



UMEA UNIVERSITET

# Navigating Data Privacy and Utility: A Strategic Perspective

**Saloni Kwatra**

## Akademisk avhandling

som med vederbörligt tillstånd av Rektor vid Umeå universitet  
för avläggande av teknologie doktorsexamen framläggs till  
offentligt försvar i BIO.A.206 Aula Anatomica, Biologihuset,  
den 4 November 2024, kl. 09:15.

Avhandlingen kommer att försvaras på engelska.

Fakultetsopponent: Professor Sébastien Gambs vid Département  
d'informatique, Université du Québec (Kanada)

Department of Computing Science

**Organization**

Umeå University  
Dept. of Computing Science

**Document type**

Doctoral thesis

**Date of publication**

14 October 2024

**Author**

Saloni Kwatra

**Title**

Navigating Data Privacy and Utility: A Strategic Perspective

**Abstract**

Privacy in machine learning should not merely be viewed as an afterthought; rather, it must serve as the foundation upon which machine learning systems are designed. In this thesis, along with the centralized machine learning, we also consider the distributed environments for training machine learning models, particularly federated learning. Federated learning lets multiple clients or organizations train a machine learning model in a collaborative manner without moving their data. Each client participating to the federation shares the model parameters learnt by training a machine learning model on its data. Even though the setup of federated learning keeps the data local, there is still a risk of sensitive information leaking through the model updates. For instance, attackers could potentially use the updates of the model parameters to figure out details about the data held by clients. So, while federated learning is designed to protect privacy, it still faces challenges in ensuring that the data remains secure throughout the training process. Originally, federated learning was introduced in the context of deep learning models. However, this thesis focuses on federated learning for decision trees. Decision Trees are intuitive, and interpretable models, making them popular in a wide range of applications, especially where explainability of the decisions made by the decision tree model is important. However, Decision Trees are vulnerable to inference attacks, particularly when the structure of the decision tree is exposed. To mitigate these vulnerabilities, a key contribution of this thesis is the development of novel federated learning algorithms that incorporate privacy-preserving techniques, such as k-anonymity and differential privacy, into the construction of decision trees. By doing so, we seek to ensure user privacy without significantly compromising the performance of the model.

Machine learning models learn patterns from data, and during this process, they might leak sensitive information. Each step of the machine learning pipeline presents unique vulnerabilities, making it essential to assess and quantify the privacy risks involved. One focus of this thesis is the quantification of privacy by devising a data reconstruction attack tailored to Principal Component Analysis (PCA), a widely used dimensionality reduction technique. Furthermore, various protection mechanisms are evaluated in terms of their effectiveness in preserving privacy against such reconstruction attacks while maintaining the utility of the model. In addition to federated learning, this thesis also addresses the privacy concerns associated with synthetic datasets generated by models such as generative networks. Specifically, we perform an Attribute Inference Attack on synthetic datasets, and quantify privacy by calculating the Inference Accuracy—a metric that reflects the success of the attacker in estimating sensitive attributes of target individuals.

Overall, this thesis contributes to the development of privacy-preserving algorithms for decision trees in federated learning and introduces methods to quantify privacy in machine learning systems. Also, the findings of this thesis set a ground for further research at the intersection of privacy, and machine learning.

**Keywords**

Privacy, Data Reconstruction Attacks, k-anonymity, Differential Privacy, Federated Learning, Decision Trees, and, Principal Component Analysis

**Language**

English

**ISBN**

Print: 978-91-8070-481-6  
PDF: 978-91-8070-482-3

**ISSN**

0348-0542

**Number of pages**

103