# Invetigation of the Elliptical Gaussian Noise in the case of multivariate normal data

Tim Sehested Poulsen

December 5, 2022

# 1 Preliminaries

## 1.1 Definitions

**The dataset** is denoted by $X$, where $X \in \mathbb{R}^{n \times d}$. I have that $n$ denotes the number of entries in the dataset and $d$ is the number of dimensions of the dataset. I will throughout the report refer to a single entry of the dataset as $x_i$ and a single dimension of the dataset as $X^{(j)}$, and therefore $x_i^{(j)}$ denotes the $j$'th dimension of the $i$'th entry.

**Differential privacy** is the heuristic of releasing a database statistic whilst limiting the impact of any one entry. It builds on the intuition that computing a statistics on a private dataset should not reveal any sensitive information about any one individual as long as that individual has little to no effect on the outcome. Differential privacy has multiple slightly different formal definitions, one such is $(\varepsilon, \delta)$-Differential Privacy refered to as $(\varepsilon, \delta)$-DP which will be introduced later on. A prerequisite for almost all of the different differential privacy definitions relies on the concept of neighbouring dataset.

**Definition 1.1** (Neighbouring dataset [2]). *Two dataset $X, X' \in \mathbb{R}^{n \times d}$ are said to be neigbouring if they differ in at most a single entry. Neighbouring dataset are denoted with the relation $X \sim X'$ and defined as followed*

$$X \sim X' \iff |\{i \in \mathbb{N} \mid i \leq n \land x_i \neq x_i'\}| \leq 1$$

**Definition 1.2** (Sensitivity [5]). *Let $f(X) : \mathbb{R}^{n \times d} \to \mathbb{R}^d$ given by $f(X) = \sum_{i=1}^n x_i$ be the sum over all vectors in a dataset. The sensitivity is then the maximal possible difference in the output of our summation from two neighbouring dataset denoted as $\Delta$. We denote the sensitivity of the $j$'th dimension as*

$$\Delta_j = \max_{X \sim X'} \left| f(X)^{(j)} - f(X')^{(j)} \right|$$

*and then the total $l_2$-sensitivity is then*

$$\|\Delta\| = \max_{X \sim X'} \|f(X) - f(X')\| \tag{1}$$

**Definition 1.3** (($\varepsilon, \delta$)-Differential Privacy [2]). *A randomized algorithm $\mathcal{M} : \mathbb{R}^{n \times d} \to \mathcal{R}$ is ($\varepsilon, \delta$)-differentially private if for all possible subsets of outputs $S \subseteq \mathcal{R}$ and all pairs of neighbouring dataset $X \sim X'$ we have that*

$$\Pr[M(X) \in S] \leq e^{\varepsilon} \cdot \Pr[M(X') \in S] + \delta$$

**Error Measure**  As I will be working exclusively with the sum of entries in a dataset, error will be defined as the expected squared $l_2$-norm between the true sum and the output of a randomized algorithm. So let $X \in \mathbb{R}^{n \times d}$ be the dataset and $f(X) = \sum_i^n x_i$ be the true sum. The error of a randomized algorithm $M : \mathbb{R}^{n \times d} \to \mathbb{R}^d$ which estimates $f(X)$ is then

$$\mathrm{Err}(M) := \mathbb{E}\left[\|M(X) - f(X)\|^2\right]$$

**extra**  If proof of elliptical ... is included then write about edp is preserved during transformatio

# 2 Algorithms

## 2.1 The Gaussian Mechanism

One of the most foundational algorithms for achieving ($\varepsilon, \delta$)-DP is the Gaussian Mechanism [3]. It computes the real value of a statistic, where the $l_2$-sensitivity is known. That is it produces a ($\varepsilon, \delta$)-DP estimate of a function $g : \mathbb{R}^{n \times d} \to \mathbb{R}^d$ where $\|\Delta\|$ for the function $g$ is known. It does so by computing the value of $g(X)$ and then adding noise to each dimension from the normal distribution $\mathcal{N}(0, \sigma_{\varepsilon,\delta}^2)$. This can be seen as adding a noise vector $\eta$ which is then distributed according to the multivariate normal distribution $\mathcal{N}(\overrightarrow{0}, \sigma_{\varepsilon,\delta}^2 I)$. The algorithm can be seen in Algorithm 1.

---
**Algorithm 1** The Gaussian Mechanism
---
   **Input**
      $\sigma_{\varepsilon,\delta}$             Standard deviation required to achieve ($\varepsilon, \delta$)-DP
      $X \in \mathbb{R}^{n \times d}$     Dataset
   **Output**
      ($\varepsilon, \delta$)-DP estimate of $g(X)$
  $\eta \leftarrow$ sample from $\mathcal{N}(\overrightarrow{0}, \sigma_{\varepsilon,\delta}^2 I)$
  **return** $g(X) + \eta$

---

It is quite apparent that the main difficulty of the mechanism lies in determining a $\sigma_{\varepsilon,\delta}$ which achieves ($\varepsilon, \delta$)-DP , and preferably the smallest such one.

The following theorem was initially proven

**Theorem 1.** [3] *Let $g : \mathbb{R}^{n \times d} \to \mathbb{R}^d$ be an arbitrary d-dimensional function with $l_2$-sensitvity $\|\Delta\| = \max_{X \sim X'} \|g(X) - g(X')\|$, and let $\varepsilon \in (0,1)$. The Gaussian Mechanism with $\sigma_{\varepsilon,\delta} = \|\Delta\| \sqrt{2\ln(1.25/\delta)}/\varepsilon$ is ($\varepsilon, \delta$)-DP .*

The proof is rather long and is therefore ommitted here.

In the case where $g(X) = f(X)$, i.e. it is estimating the sum of entries we have quite intuitively that the error is given by the norm of the noise introduced as

$$\mathbb{E}\left[\|(g(X) + \eta) - f(X)\|^2\right] = \mathbb{E}\left[\|(f(X) + \eta) - f(X)\|^2\right] = \mathbb{E}\left[\|\eta\|^2\right]$$

Which by theorem **??** is

$$\mathbb{E}\left[\|\eta\|^2\right] = \sum_i^d \sigma_{\varepsilon,\delta}^2 = d \cdot \sigma_{\varepsilon,\delta}^2$$

-**Does introduce the same variance for all dimensions regardless of individual variance**

-**Talk about finding the minimal $\sigma$ s.t. privacy is held, and the error in that case**

# 3   Problem setup

The problem consists of realeasing the sum of vectors in a dataset under differential privacy. More formally we whish to release the value of $f : \mathbb{R}^{n \times d} \to \mathbb{R}^d$ given by

$$f(X) = \sum_{i=1}^n x_i$$

under $(\varepsilon, \delta)$-DP .

The common factor for achieving $(\varepsilon, \delta)$-DP in both the Gaussian Mechanism and the Elliptical Gaussian Mechanism is the requirement that data lie within some hyperrectangle. It is formally described as each dimension of the data must lie within some range $x_i^{(j)} \in [-\Delta_j/2, \Delta_j/2]$. This requirement is needed to know the $l_2$-sensitivity $\|\Delta\|$ as defined in equation 1 of the data. In this project I will change this assumption and instead look at the case where each dimension is normally distributed. This means that for each $j \in [d]$ we have that $X^{(j)} \sim \mathcal{N}(\mu_j, \sigma_j^2)$. An equivalent formulation is that a the data is multivariately distributed but with no correlation between dimensions. This means that $X \sim \mathcal{N}(\boldsymbol{\mu}, \Sigma)$, where $\Sigma$ is a diagonal matrix with the variance of each dimension along its diagonal. It is quite apparent that determining a $\|\Delta\|$ is impossible in this setting as the Gaussian distribution is continously defined on the range $(-\infty, \infty)$. What has been done in several recent papers is that data is *clipped* by some threshhold $C$ [1,4]. Clipping is the process of limiting the norm of any one entry to be at most $C$. This means that every vector is transformed as such

$$\hat{x}_i := \min\left\{\frac{C}{\|x_i\|}, 1\right\} \cdot x_i$$

Clipping entries by a factor $C$ thus means that $\|\Delta\| = C$ as any one entry cannot have more impact on the summation than $C$. It can then be seen that if the summation $f(X)$ is instead performed on a clipped dataset $\hat{X}$ then this is equivalent to defining the summation function $\hat{f} : \mathbb{R}^{n \times d} \to \mathbb{R}^d$ as

$$\hat{f}(X) = \sum_i^n \min\left\{\frac{C}{\|x_i\|}, 1\right\} \cdot x_i$$

Then by theorem 1 the gaussian mechanism with the function $\hat{f}$ is $(\varepsilon, \delta)$-DP with $\sigma_{\varepsilon,\delta} = C\sqrt{2\ln(1.25/\delta)}/\varepsilon$. Though the mechanism is still $(\varepsilon, \delta)$-DP it will now have a larger error when regarding the true sum $f(X) = \sum_i^n x_i$ as the actual answer.

# References

[1] Biswas, S., Dong, Y., Kamath, G., and Ullman, J. Coinpress: Practical private mean and covariance estimation. *Advances in Neural Information Processing Systems 33* (2020), 14475–14485.

[2] Dwork, C., McSherry, F., Nissim, K., and Smith, A. Calibrating noise to sensitivity in private data analysis. *Journal of Privacy and Confidentiality 7*, 3 (2016), 17–51.

[3] Dwork, C., Roth, A., et al. The algorithmic foundations of differential privacy. *Foundations and Trends® in Theoretical Computer Science 9*, 3–4 (2014), 211–407.

[4] Huang, Z., Liang, Y., and Yi, K. Instance-optimal mean estimation under differential privacy.

[5] Pagh, R., and Lebeda, C. Private vector aggregation when coordinates have different sensitivity.