



KTH Stockholm
Department of Numerical Analysis

Privacy-Preserving Machine Learning

lorem ipsum upsum

Sijun John Tu

Master Thesis Report

Supervisors: Anders Szepessy (KTH) and Shahid Raza (RISE)

Abstract

STATUTORY DECLARATION

I declare that I have authored this thesis independently, that I have not used other than the declared sources / resources, and that I have explicitly marked all material which has been quoted either literally or by content from the used sources.

wOctober 23, 2023

Sijun John Tu

Contents

Notation

Mathematical conventions and notation used in this thesis:

\mathbb{R} the real numbers

Additionally, we introduce the following conventions to describe various elements from different mathematical objects to make the notations and their meaning as consistent as possible:

a, b, c scalar values

A, B, C matrices

List of Figures

1 Introduction

2 Theoretical background on Differential Privacy

In this chapter we briefly describe and derive the most important results from Cynthia Dwork's work on differential privacy that was first introduced in [1]. This summary heavily relies on her writings in her as well as lecture notes from [2].

2.1 Defining differential privacy

Differential privacy (DP) should be understood as an agreement between the data holder and the data subject: the latter should not be “affected, adversely or otherwise, by allowing [her] data to be used in any study or analysis, no matter what other studies, data sets or information sources are available” [3]. This addresses the paradox of learning something useful about a population while learning nothing about the individuals

Example 2.1.1 (Randomised response). [4] proposes the following random answering procedure: In a study where participants are asked to answer with “Yes” or “No” whether they have engaged in an illegal or embarrassing activity A , they should:

1. Flip a coin
2. If the coin shows tails, then the participant should respond truthfully.
3. If the coin shows head, then the participant should flip the coin a second time and answer “Yes” if the second coin shows head and “no” otherwise.

This procedure ensures participants' privacy by “plausible deniability”; each participant's answer has non-zero probability of being truthful or not. By understanding the probabilities of the noise generation process, the data analyst can estimate the true number of “yes” and “no” answers. To this end, let p be the true percentage of “yes” answers, N the total number of participants, n_{true} the true number of “yes” responses and \hat{n}_{obs} the observed number of “yes” responses. We assume a fair coin with equal probability of showing heads or tails. Then the expected number of “yes” answers after applying the described procedure is:

$$\mathbb{E}(\text{"Yes"}) = \frac{1}{4}n_{true} + \frac{1}{4}(N - n_{true}) + \frac{1}{2}n_{true} = \frac{1}{4}N + \frac{n_{true}}{2} \quad (1)$$

We can estimate this using the $\hat{n}_{obs} \approx \mathbb{E}(\text{"Yes"}) = \frac{1}{4}N + \frac{n_{true}}{2}$ and finally solving for n_{true} yields the estimate:

$$n_{true} = 2\hat{n}_{obs} - \frac{1}{2}N \quad (2)$$

Definition 2.1.2 (Probability Simplex). Given a discrete set B , the probability simplex over B is defined as the set

$$\Delta(B) = \left\{ x \in \mathbb{R}^{|B|}, x_i \geq 0 \text{ and } \sum_i x_i = 1 \right\} \quad (3)$$

Definition 2.1.3. (Randomised Algorithm) A randomized algorithm \mathcal{M} with domain A and discrete range B is associated with a mapping $M : A \rightarrow \Delta(B)$. On input $a \in A$ algorithm \mathcal{M} outputs $\mathcal{M}(a) = b$ with probability $(M(a))_b$

Definition 2.1.4 (Histogram representation of a data base). Given a set \mathcal{X} , the universe of all possible records“, the histogram representation of a database x is the vector

$$x \in \mathbb{N}^{|\mathcal{X}|} \quad (4)$$

in which each entry x_i represents the number of elements in database x of type $i \in \mathcal{X}$.

Definition 2.1.5 (l_1 -norm of a database in histogram representation). The l_1 -norm of a database is a measure of the size of the database and defined as:

$$\|x\|_1 = \sum_{i=1}^{|\mathcal{X}|} |x_i| \quad (5)$$

This immediately gives rise to a notion of distance between two databases x and y , namely:

$$\|x - y\|_1 = \sum_{i=1}^{|\mathcal{X}|} |x_i - y_i| \quad (6)$$

which basically counts the number of different entries.

Now we are ready to give the general definition of differential privacy:

Definition 2.1.6 ((ϵ, δ) -DP). A randomised algorithm \mathcal{M} with domain $\mathbb{N}^{|\mathcal{X}|}$ is (ϵ, δ) -differentially private if for all outcomes $S \subset \text{ran}\mathcal{M}$ and for all databases $x, y \in \mathbb{N}^{|\mathcal{X}|}$, such that $\|x - y\|_1$ (i. e. they only differ in one element) we have

$$\mathbb{P}(\mathcal{M}(x) \in S) \leq e^\epsilon \cdot \mathbb{P}(\mathcal{M}(y) \in S) + \delta \quad (7)$$

where the probability is taken over the randomness of \mathcal{M} . If $\delta = 0$, we say \mathcal{M} is ϵ -differentially private.

why e^ϵ

Example 2.1.7 (Randomised response revisited).

2.2 Important results for Differential Privacy

Theorem 2.2.1 (DP requires randomisation). *Any non-trivial DP-mechanism requires randomisation.*

Proof. TBA □

Theorem 2.2.2 (Post-processing). *Let $\mathcal{M} : \mathbb{N}^{|\mathcal{X}|} \rightarrow R$ be a randomised algorithm that is (ϵ, δ) -DP. Further let $f : R \rightarrow R'$ an arbitrary function. Then $f \circ \mathcal{M}$ is also (ϵ, δ) -DP.*

Proof. First fix data sets $x, y \in \mathbb{N}^{|\mathcal{X}|}$, s. t. $\|x - y\|_1 \leq 1$ and outcome $S' \subseteq R'$. Define a set $S = \{r \in R : f(r) \in S'\}$. Then we have:

$$\begin{aligned} \mathbb{P}(f(\mathcal{M}(x)) \in S') &= \mathbb{P}(\mathcal{M}(x) \in S) \\ &\leq e^\epsilon \cdot \mathbb{P}(\mathcal{M}(y) \in S) + \delta \\ &= e^\epsilon \cdot \mathbb{P}(f(\mathcal{M}(y)) \in S') + \delta \end{aligned} \tag{8}$$

where the inequality follows from the (ϵ, δ) - DP of \mathcal{M} . \square

Theorem 2.2.3 (Group privacy). *Let $\mathcal{M} : \mathbb{N}^{|\mathcal{X}|} \rightarrow R$ be a randomised algorithm that is (ϵ, δ) - DP, then \mathcal{M} is $(k\epsilon, ke^{k\epsilon}\delta)$ - DP for groups of size k , i. e. it holds for databases $x, y \in \mathbb{N}^{|\mathcal{X}|}$ such that $\|x - y\|_1 \leq k$ and for all $S \subseteq R$:*

$$\mathbb{P}(\mathcal{M}(x) \in S) \leq e^{k\epsilon} \cdot \mathbb{P}(\mathcal{M}(y) \in S) + k\delta \tag{9}$$

Proof. First fix data sets $x, y \in \mathbb{N}^{|\mathcal{X}|}$, s. t. $\|x - y\|_1 \leq k$ and outcome $S \subseteq R$. Now there exists a series of databases z_0, \dots, z_k , such that $x = z_0$ and $y = z_k$ and $\|z_{i+1} - z_i\|_1 \leq 1$, i. e. we can find a series of databases that transforms x into y by removing or adding one record at a time. Then we have:

$$\begin{aligned} \mathbb{P}(\mathcal{M}(x) \in S) &= \mathbb{P}(\mathcal{M}(z_0) \in S) \\ &\leq e^\epsilon \cdot \mathbb{P}(\mathcal{M}(z_1) \in S) + \delta \\ &\leq e^\epsilon (e^\epsilon \cdot \mathbb{P}(\mathcal{M}(z_2) \in S) + \delta) + \delta \\ &\leq \dots \\ &= ke^\epsilon \cdot \mathbb{P}(\mathcal{M}(y) \in S) + ke^{k\epsilon}\delta \end{aligned} \tag{10}$$

\square

Theorem 2.2.4 (Standard composition). *Let $\mathcal{M}_1 : \mathbb{N}^{|\mathcal{X}|} \rightarrow R_1$ and $\mathcal{M}_2 : \mathbb{N}^{|\mathcal{X}|} \rightarrow R_2$ be two randomised algorithms that are (ϵ_1, δ_1) - and (ϵ_2, δ_2) DP, then their composition defined by $\mathcal{M}_{12} : \mathbb{N}^{|\mathcal{X}|} \rightarrow R_1 \times R_2$, $\mathcal{M}_{12}(x) = (\mathcal{M}_1(x), \mathcal{M}_2(x))$ is $(\epsilon_1 + \epsilon_2, \delta_1 + \delta_2)$ DP.*

Proof. TBA \square

2.3 Example of DP-mechanism: Laplace mechanism

Appendix