

# Differentially Private Machine Learning for Breast Cancer Classification

Melinda Goda and Kaylan Sheally, Concord University  
Mentored by Dr. Abdur Rahman Bin Shahid



## 1. Problem

- Breast cancer among women is one of the most common and deadliest cancers worldwide.
- Recent advancement in ML is helping to develop efficient and effective intelligent systems for the early detection of breast cancers.
- Privacy vulnerability:** From model training to model deployment, privacy leakage can occur at any step in the lifecycle of machine learning.
- Therefore, protecting users' privacy is highly important in breast cancer classification, and very little works have been done to meet this requirement.
- Differential privacy (DP)-based approaches attempt to add statistical noise drawn from a probability distribution (e.g. Laplace distribution) to classifiers.

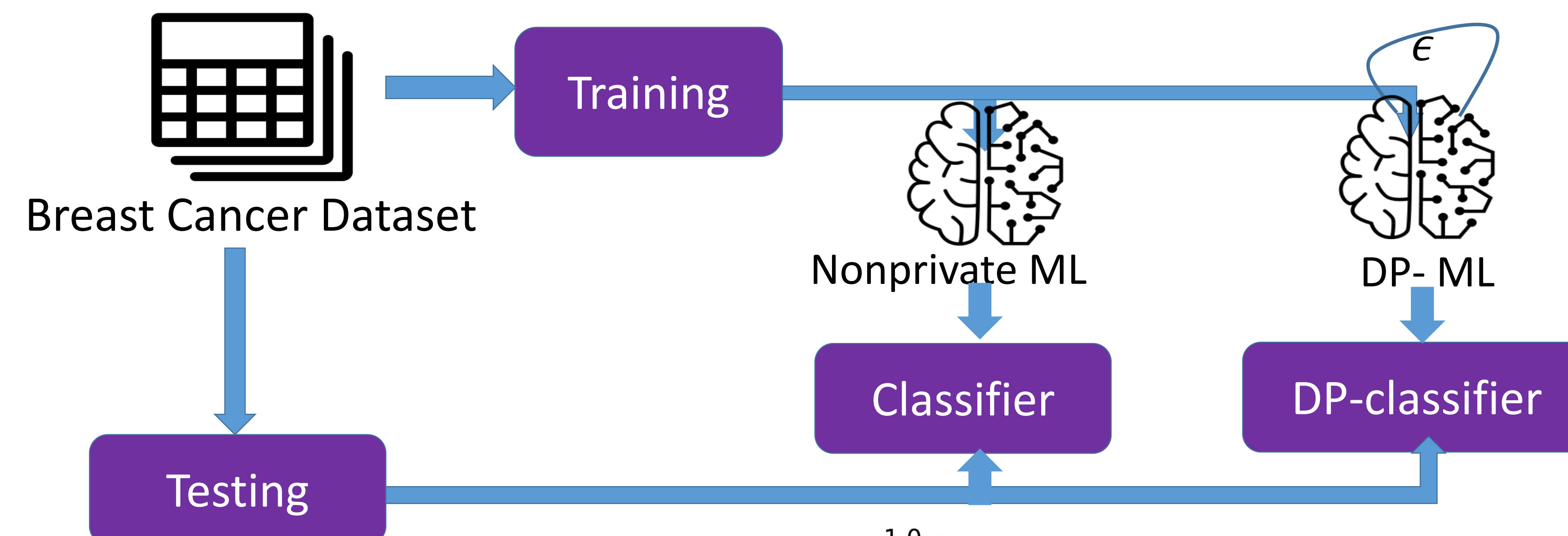
## 2. Contribution

- We present the results of our research on developing differentially-private machine learning models for breast cancer classification.
- We implemented privacy-preserving Logistic Regression and Naïve Bayes in breast cancer classification and compare them with non-private Logistic Regression and Naïve Bayes algorithms.

## 3. ML and DP-ML Models

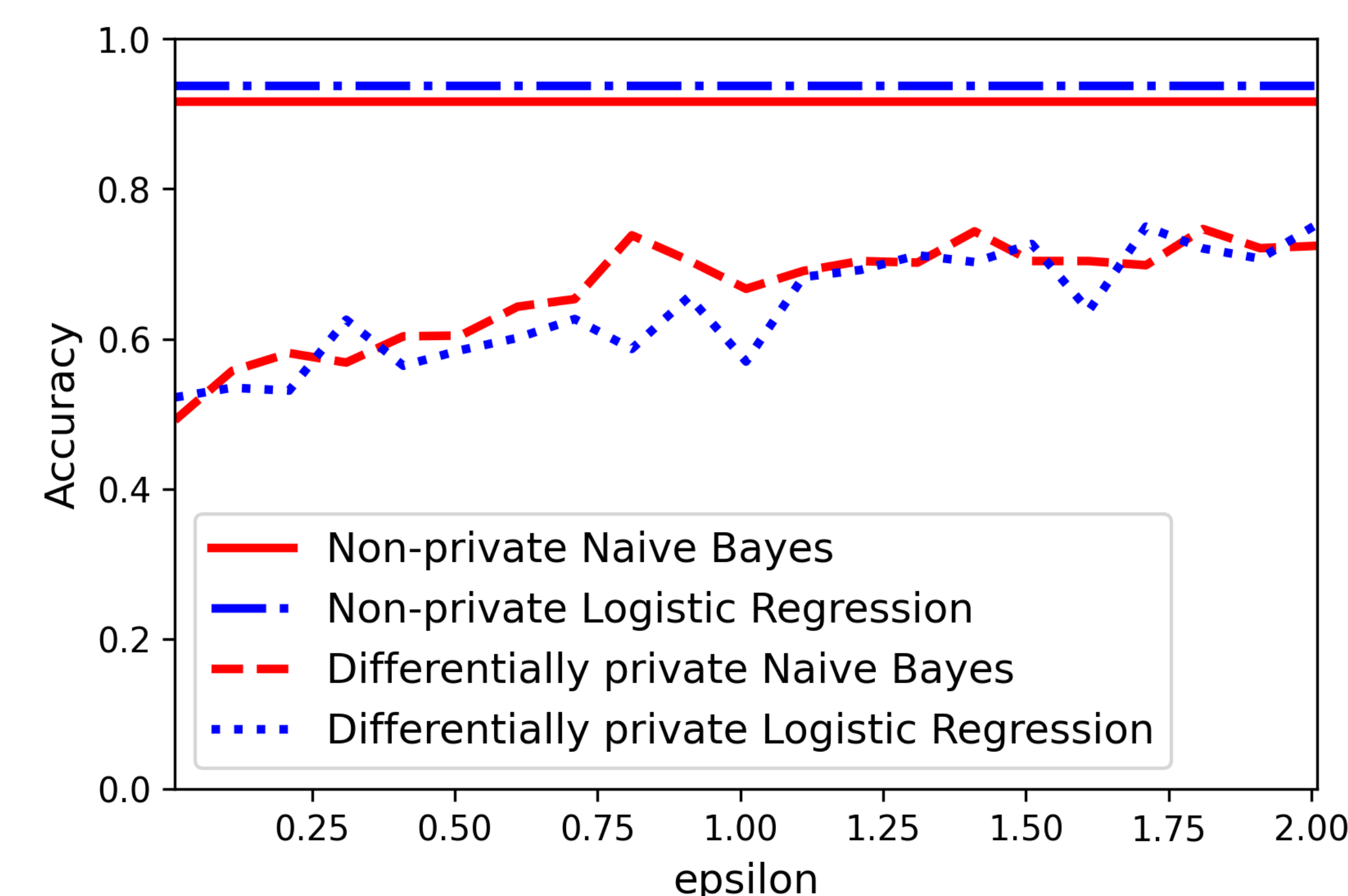
- Logistic Regression (LR)** is a model mainly used for binary classification which computes the weighted sum of the input features and output the logistic of the results.
- The Naïve Bayes (NB)** is a probabilistic classifier based on Bayes' theorem with the assumption of independence between the features.
- Differentially Private Logistic Regression (DP-LR):** Chaudhuri et al. [1] designed a privacy preserving logistic regression approach by introducing differential privacy to perturb the objective function.
- Differentially Private Naïve Bayes (DP-NB):** Vaidya et al. [2] proposed a differentially-private Naïve Bayes model where noises drawn from Laplace distribution are added to the mean and standard deviation of each attribute.

## 4. Methodology



## 5. Experimental Result

- In the experiment, we used the popular Wisconsin Breast Cancer Dataset [3].
- We implemented four different models: Logistic Regression, Naïve Bayes, differentially-private Logistic Regression, and differentially-private Naïve Bayes.
- We used the scikit-learn and IBM diffPrivLib libraries.



- The results show that it is possible to achieve high accuracy with both privacy-preserving models.

## 6. Lesson Learned, Conclusion, and Future Work

- Privacy concerns can be a big hurdle in developing intelligent systems to perform sensitive tasks.
- We have developed a differentially private Naïve Bayes and Logistic Regression classifiers for breast cancer classification.
- We have tested both classifiers on a real world dataset and results show that it is possible to achieve high accuracy with both models, compared to baseline models.
- In the future, we plan to look at how different feature engineering methods, such as data augmentation, can improve the accuracy of the models.

## 7. References

- Farzad Zafarani and Chris Clifton. Differentially private naïve bayes classifier using smooth sensitivity. arXiv preprint arXiv:2003.13955, 2020.
- Kamalika Chaudhuri, Claire Monteleoni, and Anand D Sarwate. Differentially private empirical risk minimization. Journal of Machine Learning Research, 12(3), 2011.
- Dua, D. and Graff, C. (2019). UCI Machine Learning Repository [http://archive.ics.uci.edu/ml]. Irvine, CA: University of California, School of Information and Computer Science.