:
• Medvits (0/1)
Loss: Weighted BCE

Regression Branch:
• Length of Cycle
Loss: MSE

Optimizer: Adam

Privacy Engine: Opacus
• ε = 3, δ = 1e-5, EPOCHS = 50
• Gradient Clipping (max_grad_norm=1.0)
• Gaussian Noise

## Results

### Private Aggregates



DP Histograms of MeanCycleLength by Demographic Groups (ε=0.5)



### Predictive Modeling

**Regression**



**Classification**



- Smaller ε increases privacy but leads to low accuracy.
- With lower ε, accuracy results become more variable and less predictable.
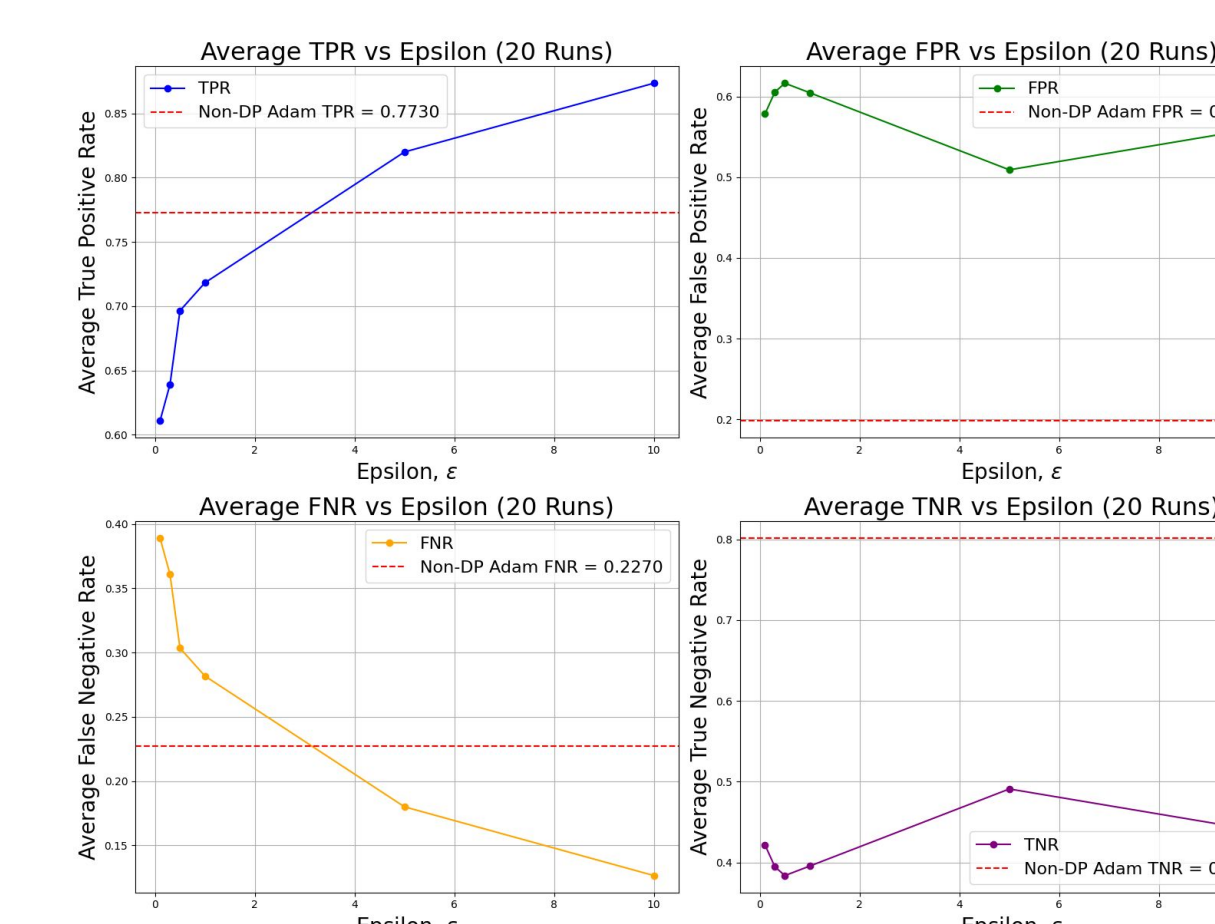


- High MSE/RMSE at low ε (0.1–0.5); unstable performance.
- Improved utility at ε ≥ 5; near-baseline regression results.
- Low ε values show wide confidence intervals, indicating high variance and instability
- Non-DP baseline outperforms all DP settings.

- High cost of missing true positives, lower ε too costly.
- DP models increases false positives, leading to inaccurate health targeting and reduced advertising effectiveness.

## Discussion and Conclusion

- **Privacy-Accuracy Trade-off:** Stronger DP settings significantly reduce model accuracy and stability in both regression and classification tasks.
- **DP-Aggregates:** ε ≥ 1 worked quite well for preserving accuracy of the histograms. However, stratifying by demographic introduced more trade-offs due to division of the privacy budget and class imbalance.
- **DP-SGD Limitations:** Models trained showed higher FPRs and variability, raising concerns for health applications.
- **Practicality:** Moderate ε values (e.g., ≥10) offer improved balance for privacy and utility, especially for our use-case.
- **Data Challenges:** Missing values, class imbalance, and demographic underrepresentation affect model generalizability.
- **Tool Gaps:** Current open-source libraries are not robust and lack proper functionality with very limited documentation
- **Future Work:** Better data collection, synthetic-data generation and hybrid privacy-preserving techniques