

# **USING KD-LAPLACE TO TRAIN PRIVACY-PRESERVING CLUSTERING ALGORITHMS ON DISTRIBUTED K-DIMENSIONAL DATA**

by

**Tjibbe van der Ende**

in partial fulfillment of the requirements for the degree of

**Master of Science**  
in Software Engineering

at the Open University, faculty of Management, Science and Technology  
Master Software Engineering  
to be defended publicly on Day Month DD, YYYY at HH:00 PM.

Student number: 852372917  
Course code: IMA0002  
Thesis committee: Dr. Ir. Mina Sheikhalishahi (chairman), Open University  
Dr. Ir. Clara Maathuis (supervisor), Open University

# CONTENTS

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Background</b>	<b>4</b>
2.1	Differential privacy . . . . .	5
2.2	Clustering . . . . .	11
2.3	Literature review . . . . .	17
<b>3</b>	<b>nD-Laplace</b>	<b>28</b>
3.1	2D-Laplace . . . . .	28
3.2	3D-Laplace . . . . .	33
3.3	kD-Laplace . . . . .	39
<b>4</b>	<b>Attacks on privacy</b>	<b>52</b>
4.1	Membership inference attacks . . . . .	52
4.2	Reconstruction attack . . . . .	54
4.3	Attack evaluation . . . . .	55
<b>5</b>	<b>Methodology</b>	<b>57</b>
5.1	Research design . . . . .	57
5.2	Data collection . . . . .	61
5.3	Data analysis . . . . .	67
<b>6</b>	<b>Results</b>	<b>68</b>
6.1	Cluster utility . . . . .	68
6.2	Mechanism utility . . . . .	81
6.3	Privacy . . . . .	87
<b>7</b>	<b>Discussion</b>	<b>99</b>
<b>8</b>	<b>Conclusion</b>	<b>100</b>
<b>Bibliography</b>		<b>i</b>
<b>Hyperparameters</b>		<b>xi</b>
.1	K-Means . . . . .	xi
.2	Agglomerative clustering . . . . .	xiii
<b>Theory</b>		<b>xx</b>
.3	Big O Notation . . . . .	xx
<b>Results</b>		<b>xxi</b>
.4	Mechanism utility . . . . .	xxi
.5	Privacy . . . . .	xxvii

# 1

## INTRODUCTION

check all sentences to have verbs

Data has become an increasingly important part of our daily lives, and its applications are continuously growing. Personal and sensitive information, such as scrolling through Facebook or Instagram, is often sourced from various structured and unstructured data origins. As many businesses see the potential of data [noa, b], the urgency of protecting sensitive information also increases. This development is why there have been several regulations regarding data and privacy. Therefore, the European Union created the GDPR [noa, a], and the United States created the privacy act [Division, 2007]. For instance, the GDPR describes in Article 34 that data should always be stored in an "incomprehensible way for unauthorized people". Due to these regulations, it is challenging for companies to store and analyze data to gain insights.

A method for gaining insights into data is using unsupervised machine learning. It can be used to find a structure or pattern in uncategorized data [Dridi, 2021]. Clustering is a type of unsupervised learning commonly applied to group similar data points and reveal patterns and structures [Dridi, 2021]. Clustering algorithms need to be trained on an adequate amount of historical data, so the data needs to be stored and processed in an accessible location.

A common solution for hiding data from unauthorized access is encryption. Although this method is widely used for databases, it is unsuitable for data analysis. The data is unreadable for examination, and it adds communication overhead [Liu et al., 2022]. For this reason, researchers have explored methods for anonymizing data to protect the privacy of individuals [Sweeney, 2002]. But, this can still reveal much information about a person, as the range of possible values is limited. A more effective way to preserve privacy is by using differential privacy. This mechanism adds noise to the original data to hide the actual value, controlled by a privacy budget [Dwork, 2006]. Due to the configurability of this mechanism, it is possible to trade off privacy and utility in such a way it is possible to preserve data patterns. In recent years, the method was extended to add noise locally before sending it to a central server and called **Local Differential Privacy (LDP)** [Kasiviswanathan et al., 2010]. Although this method reduces privacy leakage risks because no plain data is stored on the server, the amount of data is limited in the local view. This makes the noise more significant, impacting the data utility [Yang et al., 2020].

## RELATED WORK

Preserving utility in clustering while ensuring privacy is a well-known challenge in the existing literature. Previous studies have attempted to address this issue, with some focusing on preserving cluster shapes by sending feature distance information to the server rather than the feature values [Sun et al., 2019]. However, this method was found to leak distance information, compromising privacy.

Another line of research focused on improving privacy preservation using a distance-aware randomization algorithm proposed [Xia et al., 2020]. This approach was later extended by Huang et al. to minimize the difference between actual and perturbed values [Huang et al., 2021]. Nevertheless, these studies have a limited scope, mainly focusing on K-means and LDP algorithms without directly supporting other clustering algorithms.

Moreover, the existing literature concentrates on using 2-dimensional synthetic datasets, which may not fully capture the complexity and diversity of real-world data. This inflexibility and lack of practical testing in the current literature highlight the need for a more comprehensive solution with broader applicability.

Our research incorporates realistic attacks and datasets to assess privacy and utility to bridge this gap and move towards a practically deployable LDP mechanism. By doing so, we aim to develop a privacy framework that enables secure training of various clustering algorithms on distributed n-dimensional data. By addressing the shortcomings of previous studies and considering the challenges posed by real-world datasets, our goal is to advance the field and facilitate the adoption of privacy-preserving clustering techniques in practical settings.

## RESEARCH OBJECTIVE

Another take on Differential Privacy (DP) is called Geo-indistinguishability (GI) and is specifically used for geographical data [Andrés et al., 2012]. The 2D-Laplace mechanism adds noise to 2 data points (latitude/longitude) by generating random locations based on the amount of distance. The mechanism is configurable by a privacy budget and holds within a radius of the original location. An extension of this mechanism, 3D-Laplace, added support for including the altitude to support indoor navigation [Min et al., 2022]. The mechanism has also been extended to nD-Laplace to support n-dimensional data; however, its application focuses entirely on text vectors [Fernandes et al., 2019].

Because the proposed method relies on distance-based considerations, it is plausible that the shape and structure of the data will be preserved during the clustering process. Consequently, we expect this algorithm to perform well for data clustering tasks. To our knowledge, GI has never been previously applied for data clustering. Therefore, the primary research objective of this study is to extend the use of GI to optimize specifically the clustering of data.

This thesis aims at developing a new privacy framework for training clustering algorithms on distributed n-dimensional data. It does so by extending and optimizing the three variants of nD-Laplace to support clustering algorithms. Our contribution to the current literature is as follows:

1. We aim to create a new privacy framework named kd-Laplace using the three variants of Laplace to train privacy-preserving clustering algorithms on distributed n-dimensional data.

2. In addition to K-means, we also demonstrate the working of our method on multiple other clustering algorithms.
3. We introduce different optimization techniques to improve the utility of kd-Laplace using kd-tree.
4. We demonstrate the effectiveness and privacy of kd-Laplace using real-world datasets and attacks.

## RESEARCH QUESTION AND METHODS

To reach our goals and objectives, we have formulated the following main question:

*"How can the kD-Laplace algorithm be applied in training privacy-preserving clustering algorithms on distributed n-dimensional data?"*

Much literature is consulted to answer this research question, and we experiment extensively to gather enough qualitative data. Both external and internal validation methods are employed to assess the utility and privacy of the privacy mechanism. The qualitative data from these experiments are analyzed to compare the different privacy mechanisms and their performance concerning clustering tasks.

## ROADMAP

The first chapter is devoted to literature study and review. We explain the various concepts related to clustering and differential privacy. This chapter concludes with a literature review. The second chapter explicitly focuses on the kd-Laplace part, explaining the relevant algorithms and concepts. We then delve deeper into privacy attacks and how they can be evaluated.

The following chapter is dedicated to describing the research itself and the methodology. Subsequently, we present the results, and finally, we conclude with the discussion and conclusions.

# 2

## BACKGROUND

The current chapter explains the theoretical aspects of the research and evaluates similar and previous studies. First, the concept of Differential Privacy and its various types is explained. Next, we examine clustering algorithms and the methods for evaluating their performance. Regarding the related literature, we first look at studies that have used differential privacy and then the studies that have combined it with cluster methods.

## 2.1. DIFFERENTIAL PRIVACY

In practice, data is often sent to a central storage point. This requires trust, and because all data is collected in one place, the risk of private data leakage becomes very high. By applying DP, noise can be added to the data to protect it. This principle is illustrated in the following figure:

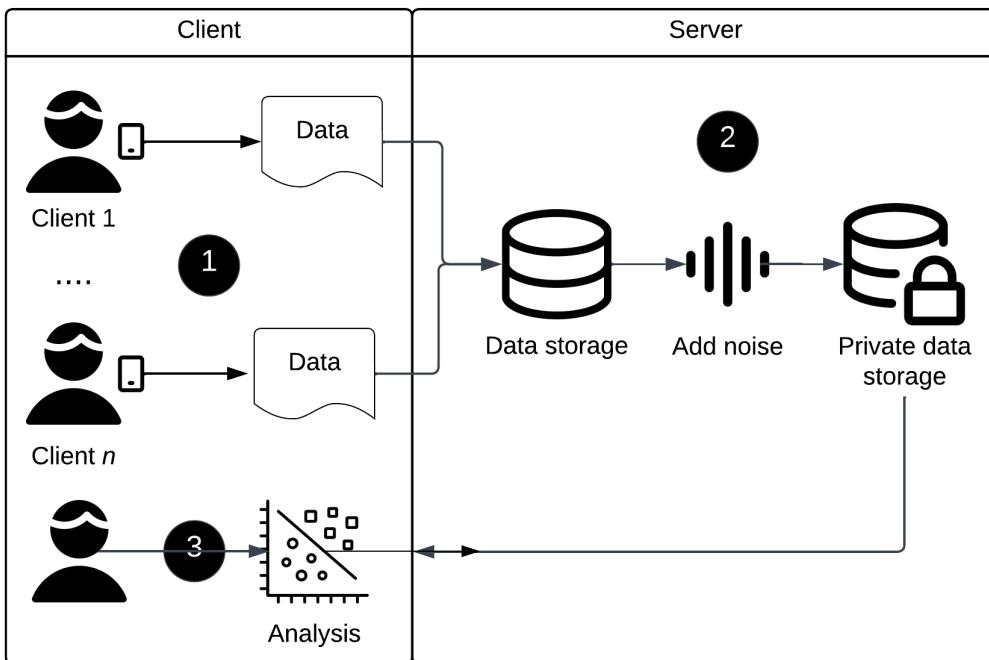


Figure 2.1: General approach for setting up (central) differential privacy.

1. The system that receives data from users. It is assumed in this setting that the system is trustworthy and that the data is securely stored.
2. The data is collected and stored in a database, and the noise is added by a differential privacy mechanism.
3. Someone interested in the data, such as a data scientist, can access and analyze the private data.

To a certain extent, a user's privacy would be ensured with differential privacy.

Although differential privacy solves privacy problems, it remains challenging to calibrate the mechanism. There is an important trade-off between utility and privacy for the adversary. For example, a data scientist wants accurate data. At the same time, the noise must be sufficient to prevent an attacker from obtaining too much information.

For this reason, a few sections will be devoted to outlining the mathematical background of differential privacy. We will examine which factors influence this calibration and whether other methods contribute. Afterward, we will further explain other types of differential privacy (local and geo-indistinguishable) similarly.

### 2.1.1. DEFINITIONS

We examine the different notations and types of differential privacy we consider in this research.

#### $\epsilon$ -DIFFERENTIAL PRIVACY

Dwork et al. formulated the notion of privacy: Participating in a database should not significantly increase the risk of an individual's privacy being compromised [Dwork \[2006\]](#). This is mathematically formulated in the same research named [DP](#). Using the definition of privacy, it is formulated as the maximal possible change when adding or removing a single record [\[Dwork, 2006; Friedman and Schuster, 2010\]](#). This is reflected using the formal mathematical formulation as formulated by Dwork et al:

$$\Pr[K(D_1) \in S] \leq \exp(\epsilon) \times \Pr[K(D_2) \in S] \quad (2.1)$$

So, given a randomization function  $K$ , it gives  $\epsilon$ -differential privacy if dataset  $D_1$  and  $D_2$  are differing at most one element [\[Dwork, 2006\]](#). The  $\epsilon$  determines the amount of noise (privacy budget) [\[Friedman and Schuster, 2010\]](#). The lower the value of  $\epsilon$ , the higher the privacy guarantee. In this regard, it is essential for a method that ensures differential privacy to consider this. For this reason, a common way to calibrate the  $\epsilon$  is to calculate the sensitivity. This value is calculated based on the impact of a function or query on the data. For example, if there is a method called *sum* for the summation of data points, the method's sensitivity is 1. This is because removing one data point would significantly affect the outcome, and  $\epsilon$ -differential privacy could no longer be guaranteed. It is also mathematically defined by Dwork et al.:

$$\Delta f = \max_{D_1, D_2} \|f(D_1) - f(D_2)\|_1 \quad (2.2)$$

#### $(\epsilon, \delta)$ -DIFFERENTIAL PRIVACY

The formal notion of differential privacy has only the privacy budget  $\epsilon$  as a parameter. This formulation is strict, but most methods relax this a little, with  $(\epsilon, \delta)$ :

$$\Pr[K(D_1) \in S] \leq \exp(\epsilon) \times \Pr[K(D_2) \in S] + \delta \quad (2.3)$$

This equation means the sensitivity ( $\delta$ ) loosens the formal definition of differential privacy. The  $\delta$  is added to the upper bound of the exponential, so the desired noise can be higher. To some extent, the delta represents the probability of the algorithm leaking information [\[Aitsam, 2021\]](#). With  $\epsilon$ -differential privacy, there would be no difference in the case of information leakage ( $\delta = 0$ ). However, with  $(\epsilon, \delta)$ -differential privacy, the information can leak up to the probability of delta.

### $\epsilon$ -LOCAL DIFFERENTIAL PRIVACY

As the name suggests, **LDP** is executed on the client side instead of on the server (See Figure 2.1). Local differential privacy removes the "trusted" curator/ server, preventing sensitive data leakage even if an attacker gains access to the dataset [Xiong et al., 2020a]. The principle of **LDP** is illustrated in the following figure:

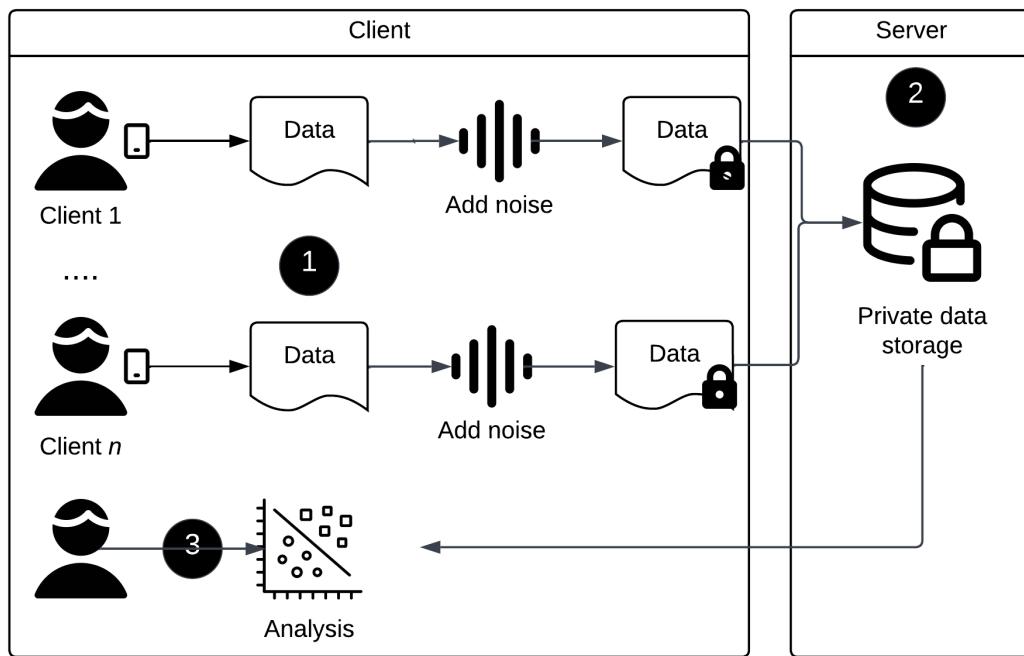


Figure 2.2: Local differential privacy, which moves the noise-adding step to the client side.

1. The clients have their data record(s) locally (client-sided) and add noise to them directly. This way, the server never receives plain data.
2. The data is sent to the server and can just be stored securely.
3. As with **DP**, a stakeholder interested in the data receives it and can analyze it.

Local differential privacy has two different frameworks. It can be set up as interactive or non-interactive [Xiong et al., 2020a].

Where interactive can be distinguished as either sequentially interactive or fully interactive <sup>1</sup>:

<sup>1</sup>Image source: [https://www.majos.net/focs\\_19\\_talk.pdf](https://www.majos.net/focs_19_talk.pdf)

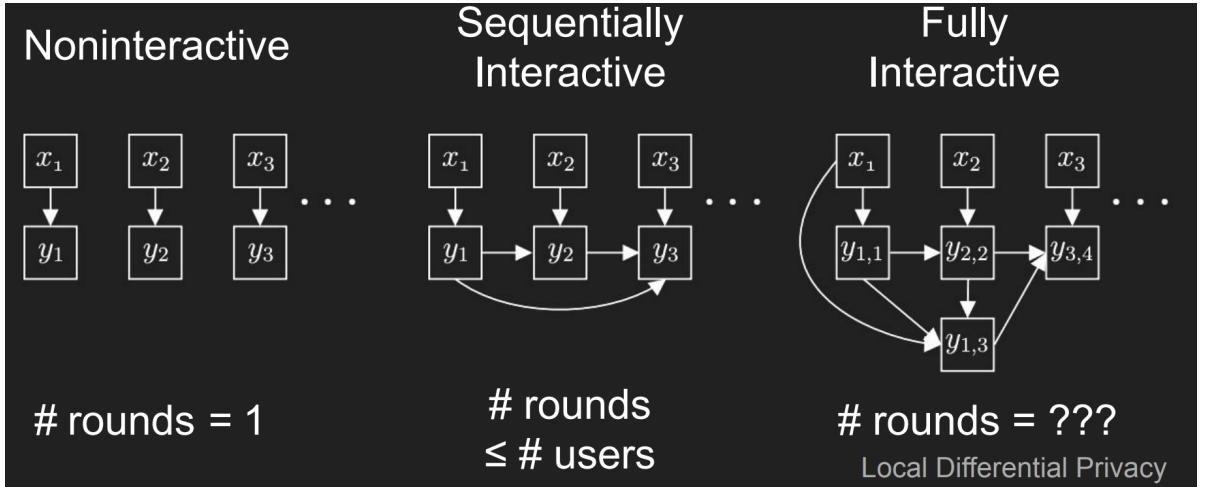


Figure 2.3: Non-interactive versus sequentially interactive versus fully interactive local differential privacy [Joseph et al., 2019]

In the non-interactive framework, the privacy mechanism operates on a single data instance (e.g., users' smartphones) without requiring interaction with other clients or between data points. This interaction means that each data point is perturbed independently, and the privacy guarantees are achieved without communication or coordination between the data points. In the interactive framework, the privacy mechanism involves a process of iterative interactions between users and clients. Each data point considers other data points' responses to determine its own perturbation or privacy level. The interactions can occur through a communication channel, where data points exchange information with each other or a trusted mediator. The iterative nature of the interactive framework allows for refining and improving the privacy guarantees based on the global knowledge of the data set.

In figure 2.3, the different interactive types are visualized, where the most notable difference is the around of rounds [Xiong et al., 2020b]. A sequentially interactive framework passes at most one message between data points. A fully interactive framework could provide an unknown and unlimited amount of messages between data points. The interactive framework has unlimited access to other data points, which is most suitable for practical appliances [Xiong et al., 2020b]. The practical applications are better because the data points' mutual correlation also contributes to utility [Wang et al., 2020].

#### $\epsilon$ -GEO-INDISTINGUISHABILITY

The last and most important type of differential privacy for this study is GI. GI can be applied to preserve privacy using a differential privacy method specific to spatial data [Andrés et al., 2012]. Consider a local-based system (LBS) that provides a map service to its users (e.g., Google Maps). The users can query the LBS for a route from their current location to a destination. Instead of the exact (private) location, the LBS should only receive an approximate location enough to provide the service at an adequate level. This approximation happens entirely locally, so no central server exists.

In addition to this, the mechanism provides an initiative way of calibrating the desired privacy, which is illustrated in the following figure:



Figure 2.4: Visualisation of radius  $r$  with  $l$  to provide Geo-indistinguishability mechanism [Andrés et al., 2012]

The figure shows two variables to configure the amount of noise:  $r$  and  $l$ . Users provide a privacy radius  $r$ , and the noise is added and preserved within this radius. A fake location  $z$  is generated within this radius based on the Euclidean distance  $d$ . The amount of noise can be configured using the privacy level  $l$ . Thus, a higher  $r$  or/and  $l$  provides a higher privacy guarantee.

In a practical scenario, the  $r$  and  $l$  variables are combined to obtain the privacy budget  $\epsilon$ . The privacy budget is obtained by using the formula  $\epsilon = \frac{l}{r}$ , and the  $\epsilon$  is used as input for the GI mechanism to calibrate the noise.

The formal definition of GI is the following [Andrés et al., 2012]:

$$K(x)(y) \leq e^{\epsilon * d(x, x')} K(x')(y) \quad (2.4)$$

$K$  is a probability method reporting  $x, x' \in X$  as  $z \in Z$ . According to the Andres et al., the mechanism provides  $(\epsilon)$ -Geo-indistinguishability for any  $z$  that is generated within the radius  $r$  [Andrés et al., 2012].

The intuition of the GI definition is that it displays the distinguishability level between two secret locations/points  $x$  and  $x'$  [Chatzikokolakis et al., 2015]. An extension of this is called  $d_x$ -privacy, a more general notation of distance-aware differential privacy. Therefore, the definition for GI is  $d_2$ -privacy, but it is essentially the same as the proof provided for GI.

### 2.1.2. LAPLACE MECHANISM

The previous sections have explained the different privacy frameworks that exist for differential privacy. In the current section, we explain an important (L)DP mechanism on which the mechanism of GI is based. We provide the definition, laying the groundwork for the mechanisms of GI.

The Laplace mechanism, initially proposed in the same paper as a foundational concept for differential privacy [Dwork, 2006], introduces noise to the data based on a privacy budget ( $\epsilon$ ). Additionally, it incorporates a sensitivity parameter ( $\delta$ ). While the privacy budget determines the amount of noise to be added, sensitivity is used to calibrate the output based on the input [Xiong et al., 2020a].

The definition of the Laplace mechanism is as follows [Xiong et al., 2020a]:

$$M(f(x), \epsilon) = f(x) + (Z_1, \dots, Z_d) \quad (2.5)$$

Here,  $Z$  is an independent collection of random variables drawn from the Laplace distribution. The distribution scale is determined by the sensitivity of the function  $f$  and the

privacy budget  $\epsilon$  [Xiong et al., 2020a]. Hence, the definition of the Laplace distribution is  $Lap(\delta f/\epsilon)$ , where  $\delta f$  is the sensitivity of the function  $f$ .

The sensitivity of  $f$  can be calculated globally and locally. Global sensitivity is independently calculated over two different datasets and is part of the original definition of DP [Dwork, 2006]. Usually, this is not the desired situation since local sensitivity always has more context of the dataset in question [Nissim et al., 2007]. As a result, the trade-off for noise is much more precise, and the balance between utility and privacy is much better. The definition of sensitivity is the following [Nissim et al., 2007].

$$LS(f, x) = \max_{x' : d(x, x') \leq 1} |f(x) - f(x')| \quad (2.6)$$

Here,  $d$  is the distance between two datasets  $x$  and  $x'$ . The equation calculates the maximum difference between the output of  $f(x)$  and  $f(x')$ , which should be at most  $LS(f)$ . Because the sensitivity is the maximum difference between output and input, it looks a lot like the definition of DP (Equation 2.2).

In a practical scenario, it can be hard to calculate the local sensitivity. To this end, the smooth sensitivity method was introduced by Nissim et al. [Nissim et al., 2007]. This method aims at smoothing out the local sensitivity by focusing on reducing the noise without the risk of revealing more information. Because it calculates many different configurations, a disadvantage of this method is that it can be computationally expensive. Also, the introduction of this method makes Laplace preserve  $(\epsilon, \delta)$ -DP instead of pure  $\epsilon$ -DP.

## 2.2. CLUSTERING

This section explains the various clustering algorithms used in the research. Subsequently, we will explain and indicate how the hyperparameters are handled in the research for each algorithm.

### 2.2.1. TYPES OF CLUSTERING ALGORITHMS

Multiple clustering types can be utilized for clustering data [Xu and Tian, 2015]. For this study, we selected the types that use some form of distance for clustering. We explain each cluster type with corresponding clustering algorithms.

#### PARTITION BASED CLUSTERING

With this type, the data is partitioned and allocated into a fixed amount of clusters. A well-known and popular clustering algorithm is K-means. Based on Loyd et al.'s original method, this algorithm randomly selects  $k$  points as cluster centroids [Lloyd, 1982]. The first step of the algorithm is to assign each datapoint  $x_i$  to the nearest centroid. Then the mean is calculated by calculating the mean of all points in the cluster [Yuan and Yang, 2019]:

$$c_i = \frac{1}{|C_i|} \sum_{x_i \in C_i} x_i \quad (2.7)$$

Finally, the K-Means algorithm aims to minimize the within-cluster sum of squares (WCSS) [Yuan and Yang, 2019]:

$$WCSS = \sum_{i=1}^k \sum_{c_i \in C_i} \|x_i, c_i\|^2 \quad (2.8)$$

The algorithm runs multiple times until the center point no longer changes or reaches the maximum iterations [Yuan and Yang, 2019].

#### DENSITY BASED CLUSTERING

The data points are partitioned for this clustering type based on nearest neighbors [Fahad et al., 2014]. A region with a high data density is considered a cluster [Xu and Tian, 2015]. A popular density-based clustering algorithm is **Density-based spatial clustering of applications with noise (DBSCAN)**. The method was introduced by Ester et al. and worked by drawing a radius around data points [Ester et al.]. For this thesis, the  $\Pi$ -symbol is used not to confuse it with the privacy budget  $\epsilon$ . It then groups all points within this radius as clusters. The main advantage is its ability to find arbitrarily shaped clusters and detect outliers [Liu et al., 2012].

Another density-based algorithm is **Ordering Points to Identify Clustering Structure (OPTICS)**, an extension of **DBSCAN**. The algorithm attempts different  $\Pi$  values to achieve the best result [Ankerst et al.]. The algorithm uses **DBSCAN**, so we explain this algorithm first in detail. It calculates the "density-reachable" property of each data point first [Ankerst et al.]:

$$N_P i(x_i) = x_i \in X : d(x_i, x_j) \leq \Pi \quad (2.9)$$

If  $N_P i(x_i) \geq minPts$  then  $x_i$  is a "core-object" [Ankerst et al.]. These core objects are used for identifying dense regions to grow clusters. Second is the "density-connected" property, where two points can be linked through a series of points close to each other (within distance  $\Pi$ ). And, each has a sufficient number of neighbors ( $minPts$ ) [Ankerst et al.]. The

final step is to build clusters based on these two properties. A cluster  $C$  is build up using the following rules [Ankerst et al.]:

1. If a core-object in  $C$  is density-reachable from  $x_i$ ,  $\Pi$  and  $minPts$ . Then  $x_i$  is added to  $C$ .
2. If a core-object is density-connected to  $x_i$ ,  $\Pi$  and  $minPts$ . Then  $x_i$  is added to  $C$ .

Finally, **OPTICS** comes into play, which makes the radius  $\Pi$  obsolete. Instead of constructing clusters, the **OPTICS** algorithm produces an ordered list. In the next step, they create a reachability plot that shows valleys (clusters) and peaks (cluster transitions). The algorithm then uses an automatic approach to determining the characteristics of the reachability plot to determine clusters.

#### HIERARCHICAL CLUSTERING

A hierarchy of clusters is generated and merged according to the closeness of data points [Meng et al., 2021]. Hierarchical clustering can be divided into agglomerative and divisive [Meng et al., 2021]. This approach means the clustering occurs from the bottom-up or top-down, respectively. Of the two approaches, agglomerative is the most common one [Meng et al., 2021]. For the implementation, there are four popular methods: Single link, Average link, Complete link, and Ward's method [Roux, 2015]. These methods are used to calculate (dis)similarities between objects.

Provide formal definition

#### RESEARCH DIRECTION

There are three different types of clustering algorithms that we have considered for this research. Therefore, we will choose one algorithm for each type to ensure a representative selection.

Based on the different types of clustering algorithms, we have chosen to use K-Means for this research. K-Mean clustering is widely used in the literature and, therefore, easily compared with other studies.

For density-based clustering, we will solely focus on **OPTICS** in this study. The parameter  $\Pi$  in **DBSCAN** is influenced by the data shape, which is further affected by our privacy mechanism. However, as **OPTICS** automatically determines the  $\Pi$ , the impact of this parameter is less significant, enabling us to emphasize the evaluation of the privacy mechanism itself.

As for hierarchical clustering, we have decided to go for agglomerative clustering.

#### 2.2.2. PARAMETER SELECTION

A list of important hyperparameters that can influence the results is provided for each clustering algorithm. Subsequently, we briefly discuss the methods for determining these hyperparameters for the algorithms.

We use the same distance function for all algorithms to have all the algorithms in the same setting. This function will be Euclidean distance, as used for **GI**.

#### K-MEANS

The most crucial parameter of the K-Means algorithm is the value of  $k$ . This value determines the number of clusters to consider and influences the results by a lot [Ahmed et al.,

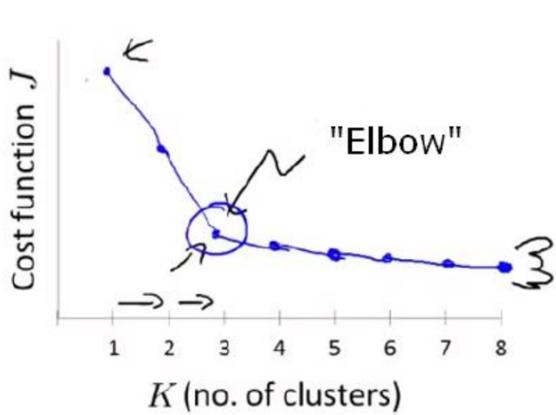


Figure 2.5: Illustration of determining  $k$  using the "elbow" method [Kodinariya and Makwana, 2013]

2020]. The first method is called an "elbow" plot [Kodinariya and Makwana, 2013]. This method finds the best  $k$  by applying the K-Means algorithm multiple times and estimating the best  $k$ . In this situation, a Silhouette plot can be used to determine the  $k$ . This method uses the Silhouette coefficient for each cluster [SAPUTRA et al., 2020]. Because this metric measures the separation and cohesion of the cluster, it can be used to choose an optimal amount of clusters [SAPUTRA et al., 2020]. By visualizing the silhouette scores, the  $K$  with the highest coefficient has the best-separated clusters.

Another method is the Gap statistic method [Yuan and Yang, 2019]. It compares the total within-cluster variation for different values of  $k$  with their expected values under a null reference distribution of the data [Tibshirani et al., 2001]. A practical appliance of this method uses a line plot for comparing the  $k$ -value and gap value [Yuan and Yang, 2019]. Based on the line's visual change, someone can select the best  $k$  (see Figure 2.5).

In summary, no fixed method exists to choose a good  $k$  for K-Means. The elbow method is standard in the existing literature and is very popular due to its simplicity. However, one disadvantage is that it can be challenging to determine the "elbow" point, as it is not always present [Kodinariya and Makwana, 2013]. In that case, the silhouette or gap statistic method is a good alternative, where the silhouette score is the simpler method.

### AGGLOMERATIVE CLUSTERING

The agglomerative clustering algorithm has two essential hyperparameters: the linkage method and the number of clusters.

**Choosing linkage:** This decision is crucial for the algorithm as it determines the similarity measurement. Choosing a linkage method that uses Euclidean distance is important for our research. From the four different types of linkage methods, we have chosen the Ward method because this method uses Euclidean distance [Roux, 2015; Seetharaman et al., 2019].

**Choosing number of clusters:** There are not many methods for determining the initial number of clusters. Therefore, the approach is similar to K-Means. The amount of clusters is selected based on the "elbow" with the highest silhouette score.

### OPTICS

With the introduction of OPTICS, we no longer need to worry about the  $radius(\epsilon)$  value. Therefore, only  $minPts$  remains an important parameter to consider.

The *minPts* is the minimum amount of points that have to be within the *radius( $\epsilon$ )* to mark it a cluster. This hyperparameter is analyzed in a paper by Sander et al. The work describes calculating this parameter by applying two times the feature amount [Sander et al., 1998]. So, using this approach, a dataset with two features will have an *minPts* of four [Schubert et al., 2017].

**DBSCAN** is a little more complicated due to the variety of datasets and noise-altering mechanisms we experiment with. This complexity is why we use **OPTICS** to determine the best value for *radius( $\epsilon$ )* and choose *minPts* based on the number of features times two.

### 2.2.3. EVALUATION METHODS

Clustering comparison measures are important in cluster analysis for external validation by comparing clustering solutions to a "ground truth" clustering [Vinh et al.]. These external validity indices are a common way to assess the quality of unsupervised machine learning methods like clustering [Warrens and van der Hoef, 2022]. A method that could be used for this is the Rand Index [Rand, 1971]. It is a commonly applied method for comparing two clustering algorithms [Wagner and Wagner]. An improvement of this method is adjusted for chance by considering the similarity of pairwise cluster comparisons [Vinh et al.]. Both the Rand Index (RI) and Adjusted Rand Index (ARI) [Hubert and Arabie, 1985] report a value between 0 and 1. Where 0 is for no-similarity and 1 for identical clusters. Alternatives for RI are the Fowlkes-Mallows Index and Mirkin Metric. However, these two methods have their disadvantages. They are, respectively, sensitive to a few clusters and cluster sizes [Wagner and Wagner]. The ARI metric suffers from cluster size imbalance as well, so it only provides not a lot of information on smaller clusters [Warrens and van der Hoef, 2022]. Instead, they recommend using the cluster index metric proposed by Fränti et al. [Fränti et al., 2014].

Another popular group of methods is the information theoretic-based measures [Vinh et al.]. This metric measures the information between centroids; the higher the value, the better [Vinh et al.]. **Mutual Information (MI)** is a metric that calculates the probability of an element belonging to cluster  $C$  or  $C'$ . But, it is not easy to interpret as it does not have a maximum value [Wagner and Wagner]. To this end, **Normalized Mutual Information (NMI)** can be used to report a value between 0 and 1 using the geometric mean [Strehl and Ghosh, 2002]. The metric also exists in an adjusted version as **Adjusted Mutual Information (AMI)**. This metric works in the same way as for the **Adjusted Rank Index (ARI)** and is mainly needed if the number of data items is small in comparison to the number of clusters [Vinh et al.].

Besides the external validity measurements for clustering, it is also possible to use internal validation methods. These metrics focus entirely on the intrinsic dataset properties instead of relying on an external baseline clustering algorithm [Craenendonck and Blockeel]. They are assessing two essential concepts of clustering: compactness and separation [Hassani and Seidl, 2017]. Both studies consider three different metrics and measure both concepts at the same time [Hassani and Seidl, 2017]:

1. **Calinski-Harabasz Index (CHI)** [Caliński and Harabasz, 1974] is used to measure the cluster variance (well-separated clusters) and low variance within the clusters (tightly coupled data). A high score indicates better clustering.
2. **Silhouette Coefficient (SC)** [Rousseeuw, 1987] This metric is similar in measuring co-

hesion within and separating clusters. However, this metric uses the pairwise distance [Hassani and Seidl, 2017]. A score of -1 indicates incorrect clustering and +1 for correct clustering [Rousseeuw, 1987].

3. Davies-Bouldin (DB) [Davies and Bouldin, 1979] uses the average distance between clusters. A lower score indicates good clustering.

## RESEARCH DIRECTION

This chapter outlines the evaluation methods used and provides a definition and brief explanation for each method.

In the literature, both internal and external evaluation methods are commonly used. While some studies choose one over the other, this thesis will perform both validations. This approach is crucial to measure how well the data shape is preserved (internal validation) and how the algorithm performs in real-world scenarios (external validation). Our research will focus on **AMI** for external validation. Because it holds the same value as **ARI**, but it is better explainable. The second metric we use is for external validation. We use **SC** because this metric returns a bounded value between -1 and 1.

Based on existing studies and the literature, we have selected metrics adjusted to compensate for the data distribution characteristics. Hence, we will focus on these metrics' "adjusted" variants:

### 1. Silhouette coefficient (Internal validation):

This metric is defined as follows [Liu et al., 2010; Rousseeuw, 1987]:

$$s(i) = \frac{b(i) - a(i)}{\max(b(i) - a(i))} \quad (2.10)$$

- (a)  $s(i)$  is the silhouette coefficient for a single datapoint  $i$ .
- (b)  $a(i)$  is the mean distance of  $i$  and all the other points in the same cluster.
- (c)  $b(i)$  is the mean distance of  $i$  and all the other points in the next nearest cluster.

Finally, the final silhouette score is the mean of all datapoint coefficients  $s(i)$ .

2. **Adjusted Mutual Information (External validation):** The first step is calculating a function  $H$ , and the second is using  $H$  to calculate **MI**. Finally, the **AMI** is calculated by adjusting **MI** for chance.

The  $H$  function is defined as follows [Vinh et al.]:

$$H(U) = - \sum_{i=1}^U p(i) \cdot \log(p(i)) \quad (2.11)$$

Here,  $p(i)$  is the probability of a random point falling into cluster  $U$ . The  $H$  function measures the information by calculating the probability of a random point falling into cluster  $U$ . If the outcome is certain (probability 1), it carries no information. But, if the outcome has a lower probability, it carries more information. The **MI** uses this property and is defined as follows [Vinh et al.]:

$$MI(U, V) = - \sum_{i=1}^U \sum_{j=1}^V p(i, j) \cdot \log \frac{p(i, j)}{p(i)p(j)} \quad (2.12)$$

It re-uses the properties of  $H$  to measure the (dis)similarity between two clusterings  $U$  and  $V$ . The  $p(i, j)$  is the joint probability of a random point falling into both cluster  $U$  and  $V$ . This same calculation is executed for the second label  $V$  and divided by each other. The **AMI** is formulated as follows [Vinh et al.]:

$$AMI(U, V) = \frac{MI - E(MI)}{\max(H(U), H(V)) - E(MI)} \quad (2.13)$$

The  $E(MI)$  is the expected value of the **MI** and serves as an upper bound for the **MI**. For the details of this calculation, we refer to the original paper [Vinh et al.]. The outcome of this formula is a value between 0 and 1, where 0 is no similarity, and 1 is identical clusters.

## 2.3. LITERATURE REVIEW

In the search for related literature, we focused mainly on (L)DP mechanisms that can be used for general purposes (e.g. not only mean estimation), which is the most comparable to our mechanism. The related literature is divided into three parts:

1. Differential privacy and clustering: We explain how differential privacy is used for clustering.
2. Differential privacy methods: We explore related literature for DP.
3. Cluster methods with (L)DP: We explore related literature for clustering with (L)DP.

Afterward, we summarize the last two literature studies in a table. This table includes components such as the type (LDP or DP) and whether a public code is available.

### 2.3.1. DIFFERENTIAL PRIVACY AND CLUSTERING

There are different ways how to combine differential privacy with clustering [Baraheem and Yao, 2022]:

1. Input perturbation: This method adds noise to the input data, so the clustering algorithms are trained on perturbed data.
2. Gradient perturbation: The term gradient is used in machine learning algorithms that solve an optimization problem like deep learning [Hassan et al., 2019]. Perturbing the gradients ensures the optimization algorithm is private [Ji et al., 2014].
3. Output/label perturbation: This method adds noise to the output of the clustering algorithm. The model trains using the original data, and the noise is added to the output [Ji et al., 2014].
4. Objective perturbation: This method slightly differs from output/label perturbation but also adds noise to the output. The difference is it is added to the output of a loss function instead of the labels (for example, logistic regression) [Baraheem and Yao, 2022].

Many approaches targeting local differential privacy mechanisms must adapt the clustering algorithm (for example, k-Cluster) [Sun et al., 2019], which we explain later in this chapter. For this reason, we concentrate on input perturbation in this study, as it allows a dataset to be used to train various clustering algorithms without modification.

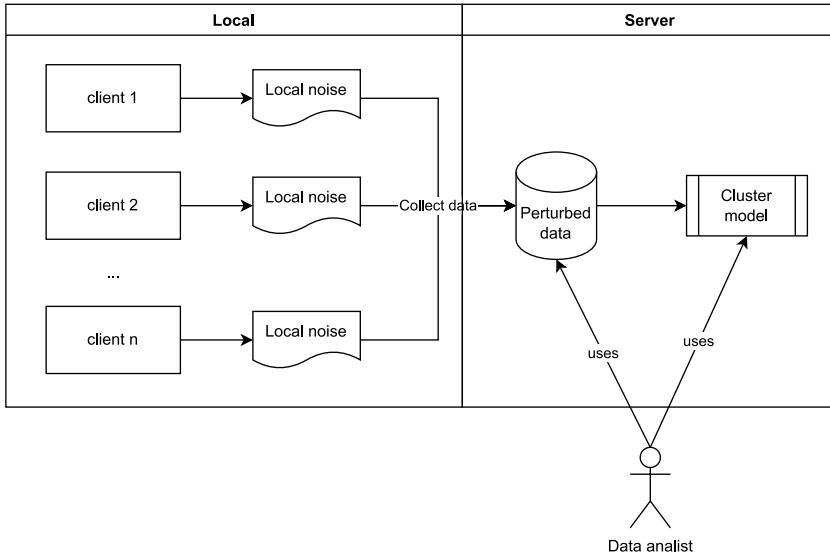


Figure 2.6: Example on how our approach of private training works for clustering and DP

Improve image

1. Left, each local dataset is perturbed using our method (locally).
2. All data is collected and used to train a clustering algorithm (on the server).
3. A cluster model is privately trained because we train the algorithm on private data.
4. Because we use input perturbation, a third party like a "data analyst" can also directly access the data privately.

The next section focuses on the differential privacy algorithms that are of interest to this type of input perturbation.

### 2.3.2. DIFFERENTIAL PRIVACY METHODS

As discussed earlier, the Laplace method was the first to establish DP [Dwork, 2006]. The first paper we discuss is provided by Soria-Comas et al. and considers the distribution of the dataset for generating noise [Soria-Comas and Domingo-Ferrer, 2013]. Their work claims the Laplace mechanism is not optimal for a univariate function and aims at improving it by extending the Laplace mechanism. The proposed method performs slightly better than Laplace on multivariate/multiple queries.

Quan et al. also proposed a new method to extend Dwork et al.'s Laplace algorithm [Geng et al., 2015]. They introduced the Staircase mechanism for 1-dimensional noise, which was later extended to support multidimensional data [Geng et al., 2015]. The mechanism aims to improve utility by adding the same level of privacy while adding less noise. It represents a staircase-shaped **Probability Density Function (PDF)**, hence the name Staircase Mechanism (SM). This mechanism accepts three configurable parameters comparable to those of the Laplace mechanism. The authors' work can handle multidimensional numerical data and preserve the  $(\epsilon)$ -differential privacy. In the previous two paragraphs, we have mainly focused on the interesting literature regarding differential privacy. Next, the succeeding paragraphs will mainly center on related literature concerning local differential privacy.

Nguyen et al. introduced a new **LDP** mechanism for working with numerical data [Nguyễn et al., 2016]. Their primary focus is estimating means and frequencies on data and applying machine learning techniques such as Support Vector Machines (SVM) and linear regression using Empirical Risk Minimization. Initially, the authors analyze Duchi et al.'s method [Duchi et al., 2013] and highlight several shortcomings. To address these issues, they introduce Harmony as a mechanism for **LDP** perturbation. This mechanism can perturb categorical and numerical data, providing high accuracy for classification and regression tasks. To compare Harmony to other methods, they create the Hybrid Mechanism (HM), a combination of two existing methods for categorical and numerical data. For this purpose, they extend other work [Bassily and Smith, 2015] for perturbing multidimensional categorical attributes and use Duchi et al.'s method for numerical data [Duchi et al., 2013]. This combination of mechanisms allows them to compare their Harmony mechanism to the hybrid mechanism and measure the utility/accuracy differences.

Duchi et al. improved their method by formalizing the trade-off between statistical utility and (local) privacy, analyzing multiple estimation problems [Duchi et al., 2017]. Examples include mean, median, and density estimation. To achieve this, they use minimax, a technique for finding the worst-case probability distribution. Additionally, they focus on existing work and propose several optimization strategies.

Duchi et al.'s method was extended a couple of years later by adding support for bounded and multidimensional data [Wang et al., 2019]. The authors introduce the **Piecewise Mechanism (PM)** to handle numeric and categorical data and the Hybrid Mechanism (HM), which combines **PM** and Duchi et al.'s method for 1-dimensional data. Their method's effectiveness is demonstrated using **Support Vector Machine (SVM)**, linear regression, and mean estimation. **PM** and HM are compared to Laplace and Duchi et al.'s solutions. Optimized Unary Encoding (OUE) [Wang et al.] is also used for comparison, but for categorical data only, as the other methods do not support this.

### 2.3.3. CLUSTER METHODS WITH (L)DP

This chapter examines the various studies conducted on clustering in combination with differential privacy. Initially, we looked at the most fundamental papers in this field. Subsequently, the focus shifted toward researching well-known papers published since 2020.

The first work we highlight was proposed by Nissim et al. and aimed at improving differential privacy methods, such as Laplace, which uses sensitivity to compensate the noise for a function [Nissim et al., 2007]. In addition to compensating the function, they also consider the dataset itself. The algorithm is called "smooth sensitivity" and is used for instance-specific noise. To apply it, the authors introduce a method/framework to calculate it effectively. They use K-means, among other clustering algorithms, to demonstrate the method's effectiveness. Their method requires the calculation of cluster distances using Wasserstein distance instead of Euclidean distance.

Another study focuses on interactive and non-interactive approaches for differential privacy in K-Means [Su et al., 2015]. The study builds upon the work done for DPLlloyd, an interactive privacy extension of K-Means described by Blum et al. [Blum et al., 2005]. The DPLlloyd mechanism partitions an n-dimensional dataset into a grid and releases the count for each grid by adding Laplacian noise to each count. Another part of their research focuses on determining the width of the data cells. The grid estimation method used in their research is called the extended uniform grid approach (EUG), and the complete K-Means method is called EUGkM. The experiment evaluates it against the DPLlloyd mechanism, which performs better in an interactive setting. Therefore, they combine their algorithm for combining both aspects into a hybrid approach (EUGkM + DPLlloyd approach) and show a better final performance.

The study of Nissim et al. researches finding the smallest possible radius in the Euclidean space  $R^d$  for a set of  $n$  points [Nissim and Stemmer, 2018]. They propose a new solution that uses locality-sensitive hashing (LSH) for differential privacy to find 1-cluster in the  $d$ -dimensional Euclidean space. This method works for differential privacy (LSH-GoodCenter), but they also extend this to the local model (LDP-GoodCenter). The algorithm to find this radius is used to count the points enclosed by the radius, and Laplace noise is added to the count to preserve differential privacy. The mechanism is combined and applied to work with the K-Means algorithm (LDP-K-Mean). This mechanism was extended in a paper proposed by Kaplan et al. and introduced a similar LDP method [Kaplan and Stemmer, 2018]. They aim to reduce the number of interactions needed between the server and users to one, instead of the  $O(k \log n)$  required for Nissam et al.'s solution [Nissim and Stemmer, 2018]. To increase the success probability, they use the same idea but extend it to have multiple centers instead of a single large one. They call it the LSH-Procedure, and the algorithm Private-Centers is applied to generate centers to use with K-Means. Then, they apply the same method to the LDP method initially proposed by Nissam et al. (LDP-GoodCenter). The most recent work by Stemmer et al. focuses on improving the work done by Kaplan et al. [Kaplan and Stemmer, 2018; Stemmer, 2021]. Because the original mechanism has a higher additive error, the noise added introduces a lot of error. To solve this, the authors aim to reduce this error by improving the original GoodCenter algorithm [Nissim and Stemmer, 2018]. Their extended method is called WeightedCenters and also adds weights to candidate centers. In the final iteration, the weights are used to create the K-Means or K-Median clusters.

Sun et al. proposed a mechanism for distributed clustering using local differential pri-

vacy (LDP) to preserve distance-based information. They claim to have the first non-interactive LDP algorithm for clustering [Su et al., 2015]. This means they can perturb the data locally at once and send it to the server to cluster with both K-Means and DBSCAN without interaction between data points. They encode the client-side data into an anonymous hamming space using Bit Vector (BV) and modify the encoding to preserve Euclidean distance. Their mechanism only shares distance information, so they could not use K-Means directly. To overcome this, they modified the algorithm and called it K-Cluster. Finally, the method is evaluated using Normalized Mutual Information (NMI) and Average Estimated Error (AEE).

Xia et al. noticed the shortcoming of Sun et al.'s work which is the need to share privacy-sensitive distance information [Xia et al., 2020]. Therefore, they create an interactive method for distributed K-means clustering using LDP. The method converts features to binary strings and uses the Random Response mechanism (RR) to perturb each feature into a feature vector. The privacy cost depends on the length of the bits of each feature transformation, meaning that a longer length yields more information at the expense of the privacy budget. In each iteration, the serverside calculates and sends K-means centroids to each user, who recalculates distances until the centroids become stable. The approach has the disadvantage of a high correlation between user data and the clusters. To solve this problem, the algorithm is improved by having the client side send the user data and a set of random zero strings. The server side then performs similar calculations to determine the actual cluster. Huang et al. propose a private distributed K-means clustering algorithm for interval data that addresses a shortcoming in Xia et al.'s work by using Condensed Local Differential Privacy (CLDP) for small-scale values and LDP for large-scale values [Huang et al., 2021]. They preserve distance using a Square Wave (SW) mechanism and apply a classical K-Means algorithm on the server side to the perturbed data.

A recent mechanism that also builds around K-Means to preserve LDP is called the LDPK mechanism [Yuan et al., 2021]. As K-Means works only with numerical data, they use K-prototypes for supporting mixed data types. The LDPK mechanism perturbs the user data first locally and interactively exchanges information with the server to complete the clustering process. The mechanism they use for perturbation is the Harmony algorithm, proposed earlier by [Nguyễn et al., 2016]. The S-Hist method is used to support categorical data, which was also introduced by Nguyen et al. The server could still infer the correct information because of the correlation between the cluster centroids and the actual data. Therefore, the author replaces S-Hist with OUE [Wang et al.] to improve accuracy. To improve their mechanism further, the authors disturb the user's cluster information with an extra extension to the LDPK method, called ELDPK. For this purpose, they perturb the clusters with the GRR (Generalized Random Response) algorithm. They show that the clustering quality increases if the data points increase.

Most existing work focuses on (L)DP in combination with K-Means. Finally, two interesting studies focus on differential privacy for Affinity Propagation (AP) or DBSCAN. A study conducted by Cai et al. focuses on AP [Cai et al., 2020]. Their method involves adding Laplace noise to the responsibility matrix. For each sample data, a neighborhood is specified using a radius around the data point. This area is called the neighborhood density, and each sample point's preference value is adjusted according to its density value. Higher density yields a higher chance of belonging to a cluster center and being ranked based on size. The perturbed responsibility matrix and densities are combined and used to run AP.

They evaluated their method using the ARI, Fowlkes-Mallows Index (FMI), and AMI.

Another study focuses on differential privacy for DBSCAN [Bozdemir et al.]. The proposed solution involves clustering data between two or more parties using two servers. Secure two-party computation (S2PC) is used to achieve this. Using S2PC, both servers receive a random-looking secret share. To recover the original data, both servers must combine their shares using S2PC, which combines the data without the servers having access to the full value. The proposed protocol is privacy-preserving DBSCAN (ppDBSCAN). The calculations in this study are based on squared Euclidean distance (SED) and are evaluated using different methods. To evaluate the performance of ppDBSCAN, the study compares its Adjusted Rand Index (ARI) to that of K-means.

Paper	Data type	Dataset	Code	Preserving	Type	Interactive	Methods
[Sun et al., 2022]	-	Synthetic dataset	-	$(\epsilon, \delta)$ -LDP	K-Means	Non interactive	-
[Huang et al., 2021]	-	-	-	LDP	K-Means	Interactive	-
[Stemmer, 2021]	numerical	-	-	LDP	K-Means	Interactive	-
[Yuan et al., 2021]	numerical	Adult dataset, US Census dataset Deer dataset,	-	LDP	K-Prototypes	Interactive	LDPK and ELDPK
[Bozdemir et al.]	-	Lsun dataset, S1	-	DP	DBSCAN	-	ppDBSCAN
[Cai et al., 2020]	-	Iris dataset, Seeds dataset	-	DP	AP	-	DP-AP
[Xia et al., 2020]	n-dimensional numerical data	3D Road Network, CarGPS	-	LDP	K-Means	Interactive	LDPKmeans
[Sun et al., 2019]	n-dimensional numerical data	Aggregation dataset , Digit dataset, Pathbased d...	-	LDP	DBSCAN, K-Means	Non interactive	Distance Aware Bit Vector (DPBV)
[Nissim and Stemmer, 2018]	n-dimensional numerical data	-	-	$(\epsilon, \delta)$ -LDP	K-Means	Interactive	LDP-GOODCenter
[Nissim and Stemmer, 2018]	-	-	-	-	K-Means	Interactive	LSH-Procedure & Private-Centers
[Su et al., 2015]	n-dimensional numerical data	Adult dataset, Gowalla dataset, Image dataset	<sup>a</sup>	DP	K-Means	Both	EUGkM and hybrid EUGkM + DPLlloyd
[Nissim et al., 2007]	n-dimensional numerical data	-	-	$(\epsilon, \delta)$ -LDP	K-Means	Non interactive	Smooth sensitivity

<sup>a</sup> <https://github.com/DongSuIBM/PrivKmeans>.

Table 2.1: Summary table of the literature review for (L)DP clustering algorithms.

Paper	Data type	Dataset	Code	Preserving	Interactive	Methods
[Wang et al., 2019]	n-dimensional and 1-dimensional	BR, MR	<sup>a</sup>	LDP	-	- Piecewise Mechanism (PM) - Hybrid Mechanism (HM)
[Duchi et al., 2017]	1-dimensional	-	<sup>a</sup>	LDP	-	-
[Nguyễn et al., 2016]	numerical, binary and categorical data.	BR	-	$\epsilon$ -LDP		Harmony
[Geng et al., 2015]	n-dimensional	-	<sup>b</sup>	DP	-	Staircase mechanism (SM)
[Geng and Viswanath, 2013]	n-dimensional	-	-	$\epsilon$ -DP	Non interactive	-

<sup>a</sup> <https://github.com/forestneo/sunPytools/>.

<sup>b</sup> <https://github.com/IBM/differential-privacy-lib>.

Table 2.2: Summary table of the literature review for (L)DP algorithms.

number	name	samples	features	target	Source	Realworld data?
1	Adult	48,842	14 numerical /categorical /boolean	income (>50k, <= 50k)	UCI	Yes
2	Seeds	210	7 numerical	type	UCI	Yes
3	Iris	150	4 numerical	class (type of iris)	UCI	Yes
4	CarGPS	17,785,500	3 geographical data	-	-	Yes
5	3D Road Network	434,874	3 geographical data	-	-	Yes
6	Pathbased	300	2 numerical	ground truth clusters	PapersWithCode	No
7	Aggregation	788	2 numerical	-	-	Unknown
8	Digit	1797	8x8 numerical	number	Scikit-learn	Yes
9	Lsun	400	2 numerical	ground truth clusters	-	No
10	S1	1500	2 numerical	ground truth clusters	-	No
11	Deer	20,033	2 numerical	-	-	Yes
13	Gowalla	6,442,890	5 (geographical, data, ids and time)	-	2	Yes

Table 2.3: The different datasets used in the related literature.

### 2.3.4. EVALUATION

It is important to compare our kD-Laplace method with other studies for this research. Studies that are similar to ours utilize LDP for clustering and add noise to cluster centroids [Huang et al., 2021; Xia et al., 2020; Yuan et al., 2021] or require modifications of the K-Means algorithm [Sun et al., 2019].

As explained in Section 2.3.1, we will use input-perturbation. So, our methodology remains independent of the specific clustering algorithm chosen. Based on the literature study, two privacy mechanisms are comparable for this approach: Harmony [Nguyêñ et al., 2016] and Piecewise [Wang et al., 2019].

Both mechanisms focus on solving the problem in Duchi et al.'s paper. However, the Harmony mechanism focuses on solving a paper from 2013 [Duchi et al., 2013] and the Piecewise mechanism on Duchi et al.'s paper from 2017 [Duchi et al., 2017]. The Piecewise mechanism is newer, and given the similarities, Wang et al.'s paper appears much like a more recent version of Nguyen et al.'s paper. Therefore, we only selected Piecewise to compare with our mechanism.

#### DUCHI ET AL.'S MECHANISM

The Piecewise mechanism is based on Duchi et al.'s mechanism for one-dimensional data. Therefore, we start by explaining this mechanism first. Duchi et al.'s mechanism is a relatively simple method based on the Bernoulli distribution [Duchi et al., 2017]. This distribution yields either a 0 (negative) or a 1 (false) based on specified probabilities <sup>3</sup>. The mechanism works on a domain tuple [-1, 1] and returns -1 or 1. So, the PDF of this function is as follows <sup>4</sup>:

$$[H]P(n) = \begin{cases} 1-p & \text{if } n = 0 \\ p & \text{if } n = 1 \end{cases} \quad (2.14)$$

The probabilities are calculated based on the input value and can then be used to estimate a mean value. Compared to Laplace, Duchi et al. solution performs better in variance for epsilons smaller than 2 [Wang et al., 2019]. However, it performs worse when the epsilon value is higher because the algorithm does not consider the privacy budget for values that are 0. The multidimensional variant looks like the one-dimensional version but samples each data point's noise independently.

#### PIECEWISE MECHANISM

The authors aim to create a method that combines the advantages of the Laplace mechanism and Duchi et al.'s methods [Wang et al., 2019]. The goal is to reduce the variance for a broader range of privacy budgets. Like the mechanism above, the Piecewise mechanism only accepts data from -1 to 1. We first explain the one-dimensional variant and then the multidimensional variant.

---

<sup>3</sup><https://mathworld.wolfram.com/BernoulliDistribution.html>

<sup>4</sup><https://mathworld.wolfram.com/BernoulliDistribution.html>

The mechanism takes an input of the range  $[-1, 1]$  and returns a perturbed value in the same range. It uses a piecewise function (hence the name), which is a function that is defined by multiple sub-functions. The Piecewise mechanism starts with randomly selecting a value  $x$  between 0 and 1 and a value  $t \in R^1$ . Then, based on this value the mechanism re-assigns a random value to  $x$ , but with a modified domain (based on the initial value of  $x$ ) [Wang et al., 2019]:

$$t = \begin{cases} [l(t), r(t)] & \text{if } x < \frac{e^{\epsilon/2}}{e^{\epsilon/2} + 1} \\ [-C, l(t) \cup r(t), C] & \end{cases} \quad (2.15)$$

1. Calculated using the privacy budget  $\epsilon$ , the  $C$  determines the amount of noise by bounding the lower/ upper limit of the domain.
2.  $l(t)$  and  $r(t)$  are the left and right boundaries of the domain. These are calculated based on the value of  $t$  and the privacy budget.

Depending on the value of  $t$ , the domain is shifted to the left or right on a domain of three "pieces":  $[-1, 0, 1]$  [Wang et al., 2019]. According to Wang et al., when  $t = 0$ , the probability of  $t \in [l(t), r(t)]$  is higher. In this case, the mechanism generates a perturbed value close to 0 (depending on the privacy budget). This phenomenon is the same for  $t = 1$  and  $t = -1$ .

For multidimensional data, the authors use the one-dimensional version multiple times. However, they compensate for the number of dimensions by generating a  $k$  based on the following formula [Wang et al., 2019]:

$$k = \max(1, \min(d, \lfloor \frac{\epsilon}{2.5} \rfloor)) \quad (2.16)$$

Here,  $d$  is the number of dimensions, and  $k$  is used to sample a  $d$  number of random samples. Finally, the Piecewise mechanism (see Equation 2.15) is executed for each sample.

Might need to recheck

# 3

## ND-LAPLACE

In this chapter, we delve deeper into geo-indistinguishability and the various mechanisms that work with it. We also explain the theory behind the nD-Laplace mechanism and how we modify it and apply it for clustering. Because we frequently discuss 2, 3, and dimensional data in our thesis, we decided to prefix the respective Laplace method with this information. Therefore, for the rest of the thesis, we use the following names:

1. 2D-Laplace: This refers to the planar Laplace mechanism [Andrés et al., 2012].
2. 3D-Laplace: This refers to the spherical Laplace mechanism [Min et al., 2022].
3. nD-Laplace

For each mechanism, we explain the equation for GI, the mechanism, and the data truncation.

### 3.1. 2D-LAPLACE

The idea of GI was introduced to address the issue of privacy and location data [Andrés et al., 2012] (See Equation 2.4). It offers an alternative approach for achieving (local) differential privacy for geographical data (latitude/longitude). The mechanism achieves this by locally adding noise to the location before sending it to a location-based system (LBS). This section starts with an introduction to mathematics, and for each of the different subsections, we visualize and explain open challenges and theoretic for applying them for clustering.

#### 3.1.1. PLANAR AND POLAR LAPLACE

In section 2.1.1, an explanation of the concept of GI has been given. As indicated, the method works on 2-dimensional data, and when visualized, this can be done with a so-called plane [Andrés et al., 2012]. From there, the term "Planar" Laplace (2D-Laplace) originated and was used thoughtfully by Andres et al.

The idea of 2D-Laplace is to generate an area around  $x_0 \in X$  according to the multivariate Laplace distribution. The mechanism of 2D-Laplace is a modification of the Laplace algorithm to support distance [Andrés et al., 2012]. This distance method  $dist(x, x')$  is a method to calculate the Euclidean distance between two points  $x$  and  $x'$ . Recalling the

definition of Laplace, this method  $|x - x'|$  is replaced by the distance metric. Hence, the definition of the PDF by Andrés et al. is:

$$\frac{\epsilon^2}{2 * \pi} e(-\epsilon d(x_0, x)) \quad (3.1)$$

Which is the likelihood a generated point  $z \in Z$  is close to  $x_0$ . The method works for Cartesian coordinates but was modified to support polar coordinates by including  $\theta$ . So each point is reflected as  $(r, \theta)$ . A point  $z \in Z$  where  $z = (r, \theta)$  is randomly generated using two separate methods for calculating  $r$  and  $\theta$ . This idea is visualized in the following figure:

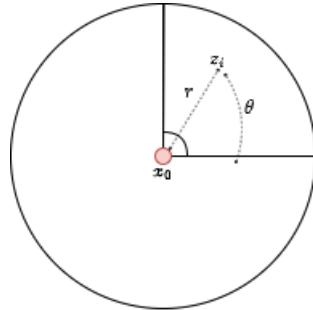


Figure 3.1: Representation of the generated  $z = r\theta$  and original point  $x_0$ .

**Calculating  $r$ :** This variable is defined as  $dist(x_0, z)$  and can be randomly drawn by inverting the Cumulative Distribution Function (CDF) for the Laplace distribution:

$$C_\epsilon^{-1}(p) = -\frac{1}{\epsilon}(W_{-1}(\frac{p-1}{e}) + 1) \quad (3.2)$$

For this equation,  $W_{-1}$  is a Lambert W function with a -1 branch.

Explain -1 branch?

The Lambert w function (also called the product logarithm) is defined as  $W(x)e^{W(x)} = x$  [Lehtonen, 2016]. The purpose of the Lambert w function is to invert the CDF of the Laplace distribution to generate random noise for one of the coordinates ( $r$ ) using the random value of  $p$ .

**Calculating  $\theta$ :** The other variable ( $\theta$ ) is defined as a random number  $[0, 2\pi]$ . To visualize the data, it is necessary to convert the polar coordinates for  $z = (r, \theta)$  back to planar coordinates  $z = (x, y)$ . This conversion is described as step 4 of the planar Laplace algorithm [Andrés et al., 2012] and visualized using figure 3.2.

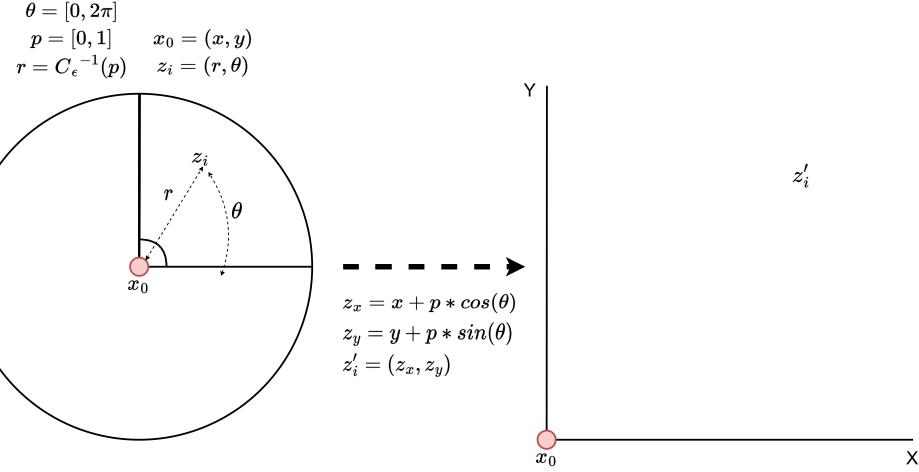


Figure 3.2: Representation of converting the perturbed point  $z = (r, \theta)$  to a point  $z_x, z_y$

### 3.1.2. TRUNCATION

After adding the noise to the data, it cannot be ensured the data is within the original domain (figure 3.3). If this is not the case, the data is easily distinguished by an unwanted adversary [Andrés et al., 2012; Min et al., 2022]. The truncation is an essential part of the mechanism to ensure the data is contained within the domain of the original data  $X$ .

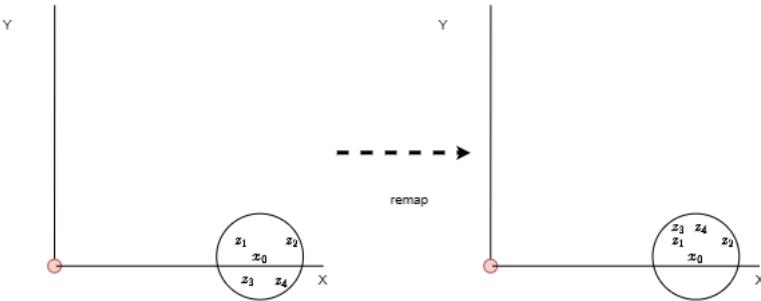


Figure 3.3: Representation of truncation of data points for 2-dimensional Laplace mechanism.

The approach Andres et al. introduces is to remap a perturbed point  $z$  to the closest point  $x'$  in  $G$  [Andrés et al., 2012]. Here,  $G$  is a grid generated using constant cell width. Let the below equation be the collection of probabilities for a point being remapped to a grid cell  $x'$  [Andrés et al., 2012]:

$$R(x) = y \in R^2 | \forall x' \in G \cdot d(y, x') \leq d(y, x') \quad (3.3)$$

The original GI definition contains  $K$ , which is the probability of  $z$  being reported as  $x$  (See Equation: 2.1.1). However, this probability is no longer guaranteed because  $z$  can also be part of  $G$  [Andrés et al., 2012]. Hence the probability is now  $R_X(x) = R(x) \cup X$ .

So, the  $R_X(x)$  has a different shape depending on the distance  $x_0$  and  $x'$  (ergo, it depends on the grid cell width). The probability may be 0, and GI cannot be satisfied.

To overcome this issue, Andres et al. propose a way of calculating  $\epsilon'$ , depending on the grid-cell width  $u \cdot w$ . [Andrés et al., 2012]. They proved this in theorem 4.1 [Andrés et al., 2012]:

**Theorem 1 (Discretization 2D-Laplace)** Assume  $r_{max} < \frac{u}{\delta_\theta}$ , and let  $q = \frac{u}{r_{max}}\delta_\theta$ . Let  $\epsilon, \epsilon' \in R^+$  such that

$$\epsilon' + \frac{1}{u} \ln \frac{q+2e^{\epsilon'u}}{q-2e^{\epsilon'u}} \leq \epsilon$$

Then  $K_{\epsilon'}$  provides  $\epsilon$ -geo-indistinguishability within the range of  $r_{max}$ . Namely, if  $d(x_0, x), d(x'_0, x) \leq r_{max}$  then:

$$K_{\epsilon'}(x_0)(x) \leq e^{\epsilon d(x_0, x')} K_{\epsilon'}(x')(x).$$

Here,  $\delta_\theta$  is the machine's precision.

Do we need to do something with this still?

The theorem states that  $\epsilon'$  is the additional noise needed to satisfy GI with the introduction of discretization. It is sufficient to take  $r_{max}$  as  $diam(X)$ , which is the diameter of the set of points  $X$  if it satisfies theorem 1 [Andrés et al., 2012]: So,  $r_{max}$  is the maximum distance between points in  $X$ .

In the context of our research, it is more convenient to find the  $u$  based on a given  $\epsilon$ . So the challenge is to find the  $u$  that satisfies GI. This can be solved by finding the root of the function:

$$\epsilon' + \frac{1}{u} \ln \left( \frac{\frac{u}{r_{max}} + 2e^{\epsilon'u}}{\frac{u}{r_{max}} - 2e^{\epsilon'u}} \right) - \epsilon = 0 \quad (3.4)$$

The idea of finding the root of a function is to try different values of  $u$  until  $f(u) = 0$ .

Have some problem

Maybe extend this with some more explanation?

### 3.1.3. FINAL MECHANISM

Finally, we provide as means of a summary the final algorithm for the Laplace mechanism for 2D space

---

**Algorithm 1** Full mechanism for perturbing training data for 2D-clustering using planar/2D-Laplace [Andrés et al., 2012]

---

```

Input:  $x \in X$                                      ▷ 2D array of points
 $\epsilon$                                                  ▷ should satisfy Theorem 1
Output:  $z \in Z$                                      ▷ 2D array of perturbed points

 $x_{min} \leftarrow min(X)$ 
 $x_{max} \leftarrow max(X)$ 
 $Z \leftarrow []$ 
for  $point_i \in X$  do
     $\theta \leftarrow [0, \pi/2]$                                 ▷ Random noise for  $\theta$ 
     $p \leftarrow [0, 1]$ 
     $z_i \leftarrow C_\epsilon^{-1}(p)$                          ▷ formula 3.2
     $x_{perturbed} \leftarrow point_{i_x} + (z_{i_x} * \cos(\theta))$  ▷ add noise to x-coordinate
     $y_{perturbed} \leftarrow point_{i_y} + (z_{i_y} * \sin(\theta))$  ▷ add noise to y-coordinate
    append  $x_{perturbed}, y_{perturbed}$  to Z
end for
return Z

```

---

## 3.2. 3D-LAPLACE

The previous sub-section described the use of 2-dimensional noise on geographical data. This approach has recently been extended to support 3-dimensional data, which benefits indoor navigation [Min et al., 2022]. The method is similar to the 2D approach but includes the azimuth angle  $\psi$  besides the polar angle  $\theta$  and radial distance  $r$ .

### 3.2.1. GEO-INDISTINGUISHABILITY

To establish the same privacy guarantees for 3-dimensional data as for 2-dimensional data, the original equation 2.4 is extended [Min et al., 2022].

$$K(x_1)(z) \leq e^{\epsilon * d_3(x_1, x_2)} K(x_2)(z) \quad (3.5)$$

Where  $x_1$  and  $x_2$  are two real data points in the same dataset  $X$ , and  $K$  is the probability.

### 3.2.2. SPHERICAL LAPLACE

Implementing the 3D-Laplace is based on the same method as the 2D-Laplace mechanism. Therefore, the PDF where the noise is drawn from is almost identical [Min et al., 2022]:

$$D_\epsilon(x_i)(x_0) = A \cdot e^{-\epsilon \cdot d_3(x_i, x_0)} \quad (3.6)$$

Where  $A$  is the normalization factor, and  $d_3$  is the Euclidean distance in a 3-dimensional space. Instead of only consisting of a polar angle  $\theta$  and radial distance  $r$ , the third dimension (altitude) is the azimuth angle  $\psi$  [Min et al., 2022]:

$$D_\epsilon(x_i)(x_0) = \frac{1}{4 \cdot \pi^2} \cdot \epsilon^3 \cdot r^2 \cdot e^{-\epsilon \cdot r} \quad (3.7)$$

The distance between  $x_i$  and  $x_0$  is replaced by  $r$ , leaving us with the normalization factor. As with the 2D-Laplace, we are most interested in the normalization factor, which uses the three coordinates [Min et al., 2022]. Because these can be sampled independently, Min et al. formulate three marginal distributions  $D$  for each coordinate:  $R$ ,  $\Theta$  and  $\Psi$  [Min et al., 2022]:

$$D_\epsilon, R(r) = \frac{1}{2} \epsilon^3 \cdot r^2 \cdot e^{-\epsilon \cdot r} \quad (3.8a)$$

$$D_\epsilon, \Theta(\theta) = \frac{1}{\pi} \quad (3.8b)$$

$$D_\epsilon, \Psi(\psi) = \frac{1}{2 \cdot \pi} \quad (3.8c)$$

- The radius  $r$  is sampled randomly from the Gamma distribution (Equation 3.8a) with a shape of 3 and a scale of  $1/\epsilon$  [Min et al., 2022].
- The  $\theta$  and  $\psi$  are also random sampled from the uniform distribution with the range given in Equations 3.8b and 3.8c [Min et al., 2022].

Finally, the noise is added to the original location  $x$  to obtain the perturbed location  $z = x + U * r$  (See Figure 3.5). The spherical coordinates  $(r, \theta, \psi)$  are converted to the Cartesian

coordinates  $(x, y, z)$  to obtain the perturbed location  $z$ :

$$\begin{aligned} z_x &= r * \sin(\theta) * \sin(\psi) \\ z_y &= r * \sin(\theta) * \cos(\psi) \\ z_z &= r * \cos(\theta) \end{aligned}$$

The complete overview is visualized in figure 3.4.

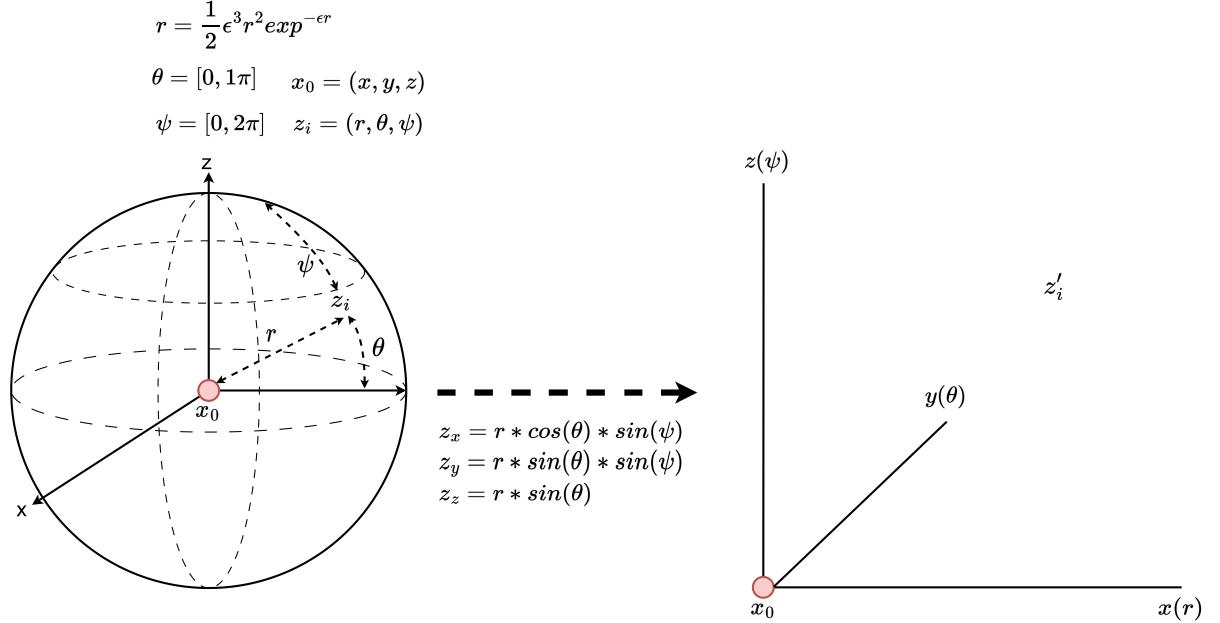


Figure 3.4: 3D-Laplace noise distribution according to the method proposed by Min et al. [Min et al., 2022]

The above figure shows the whole process of generating noise for a single point  $x_0$  (left) to a new point  $z$  (right). As observed, the noise is generated using a sphere and converted to a cartesian coordinate system.

Kernel density plot of a point  $x_0$  (centre) with 50 random generated points according to the formula for 3D-Laplace in a projected sphere with epsilon 3

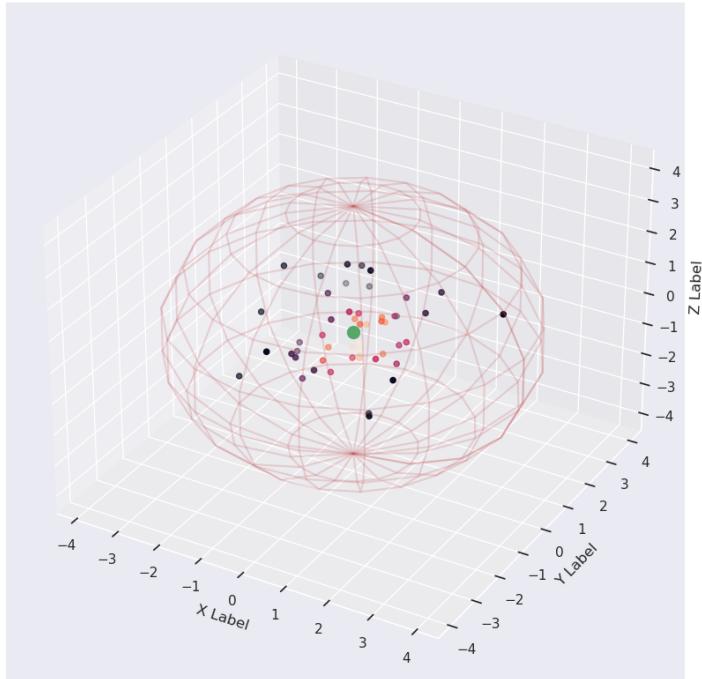


Figure 3.5: 50 random noise samples generated around point  $x_0$  (green dot) using the 3D-Laplace noise method [Min et al., 2022] plotted on a sphere.

This figure shows an example of the noise generated around a point  $x_0$  (green dot). For this purpose, 50 random samples were generated using the 3D-Laplace noise method. As is visible, most points are plotted close to  $x_0$  and fewer at the sphere's surface, according to the Gaussian distribution.

### 3.2.3. TRUNCATION

As with the 2D-Laplace method, the 3D-Laplace method also needs a discretization method that is used to truncate the data. Instead of a plane grid  $G$ , Min et al. extended this to a cuboid grid  $G_c$  to support 3-dimensional data [Min et al., 2022]. This impacts the probability  $K$ , as was the case with the 2D-Laplace method (See Equation: ??). The equation is extended by Min et al. to support 3-dimensional data [Min et al., 2022]:

$$N(x_0) = \{y \in R^3 | \forall'_0 \in G_c : d_3(x'_0, y) \leq d_3(x_0, y)\} \quad (3.9)$$

So, the complete probability for a point  $x_0$  being generated by the 3D-Laplace method to be part of  $G_c$  is:

$$N_B(x_0) = N(x_0) \cap Z \quad (3.10)$$

Here,  $Z$  is the discrete set of points of  $X$  generated according to the 3D-Laplace mechanism. According to the mechanism,  $Z$  depends on the radial distance  $r$ , but now also needs to consider  $G_c$ . The latter stays the same size, while the radius  $r$  differs according to 3.8a. This jeopardizes geo-indistinguishability, and therefore Min et al. provided a modified version of Theorem 1 [Min et al., 2022]. Moreover, they extend the original theorem (See Theorem 1) to support 3-dimensional data. They define the cuboid  $G_c$  using the units  $u, w, h$  (length, width, and height). Now these are established, we provide the extended theorem as provided by Min et al. to show GI holds for the maximum radius of  $r_M$  [Min et al., 2022]:

**Theorem 2** Assume  $R_M < \frac{w}{d_0}$ ,  $d_r \leq \frac{r_M \cdot d_\theta \cdot h}{w}$ , and let

$$w = \frac{u \cdot w}{r_M^2 \cdot \sin \theta \cdot d_\psi \cdot d_\theta} > \frac{v}{r_M \cdot d_\psi}. \text{ Let } \epsilon, \epsilon' \in R^+ \text{ such that}$$

$$\epsilon' + \frac{1}{h} \ln \frac{w+3 \cdot \epsilon' \cdot h}{w-3 \cdot \epsilon' \cdot h} \leq \epsilon$$

Then  $K_{\epsilon'}$  satisfies  $\epsilon$ -geo-indistinguishability within the range

of  $r_M$ . That is to say, if  $d_3(x_1, x')$  and  $d_3(x_2, x') \leq r_M$  then,

$$K_{\epsilon'}(x_1)(x') \leq e^\epsilon \cdot d_3(x_1, x_2) \cdot K_{\epsilon'}(x_2)(x')$$

This theorem is slightly longer than the 2-dimensional variant, but the principle stays the same.

1. The device precision variables are defined as  $d_0, d_\theta$ , and  $d_\psi$ .

Not sure what these mean

2. The cuboid  $G_c$  is defined using the units  $u, w, h$  (length, width, and height). The  $h$  was used for the theorem as minimum grid-unit value, assuming that  $u > v > h$  holds [Min et al., 2022].

3. Finally, the values  $x_1, x_2$  should be lower or equal to the maximum radius  $r_M$ .

What to do with device precision

Maybe explain that we use the root of the function  $f(x) = 0$  for the minimum value  $h$  (height)

We plotted example data points on a 3-dimensional grid in Figure 3.6 to demonstrate this:

Example of generating noise for a dataset  
and remapping it to  $X \subset G$  when outside the domain

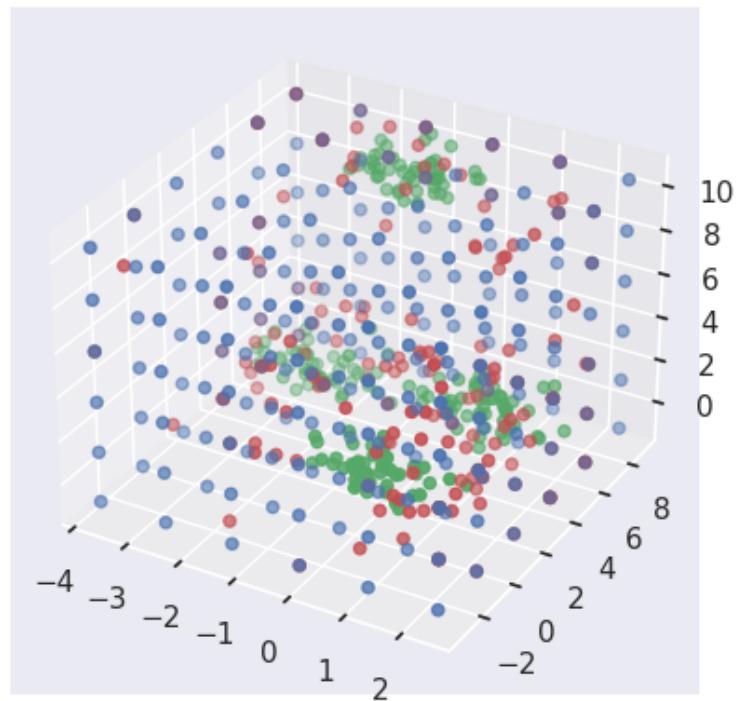


Figure 3.6: Applying 3-dimensional noise with  $\epsilon = 1$  (red dots) to a dataset  $X$  (green dots). Demonstrating remapping to the closest grid point (blue) or  $X$ .

### 3.2.4. FINAL MECHANISM

Finally, we provide as means of a summary the final algorithm for the Laplace mechanism for 3D space

---

**Algorithm 2** Full algorithm for perturbing training data for 3D-clustering using planar/2D-Laplace [Andrés et al., 2012]

---

```

Input:  $x \in X$                                      ▷ 3D array of points
Input:  $l \in R^+$ 
Input:  $r \in R^+$ 
Output:  $z \in Z$                                      ▷ 3D array of perturbed points
 $\epsilon \leftarrow \frac{l}{r}$                            ▷ Calculating privacy budget [Andrés et al., 2012]
 $Z \leftarrow []$ 
for  $point_i \in X$  do
     $\theta \leftarrow [1, \pi 2]$                          ▷ Random noise according to equation 3.7
     $p \leftarrow \frac{1}{2\pi}$                           ▷ Random noise according to equation 3.8
     $r \leftarrow \frac{1}{2}\epsilon^3 * r^2 * e^{-\epsilon * r}$  ▷ Draw  $r$  based on equation ???
     $z_x \leftarrow r * \sin(\theta) * \sin(p)$ 
     $z_y \leftarrow r * \sin(\theta) * \cos(p)$ 
     $z_z \leftarrow r * \cos(\theta)$ 
     $Z ::= Add(z)$                                  ▷ Adds z to the list Z.
end for
return Z

```

---

### 3.3. KD-LAPLACE

Introduce  $d_x$  privacy

As mentioned in the previous chapter, the paper introduced by Min et al. can handle 3-dimensional data. A small recap: a point  $(r, \theta, \psi)$  gives us the three coordinates of a location on the 3-dimensional sphere. An important property is that these coordinates can be generated separately [Andrés et al., 2012; Min et al., 2022]. The  $r$  gives us the radius or distance from  $(\theta, \psi)$  to the center of the sphere<sup>1</sup>. So, instead of having just these two coordinates, we can extend this to n-dimensions by considering an n-hypersphere [Fernandes et al., 2019; Min et al., 2022]. To this end, in addition to points  $\theta$  and  $\psi$ , we consider  $\theta \in S^n$ , where  $S$  is a unit hypersphere. The theorem on extending  $\epsilon$ -geo-indistinguishability to  $\epsilon - d_x$ -privacy is as follows [Fernandes et al., 2019]:

**Theorem 3** Let  $d_x$  be a pseudo-metric on  $X$  and let  $K : X \rightarrow Z$  be a mechanism satisfying  $\epsilon - d_x$ -privacy.

$$K(x)(Z) \leq e^{\epsilon \cdot d_x(x, x')} \cdot K(x')(Z), \forall x, x' \in X, Z \subseteq X.$$

Because  $d_x$  is a pseudo metric, we can also provide  $d_E$  to denote the Euclidean distance. The 2D-Laplace mechanism is referenced as a plane, and the coordinates can be generated separately as  $(r, \theta)$  [Andrés et al., 2012; Fernandes et al., 2019]. Hence, generating the nD-variant can be seen as generating multiple  $(r, \theta)$  pairs [Fernandes et al., 2019]. Where  $r$  is the radial distance from the origin  $x_0$  and  $\theta$  is the angle randomly selected from a unit hypersphere. The selection of  $r$  is according to the Gamma distribution, with shape  $n$  and scale  $\delta > 0$  Fernandes et al. [2019]:

$$Gam_{\delta}^n(r) := \frac{r^{n-1} \cdot \frac{r}{\delta}}{\delta^n (n-1)!} \quad (3.11)$$

Then,  $\theta$  is drawn from a uniform distribution of a unit hypersphere  $S$  [Fernandes et al., 2019]:

$$Uniform^n(\theta) := \frac{\gamma(\frac{n}{2})}{n \cdot \pi^{\frac{n}{2}}} \quad (3.12)$$

This approach is similar to the other variants of the Laplace mechanism. Only now, we select points from an  $n$ -dimensional hypersphere<sup>2</sup> instead of a unit sphere 3.8b.

In the next step, the spherical coordinate representation is converted to Cartesian coordinates [Fernandes et al., 2019]: It is comparable to the way it was done in the previous chapters; however, as there are an  $n$ -amount of angles, the equation is repeated and slightly different:

$$x_1 = r * \cos(\theta_1) \quad (3.13)$$

$$x_2 = r * \sin(\theta_1) * \cos(\theta_2) \quad (3.14)$$

$$x_n = r * \sin(\theta_1) \dots \sin(\theta_{n-2}) * \cos(\theta_{n-1}) \quad (3.15)$$

$$x_n = r * \sin(\theta_{n-1}) * \sin(\theta_{n-2}) * \sin(\theta_{n-1}) \quad (3.16)$$

---

<sup>1</sup><https://mathworld.wolfram.com/SphericalCoordinates.html>

<sup>2</sup><https://mathworld.wolfram.com/SpherePointPicking.html>

The combination of sections 1 and 2 of this chapter provides a good overview of the solution using a similar image as the 2D and 3D variants (figure 3.7).

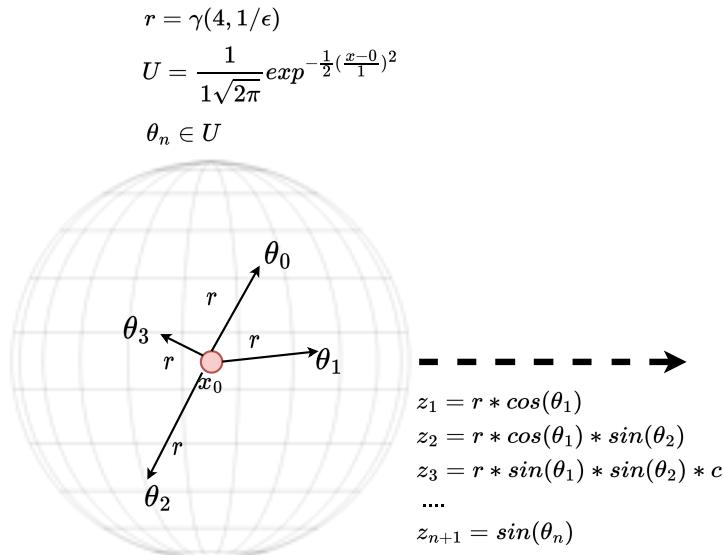


Figure 3.7: Overview of the kD-Laplace mechanism

### 3.3.1. PRIVACY VERSUS UTILITY

Needs refactor

If we continue adding dimensions, we notice the noise is shrinking proportionally. We must first examine the formula for a hypersphere's volume to understand this behavior.

$$S_n = \frac{2\pi^{n/2}}{\gamma(\frac{1}{2}n)} \quad (3.17)$$

$\gamma$  is the Gamma distribution determined based on the number of dimensions  $n$ <sup>3</sup>.

What are these numbers?

As the amount of dimensions increases, the most volume is located on the hypersphere surface. When we convert the points to Cartesian coordinates, some will be located at the center (e.g., 0.5), while others will be close to the surface (e.g., 0.0). However, as the number of dimensions increases, most will be close to the surface (e.g., 0.99). The decreasing amount of volume is illustrated using this figure:

---

<sup>3</sup><https://mathworld.wolfram.com/Hypersphere.html>

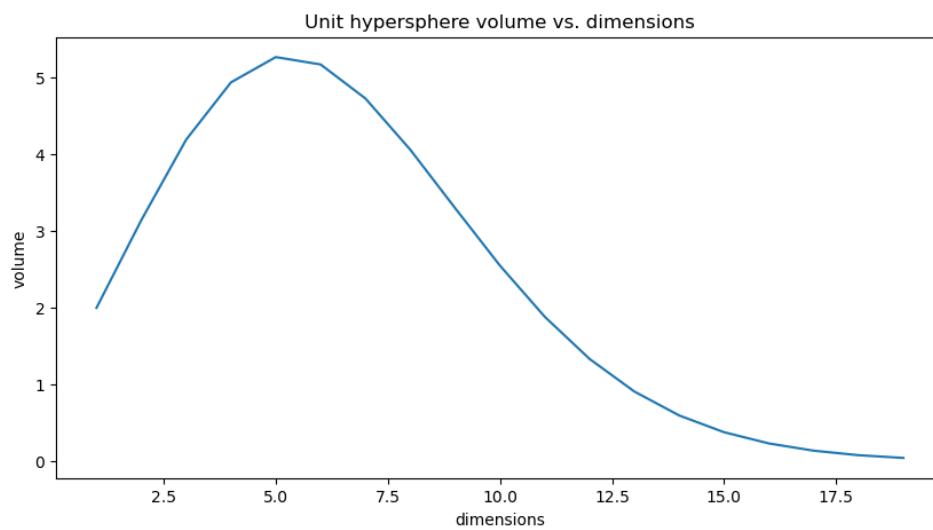


Figure 3.8: Illustration of the decreasing volume while increasing the number of dimensions

The noise decreases as the dimensions increase, increasing utility. Hence, observing the behavior of privacy relative to utility is intriguing. This behavior will be further emphasized in a later stage of this research.

### 3.3.2. GRID-REMAPPING

Needs refactoring with proof and formal representation

In the previous sections, we introduced the 2D, 3D, and nD Laplace mechanisms. The current section will be used to formalize nD-Laplace and introduce modifications to improve the mechanisms for clustering.

#### DISCRETIZATION AND TRUNCATION

Recalling both mechanisms, the 2D version operates on a plane and approximates on a grid  $G$ , while the 3D version works in a 3D space using a cuboid grid. Given a set of input points  $X \subset R^2$ , we can truncate points that are outside the domain by remapping them to points within  $G$  ( $Z = X \cap G$ ) [Andrés et al., 2012]. Here,  $X$  represents other data points reported locally by the same user.

proof & formalize grid-remapping / discretization for nD-Laplace

The utility of this method depends on the number of grid cells in  $G$  since a smaller distance will result in more frequent mapping to the surface of the grid. When  $\epsilon$  is very low (and thus farther away), the data points are more likely to map to the grid surface (3.5, 3.6). Increasing the number of grid cells can improve the utility, but this comes at the cost of significantly increased space complexity. Also, when the number of dimensions increases, the number of grid cells grows quadratic.

#### OPTIMIZATION

To extend Andres et al.'s approach to n-dimensional data, we need an efficient way to search points in an n-dimensional hypersphere. To do this, we adopt the idea proposed by Chatzikokolakis et al. of using a kd-tree to search the grid efficiently [Chatzikokolakis et al., 2017]. Their research describes the theoretical utilization of a kd-tree for searching nearby points for a given point. For this reason, we aim to apply this practically by using a kd-tree for the following tasks:

1. Finding nearby points for  $z \in G$  (Section: 3.3.2).
2. Finding nearby points for  $x \in X$  and  $z \in Z$  (Section: 3.3.2).

For visualization purposes, this section will primarily focus on 2D data. However, it is essential to emphasize that the same algorithm will also be applied to 3D and nD data. The underlying principles and steps of the algorithm remain the same.

The next section explains the concept and underlying idea of kd-trees. Then, we delve into its application for grid remapping. Finally, we expand this to a practical application

## KD-TREES

A kd-tree algorithm can search a grid for nearby points [Bentley, 1975]. It can do so by recursively splitting the grid into a binary tree to search for grid coordinates [Washington, 2002]. In addition, it preserves spatial information of the data so it can be utilized to find nearby points using Euclidean distance (nearest neighbor search). The following example provides an idea of how this works (Figure 3.9):

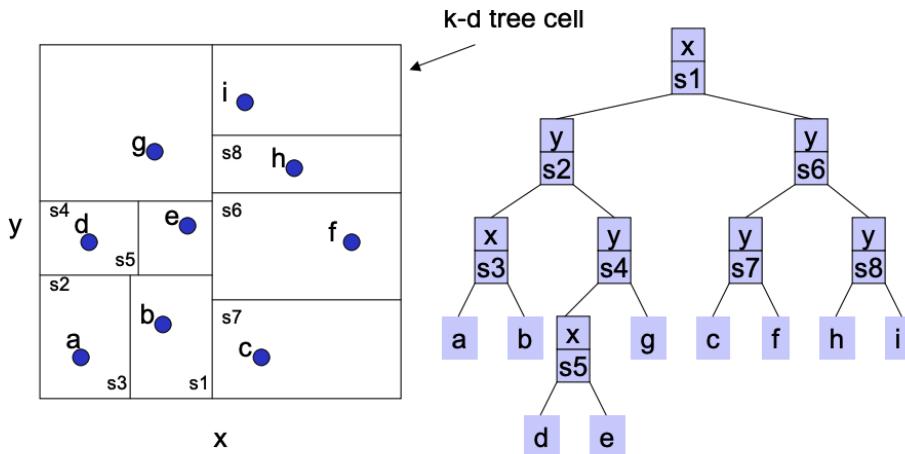


Figure 3.9: Representation of constructing a kd-tree with 2 dimensions [Washington, 2002].

Take, for example, the 2D Laplace algorithm that utilizes a plane (left side). The data points can be divided based on their x and y coordinates. Each coordinate becomes a node in the binary tree, and the grid is divided based on these splits. The binary tree allows us to search the grid efficiently. An example of this is provided in the following image (Figure 3.10):

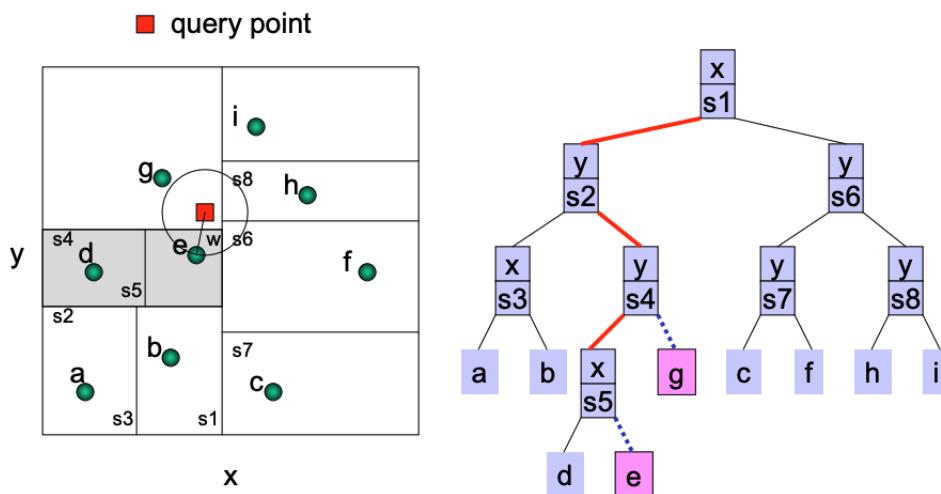


Figure 3.10: Representation of searching a kd-tree with 2 dimensions [Washington, 2002].

In the example, we are searching for all points that fall within the radius of a random query point. The most significant advantage is that this dramatically reduces the complex-

ity of searching. Constructing the kd-tree costs equal to grid-remapping  $O(kn)$ , where  $k$  is the number of dimensions and  $n$  is the dataset size. But, searching for the nearest neighbor is a logarithmic function with time complexity of  $O(\log n)$  [Washington, 2002] (See Figure 17 for an overview of the Big-O notation). This approach is a significant improvement over searching manually, which would have the same complexity as constructing the tree.

### GRID REMAPPING WITH KD-TREE

The kd-tree search method is beneficial for the optimizations we are striving for. This optimization extends grid-remapping for 2d [Andrés et al., 2012] and 3d data [Min et al., 2022] to  $n$ -dimensions.

We have illustrated the three steps required for this below:

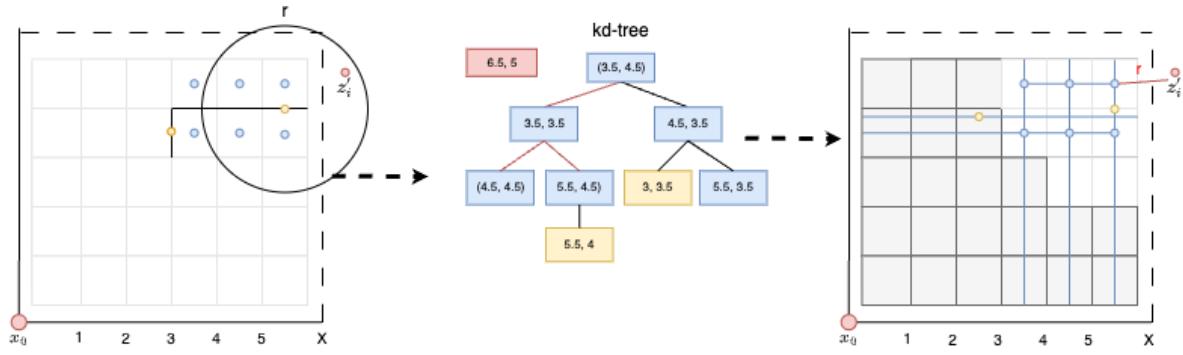


Figure 3.11: Representation utilizing a kd-tree for applying grid-remapping [Andrés et al., 2012]

For easy visualization, the process shown above is with 2-dimensional data. However, this can be extended to  $n$ -dimensional using the kd-tree optimization. The goal is to remap a perturbed point  $z \in Z$ , outside the original domain  $X$ , to the closest point in  $X$  or  $G$  [Andrés et al., 2012].

Firstly, a grid is generated, where each (blue) point represents the center of a grid cell. Together, these centroids form the grid dataset, denoted as  $G$ . The yellow points and  $x_i$  are part of the original collection, denoted as  $X$ . Here,  $r$  represents the radius used to generate a private version of the data point  $x_i$ , named  $z_i$ , based on 2D-Laplace. In the illustration, you can observe that  $z_i$  falls outside the original domain of  $X$ , so it needs to be remapped. We accomplish this by utilizing the nearest-neighbor search from the kd-tree algorithm, allowing us to search in  $X \cup G$ . Using this algorithm, we can effectively remap point  $z \in Z$  to either  $X$  or  $G$  based on the closest Euclidean distance (Algorithm 4 and Algorithm 3).

Although the methodology works for scaling grid-remapping for  $n$ -dimensional data, it can still cause scalability issues. It still does not scale well with the number of dimensions and a more significant number of grid cells. Therefore, we are looking for a solution that might work better in practice

---

**Algorithm 3** Algorithm for finding points outside the domain of  $X$ .

---

**Input:**  $x \in X$  ▷ original dataset  
**Input:**  $z \in Z$  ▷ perturbed dataset

$X_{domain} \leftarrow \text{KDTREE::QUERY}(Z)$  ▷ find the closest points.  
 $X_{features} \leftarrow \text{LIST::GETFEATURES}(X)$  ▷ retrieve dataset dimensions

$X_{outside-domain} \leftarrow []$

**for**  $feature \in X_{features}$  **do** ▷ iterate over all features and check if any points are outside the domain.

**if**  $feature \leq X::\text{MIN}(Z)$  **then**  
        ROW::APPEND( $X_{outside-domain}$ )  
    **end if**  
    **if**  $feature \geq X::\text{MAX}(Z)$  **then**  
        ROW::APPEND( $X_{outside-domain}$ )  
    **end if**  
**end for**

**return**  $X_{outside-domain}$  ▷ The index of points outside the domain of  $X$ .

---

**Algorithm 4** Algorithm for generating and remapping to a grid.

---

**Input:**  $x \in X$  ▷ original dataset  
**Input:**  $z \in Z$  ▷ perturbed dataset  
**Input:**  $grid$  ▷ grid structure ( $n * m$ )

$d_X = Z::\text{DIST}(X)$  ▷ euclidean distances between  $X$  and  $Z$ .  
 $d_{grid} = Z::\text{DIST}(grid)$

$Z_{out-domain} \leftarrow X::\text{FINDPOINTSOUTSIDEDOMAIN}(Z)$  ▷ Algorithm 3

$grid_{tree} \leftarrow \text{KDTREE::BUILDTREE}(grid)$  ▷ create kd-tree

$grid_{mask} \leftarrow \text{KDTREE::QUERY}(Z)$  ▷ find indices of  $z \in Z$  that are closeby grid cells.

$Z_{grid-mask} \leftarrow Z_{out-domain} \cup d_{grid} < d_X$  ▷ All points  $z \in Z$  that are closeby grid cells and are outside domain.

$Z' \leftarrow Z[grid[grid_{mask}][Z_{grid-mask}]]$  ▷ combinate masks to set appropriate indexes to  $g \in grid$ .

**return**  $Z'$

---

### 3.3.3. DENSITY REMAPPING

As discussed, the remapping will be performance intensive to provide good utility, so we adopt the optimal remapping [Chatzikokolakis et al., 2017]. This thesis will refer to this as density remapping because this better explains the application.

According to the grid-remapping methodology, a point  $z_i$  is mapped to the center of the grid cell. The utility is improved by reducing the distance between  $z_i$  and  $x_i$ , as this is closer to the original data point. But, to consider the privacy of the data, this can only be done if there is enough indistinguishability between  $z_i$  and  $x_i$ . To this end, we adopt the remapping approach of Chatzikokolakis et al. to remap based on the data density [Chatzikokolakis et al., 2017]. The remapping algorithm works on the idea of crowded places ??, with the intuition that a crowded place leverages indistinguishability by crowdedness (density) [Chatzikokolakis et al., 2017].

The approach is visualized using the following figure:

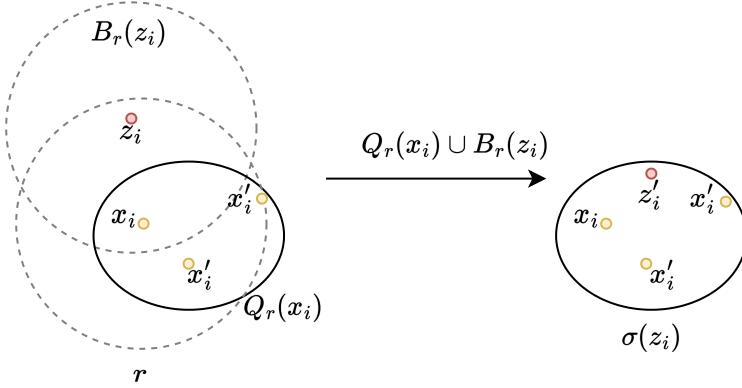


Figure 3.12: Representation of density remapping [Chatzikokolakis et al., 2017]

The first step is calculating  $B_r(z_i)$ , which refers to all the data points within the radius  $r$  around the data point  $z_i$  [Chatzikokolakis et al., 2017]. Then, the algorithm collects the data points around  $x_i$ , again using  $r$ , to determine the density of  $x_i$ . This is the convex hull (collection) of all the original data points within the radius  $r$  around  $x_i$  and is denoted as  $Q_r(x_i)$ . Finally, we obtain  $Q_r$  which is a union of  $Q_r(x_i)$  and  $B_r(z_i)$ . Now that we have the sets of points around  $x_i$  and  $z_i$ , we can calculate the density [Chatzikokolakis et al., 2017]:

$$\forall x'_i \in Q_r \quad \sigma(x_i) = \frac{w(x'_i) e^{-\epsilon d(x'_i, z_i)}}{\sum_{q \in Q_r} w(q) e^{-\epsilon d(q, z_i)}} \quad (3.18)$$

$w(q)$  is the weight of a point  $q \in Q_r$ , and  $w(x_i)$  the weight of a point  $x_i \in X$ . Here, the weight can be a point of interest [Chatzikokolakis et al., 2017]. But, in the context of our study, we consider the points within  $r$  of  $x_i$  to be the weight.

In equation 3.18, the  $w(x)$  and  $w(q)$  are normalization factors for estimating the likelihood of a point  $x'_i$  being close to  $z_i$  and  $q$  being close to  $z_i$ . So, this remapping is according to the original definition of the 2D Laplace PDF definition 3.1. The outcome of the formula is a collection of values that indicate the degree of density  $\sigma(x_i)$ . The new value  $z'_i$  is calculated by taking the mean of the  $\sigma(x_i)$  [Chatzikokolakis et al., 2017]:

$$z'_i = \sum_{x'_i \in Q_r} \sigma(x'_i) * x'_i \quad (3.19)$$

By applying this formula, the new  $z'_i$  is closer to  $x_i$  to minimize the expected loss of utility [Chatzikokolakis et al., 2017]. Much like the 2D and 3D Laplace methods, which use the likelihood of a point  $z_i$  being close to  $x_i$  [Andrés et al., 2012; Min et al., 2022].

---

**Algorithm 5** Algorithm to implement the density remapping of  $z \in Z$  to be in the domain of  $x \in X$

---

**Input:**  $x \in X$

**Input:**  $z \in Z$

**Input:** *epsilon*

**Output:**  $z' \in Z$

```

1:  $Z' = FindRemappedPoints(Z)$                                      ▷ Algorithm: 3
2:  $tree \leftarrow KDTree(X)$ 
3: for  $z' \in Z'$  do
4:    $r = \text{FINDRADIUS}((z'))$                                      ▷ Get original radius  $r$ .
5:    $X_r \leftarrow \text{KDTree::QUERY}(x)$                                 ▷ find  $q \in X$  around  $x$  with radius  $r$ .
6:    $B_r \leftarrow \text{KDTree::QUERY}(z')$ 
7:    $\sigma(x) = []$ 
8:    $Q_r = X_r \cap B_r$ 
9:    $w_x = \text{LENGTH}(X_r, B_r)$                                      ▷ weight is simply adding density of  $X_r$  and  $B_r$ .
10:  for  $q \in Q_r$  do
11:     $q \leftarrow \text{KDTree::QUERY}(q)$ 
12:     $w_q = \text{LENGTH}(Q_r)$                                          ▷  $w_q$  is the density of each point  $q$  within  $Q_r$ .
13:     $\sigma(w_x) \leftarrow \text{REMAP}(w_x, \epsilon)$                          ▷ Equation 3.18.
14:     $\sigma(x) \leftarrow \text{APPEND}(\sigma)$                                 ▷ add to the list  $\sigma(x)$ .
15:  end for
16:   $z' \leftarrow \text{AVERAGE}(\sigma(x))$                                ▷ Equation: 3.19.
17: end for

```

---

- 1: Receives all the points  $z \in Z$  that are outside the domain of  $X$ .
- 2: Constructs a KDTree from the original data  $X$ , so it can be queried.
- 3: Receives the original privacy radius  $r$  that was used to generate  $z$  from  $x$ .
- 4:  $X_r$  is the set of points  $x \in X$  that are within the radius  $r$  of  $x$ .
- 5:  $B_r$  is the set of points  $z \in Z$  that are within the radius  $r$  of  $z$ .
- 7:  $w_x$  is the "popularity" of  $x$  and  $z$ , for the number of points is counted within the radius  $r$ .
- 12:  $w_q$  is the weight of each point  $q \in Q_r$ , so  $q$  is re-assigned to be the number of points within the radius  $r$  of  $q$ .

The equation for density remapping 3.18 was created for 2-dimensional data. However, we aim to use this same approach for 3-dimensional and n-dimensional data. To do so, we first revisit the PDF for 3D-Laplace, which was defined using this Equation: 3.6. Few modifications are necessary, as the normalization factor  $A$  is now defined using Equation: 3.18.

Formulate the definition of  $A$

Check if this can be extended to nD-Laplace

### 3.3.4. PUTTING IT TOGETHER: KD-LAPLACE

Now that we have defined everything, we can write the algorithm in a step-by-step manner:

---

**Algorithm 6** Full algorithm for perturbing training data for nD-clustering using kd-Laplace

---

**Input:**  $x \in X$  ▷ n-dimensional array of original points  
**Input:**  $\epsilon$  ▷ privacy budget  
**Output:**  $z \in Z$  ▷ n-dimensional array of optimal-remapped perturbed points  
 $sphere = \text{GENERATEUNITSPHERE}(x)$  ▷ construct a sphere around  $x$ .  
**for**  $row \in X$  **do**  
     $d \leftarrow \text{LENGTH}(row)$  ▷ amount of dimensions  
     $r \leftarrow \text{GENERATERADIUS}(d)$  ▷ generate radius  $r$  using Equation ??.  
     $sphere \leftarrow \text{GENERATETHETA}(d)$  ▷ perturb the sphere Figure 3.7.  
     $noise \leftarrow \text{CARTESIAN}(row, epsilon)$  ▷ use Equation 3.16 to perturb  $row$  and convert to cartesian.  
     $z = x + noise$   
    APPEND( $Z, z$ )  
**end for**  
**return**  $Z$

---

It is important to note that with the introduction of this algorithm, it is no longer a non-interactive method, as the data points now interact with each other. Therefore, we expect the introduction of optimal grid remapping with kd-tree will bring more utility [Wang et al., 2020; Xiong et al., 2020b]. The introduction of the kd-tree is also why we present our method as *kd-Laplace*, where  $k$  is the number of dimensions.

Also refactor this

### 3.3.5. MECHANISM FLOWCHART

All formulas and theories are established for 2D, 3D, and nD-Laplace, so the mechanism design applies to all three variants: For easy navigation, we provide a list of all algorithms:

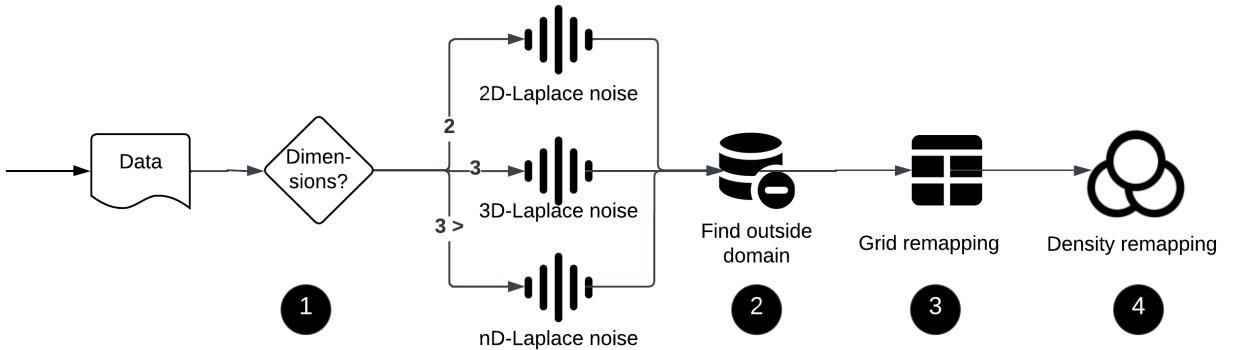


Figure 3.13: Non-interactive mechanism design for kD-Laplace.

1. Based on the number of dimensions, the algorithm decides the correct Laplace mechanism to use:
  - 2D-Laplace: 1
  - 3D-Laplace: 2
  - nD-Laplace: 6
2. Find points outside domain: 3
3. Grid remapping: 4
4. Density remapping: 5

#### PRACTICAL EXAMPLE

The shape of the dataset is necessary for the usefulness of clustering. With our algorithm, there are four different shapes/variants of the dataset. For example, this has been visualized using a 3D dataset based on the heart dataset (5.1.3). Our mechanism aims to provide privacy and preserve the dataset's shape to benefit the utility of clustering. Grid remapping and optimal remapping are used to achieve this goal.

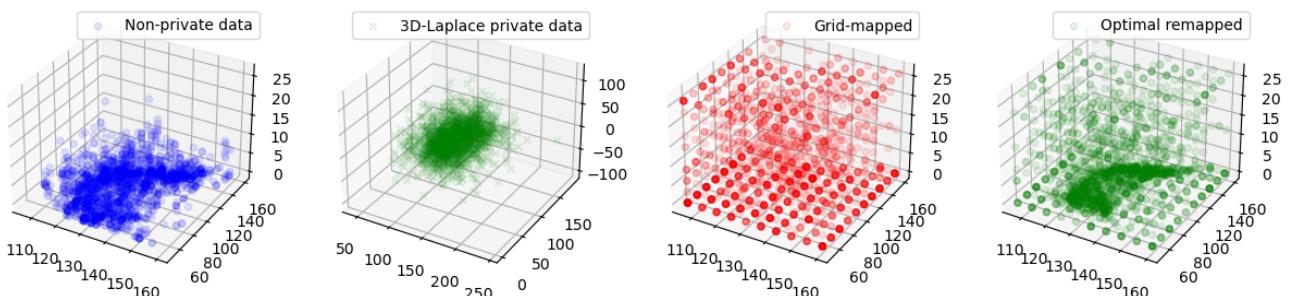


Figure 3.14: Example of optimal remapping for the 3D-dataset: Cardiotocography. The example shows the different steps of the mechanism in sequence for a dataset perturbed with a privacy budget of 0.1.

1. Dataset: the blue dots represent the original dataset without any modifications.
2. Adding noise: the green crosses represent the dataset after adding noise; for this particular example, this is 3D-Laplace (Algorithm 2): As can be observed, the data is generated from the center, causing many data points to fall outside the original domain of the dataset.
3. Grid-remapping: the red dots represent the dataset after grid-remapping (Algorithm 4) After performing the grid remapping algorithm, all points within the domain are plotted. However, the original shape of the data is mostly lost. This makes it challenging to cluster the data as was possible with the original data.
4. Optimal-remapping: the green dots represent the dataset after optimal-remapping (Algorithm 5). After completing the previous step, the data points are again remapped based on the (original) density. This results in restoring the original shape of the data and, consequently, the clusters.

# 4

## ATTACKS ON PRIVACY

This chapter is devoted to investigating and evaluating attacks on machine learning models. Differential privacy protects the centrally stored dataset from leaking sensitive information. Therefore, assessing the mechanism's privacy is best measurable using common attacks [Jayaraman and Evans]. In this context, we evaluate attacks that explicitly uncover training data from a privately trained model. We consider two types of attacks:

1. **Membership inference attack:** An adversary attempts to infer whether a data point was used for training.
2. **Reconstruction attack:** An adversary attempts to reconstruct the training data using the model.

The knowledge of the attacker (adversarial knowledge) is an important factor to consider. This knowledge can be divided into white-box and black-box approaches [Hu et al., 2022].

1. **White-box:** The attacker has all the needed data. Including target model parameters or the architecture [Hu et al., 2022].
2. **Black-box:** The attacker has limited information, like training data distribution and the trained model [Hu et al., 2022].

We will discuss both types of attacks and types in the next two sections.

### 4.1. MEMBERSHIP INFERENCE ATTACKS

Formalize this

An attack model that plays a significant role in machine learning is **Membership Inference Attack (MIA)**. With this attack, an adversary attempts to infer the training data  $x \in X$  (member) from a given data point  $z \in Z$  (non-member). The attack happens exclusively on supervised learning models, which predict labels or probabilities. Most attacks on models trained on a centralized dataset occur during the inference phase, where the trained model is used to make predictions. [Rigaki and Garcia, 2021]. This is also why we are primarily interested in this phase, as we are not using a distributed learning model.

The most well-known member inference attack is training shadow models [Rigaki and Garcia, 2021]. In this attack, an attacker trains multiple models. These models do not necessarily have to be the same as the original model, and the focus is mainly on the data input/output. It is a black-box attack, but the attacker often needs knowledge of the data distribution to create a good shadow dataset [Rigaki and Garcia, 2021].

One of the earlier works that used this attack was Shokri et al. [Shokri et al., 2017]. An attacker trains multiple models (shadow model) to overfit the original modal. This idea is based on the model giving higher scores to the data on which it was trained (overfitting). Using this approach, attackers can retrieve the model's training data (member data) by injecting many fake data (non-membership data).

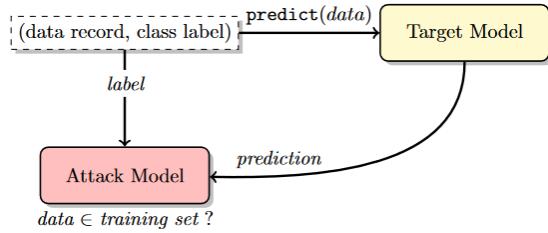


Figure 4.1: Black-box MIA attack on a machine learning model [Shokri et al., 2017]

Another approach to a black-box attack was introduced by Peng et al. and only considers that the attacker has access to the already trained model. They first rescale the probabilities using temperature scaling to compensate for overconfident models [Peng et al.]. So instead of having a probability between two classes with, for example, 99% against 1%, it will be more evenly distributed based on the training data. They then cluster the probabilities into two clusters using K-Means and label the higher confidence scores as members. Another take on this is prediction and confidence-based MIA, both proposed by [Yeom et al., 2018]. They assume that an attacker knows the standard error and has access to the perturbation dataset. The algorithm can then extract the truth label by minimizing the loss.

## 4.2. RECONSTRUCTION ATTACK

Also formalize

The concept of reconstruction attacks predates differential privacy, as this principle also gave rise to the idea of necessary database privatization [Dinur and Nissim, 2003]. An adversary could reconstruct training data from a given (classifier) model using a reconstruction attack. Their research evaluates the perturbation that needs to be added to a database to protect it versus a reconstruction attack.

A general reconstruction attack for our use-case is the attribute inference attack [Dwork et al., 2017] or model inference [Rigaki and Garcia, 2021]. Both terms are essentially the same, and we have chosen to name attribute inference attack, as it is the most common one in most literature [Jegorova et al., 2022]. Attribute inference focuses on a reconstruction attack with the adversary's goal of retrieving the secret from each user [Dwork et al., 2017]. For example, an attacker may attempt to reconstruct information about someone's heart disease using the individual's properties.

A practical implementation of the attack was provided by Fredrikson et al. as a way to infer sensitive features [Fredrikson et al., 2015]. To accomplish this, they used a decision tree attack, a white-box approach, as they also accessed the count of instances for each decision tree branch. They also considered a black-box approach with access only to the target model (ML-as-a-service in this example). The attack targets gradient descent, which is used to optimize the input data of an attack to mimic the original data. It has only been shown to work on a neural network for face identification images (named MIFace) [Fredrikson et al., 2015], but it could be extended to another machine learning classifier if it uses gradient descent (e.g. Support Vector Machines) [Nicolae et al., 2019]. A more general approach was undertaken by Yeom et al., building upon the same research discussed in the previous section (Member inference). In their work, Yeom et al. proposed a membership inference attack, which can also be utilized for attribute inference in a white-box setting. The attribute inference attack follows a similar methodology, wherein the target model is queried to assess the loss of synthetic data based on adversarial knowledge. Repeating the process, the attack identifies and selects the value with the highest prior probability, incorporating membership information [Jayaraman and Evans, 2022; Yeom et al., 2018]. Evaluating the effectiveness of such attacks is challenging, as it relies on the correlation between attributes, irrespective of whether the data belongs to the training or test dataset [Zhao et al., 2021]. Therefore, Yeom et al. focus on the similarity between membership/attribute-inference and combine them to evaluate the attribute inference attack [Yeom et al., 2018].

Differential privacy aims to introduce sufficient noise to mitigate the risk of various attacks [Dwork et al., 2017; Jayaraman and Evans]. Assessing the privacy leakage of our model can be effectively accomplished by employing attribute/membership inference techniques. Within this thesis, our focus is mainly on MIA. We aim to establish a quantifiable measure for our mechanism, and combining both attribute and membership attacks is not beneficial since they essentially capture the same information. Therefore, it is sufficient to concentrate on membership inference attacks.

## 4.3. ATTACK EVALUATION

In this section, we evaluate the membership inference attack and evaluate it as it is the most appropriate for this study (See previous chapter). We assess whether differential privacy provides protection against the attack and discuss how this can be measured.

### 4.3.1. MEMBER INFERENCE ATTACKS

Most current research for **MIA** is evaluated for neural networks [Rigaki and Garcia, 2021]. A tiny percentage evaluates this attack for supervised learning, with the majority using classification with decision trees. Most studies have used a black-box approach [Rigaki and Garcia, 2021] for these attacks. This approach is not surprising, as these attacks have a high success rate and pose a greater risk of exploitation.

Introducing differential privacy reduces the impact of a member inference attack [Hu et al., 2022; Rigaki and Garcia, 2021]. This is because the input to the model is perturbed. While it is still possible to retrieve the training data, the leaked privacy is significantly reduced. A simple but effective way to measure the privacy leakage is by calculating the accuracy of correctly predicting membership by the adversary [Choquette-Choo et al., 2021].

Yeom et al. created a metric specifically for membership inference attacks which can be measured using an "adversarial advantage." This metric describes the percentage of privacy compromised during a member inference attack [Yeom et al., 2018]. This metric is calculated by subtracting the **False Positive Rate (FPR)** from the **True Positive Rate (TPR)**. The **TPR** represents the number of correctly predicted member data (training data), and the **FPR** represents the number of correctly predicted non-member data. Although these metrics are commonly applied in the literature for **MIA**, they do not provide enough information [Carlini et al., 2022]. Both metrics do not consider the imbalance between the **TPR** and **FPR**. The metric should emphasize the **TPR**, as this is the percentage of correctly predicted member data. Therefore, they propose to use a ROC curve to show the effectiveness of a membership inference attack [Carlini et al., 2022].

In conclusion, the attacks that use **MIA** are all (mis)using supervised machine learning. However, in this study, we use clustering algorithms. Therefore, a semi-supervised approach can be used, as illustrated in Figure 4.2.

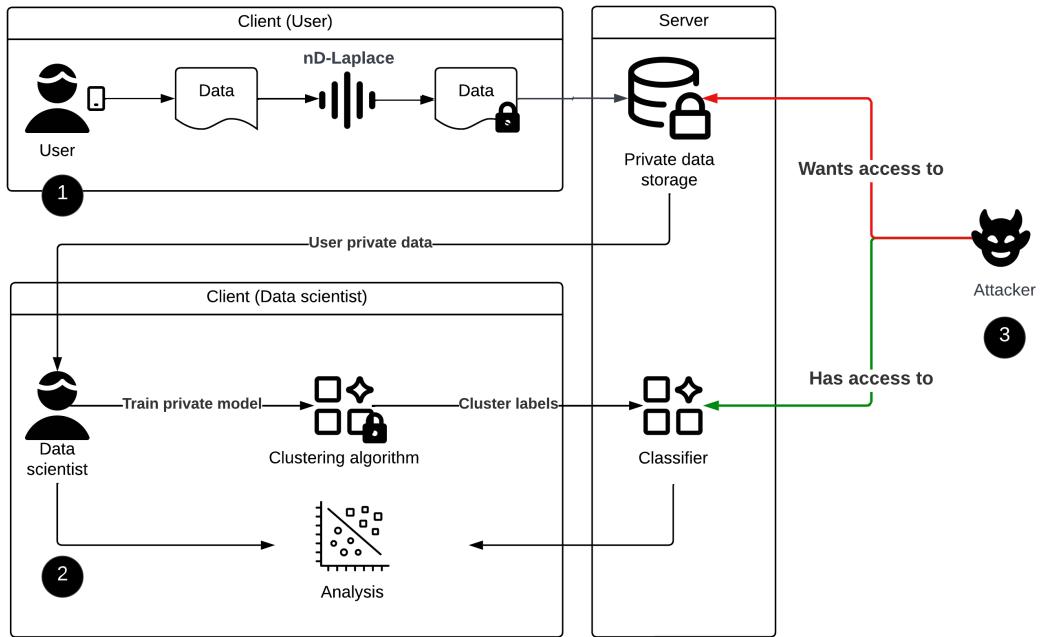


Figure 4.2: Semi-supervised black-box approach to execute a member inference attack.

1. The user uses a client (e.g., mobile app), where the nD-Laplace noise is added locally to the data.
2. A data scientist trains the clustering algorithm with the privatized data. Therefore, the clustering algorithm's input and output (labels) are private. These labels are used to train a classifier (semi-supervised setup).
3. The attacker wants access to the private data on the server. In this approach, the attacker has access to the classifier on the server (black-box setting). So, the attacker can use the classifier's output to conduct a **MIA**.

# 5

## METHODOLOGY

We conducted experiments to gain insights into the proposed methods for researching the appliance of nD-Laplace for clustering algorithms. The experiment results are used to evaluate our method against other literature. In this chapter, we explain:

1. Research design: The background and research questions are explained. Furthermore, we explain the datasets and privacy mechanisms used in the experiments.
2. Data collection: The methods used to collect the data necessary to answer the research questions. These are methods for running experiments to collect data about utility and privacy.
3. Data analysis: We explain how we analyze the data collected in the previous step.

### 5.1. RESEARCH DESIGN

This section describes the overall setup of the research.

#### 5.1.1. RESEARCH BACKGROUND

Clustering data is a crucial part of machine learning, aiming to discover patterns in the data that may not be immediately visible to the human eye. To achieve this, the data is trained on large datasets. However, storing such data is unsafe, as unauthorized access could lead to breaches.

LDP has been introduced to safeguard user privacy while preserving data patterns. Adding local noise means the plain data is never stored on a central server. But, the limited available data in the local view distorts the patterns excessively in each noise addition.

Current literature has attempted to address this challenge, but the solutions fall short. They either require modifications to the clustering algorithm or lack practical applicability due to the use of synthetic data. Existing literature focuses mainly on 2-dimensional data with K-Means as the clustering algorithm.

To address these limitations, we propose a new LDP method called nD-Laplace, utilizing distance to preserve data patterns using the concept of GI.

Our method enables secure training of clustering algorithms based on input-perturbation. Due to this approach, we can train multiple clustering algorithms without modifications

and extend them to n-dimensional data. We will test three variants of the nD-Laplace mechanism to assess their utility and privacy characteristics.

1. nD-Laplace: Plain mechanism without any additions for truncation.
2. grid-nD-Laplace: Mechanism with truncation using grid-remapping (See Equation: 4).
3. density-nD-Laplace: Mechanism with truncation using density-remapping (See Equation: 5).

We will compare the kD-Laplace mechanism to Piecewise, which is another LDP mechanism (See Section 2.3) and assess the effectiveness of the proposed kd-Laplace variants.

This research combines real-world and synthetic data to demonstrate practical applicability. Extensive external and internal validation experiments will be conducted to establish the method's practical relevance and effectiveness.

### 5.1.2. RESEARCH QUESTIONS

The main question of the research is as follows:

*"How can the nD-Laplace mechanism be applied in training privacy-preserving clustering algorithms on distributed n-dimensional data?"*

This question is divided into four smaller sub-questions:

1. *RQ1: How can the 2D-Laplace, 3D-Laplace, and nD-Laplace approaches be adapted to train privacy-preserving clustering algorithms?*  
This research question investigates the applicability of the different variants of the nD-Laplace mechanism for training clustering algorithms.
2. *RQ2: How can the noise generated by nD-Laplace outside the data boundary be remapped to inside the data domain?*  
This research question extends grid-remapping (See Equation: 4) and density-remapping (See Equation: 5) to be used with nD-Laplace for clustering.
3. *RQ3: How can one optimize the computational complexity when implementing grid and density-remapping using the nD-Laplace mechanism?*  
This research question investigates the computational complexity of the grid-remapping and density-remapping algorithms.
4. *RQ4: How do dataset characteristics impact the nD-Laplace mechanism for utility and privacy?*  
In this research question, we evaluate the dataset characteristics and their influence on nD-Laplace. Also, we investigate if we can improve the utility and privacy of kd-Laplace by adopting grid-nD-Laplace and density-nD-Laplace (??). We use several hypotheses questions to answer research question 3.  
In the last two hypotheses, we specifically focus on the K-Means clustering algorithm to test our hypotheses. The choice of K-Means is due to its simplicity, making it suitable for our hypothesis tests. Also, we assume that the K-Means results are representative of the other clustering algorithms.

- (a) *H1: Adding remapping based on density improves utility without sacrificing privacy:* Since the noise is updated based on the density of the data points, the clusters stay preserved (5). As a result, the utility increases while the data points remain indistinguishable. Therefore, we have extended our mechanism with density remapping and named it density-kD-Laplace.
- (b) *H2: The privacy leakage (adversary advantage) and utility increases for more dimensions:* The hypothesis arose from our observation while investigating/implementing kD-laplace(3.3.1). Adding more dimensions (5 >) decreases the added noise gradually, increasing utility. On the other hand, privacy is expected to decrease because the noise is less.
- (c) *H3: The shape of the data negatively impacts the kd-Laplace mechanism in terms of privacy and utility:* Because the mechanism heavily relies on Euclidean distance, the shape (distance between the data points) can ultimately harm utility and privacy. We have already observed this difference when comparing the heart and seed datasets in research questions 1 and 2. But, to rule out the potential impact of the number of data points (the heart dataset has 2126 samples and the seed-dataset 210), we generate three synthetic datasets, each with 1000 samples (See section 5.1.3).

### 5.1.3. DATASETS

For this research, we selected two real-world datasets based on the related papers (2.3). The datasets are sourced from the UCI Machine Learning Repository [noa, c].

1. **Seeds dataset**<sup>1</sup>: This dataset was used in several related works and contains 210 samples with 7 (numerical) attributes. The dataset contains information about seeds, like kernel width and density. We conducted experiments with 2, 3, and 7 dimensions and decided to use the following features (based on the correlation between the features):
  - (a) 2-dimensional data: area and perimeter.
  - (b) 3-dimensional data: the kernel's area, perimeter, and length.
  - (c) 7-dimensional data: All numerical features.
2. **Heart dataset**<sup>2</sup>: This dataset is selected because of the mixed data and amount of instances. It has 23 attributes, including ten numerical attributes and 2126 samples. The dataset contains information about measurements of fetal heart rate (FHR) and uterine contraction (UC). We conducted experiments with 2, 3, and 10 dimensions and decided to use the following features (based on the correlation between the features):
  - (a) 2-dimensional data: FHR baseline and histogram-min.
  - (b) 3-dimensional data: FHR baseline, histogram-min, and accelerations.
  - (c) 10-dimensional data: All numerical features.

---

<sup>1</sup><http://archive.ics.uci.edu/ml/datasets/seeds>

<sup>2</sup><https://archive.ics.uci.edu/ml/datasets/cardiotocography>

The following datasets are for the H3 (shape) hypothesis investigated in research question 3. For this purpose, we generated a 2-dimensional and 3-dimensional dataset with 1000 samples.

**Skewed looks like circle**

**Extend to n-dimensions**



Figure 5.1: Synthetic datasets with 1000 samples and 2-dimensions

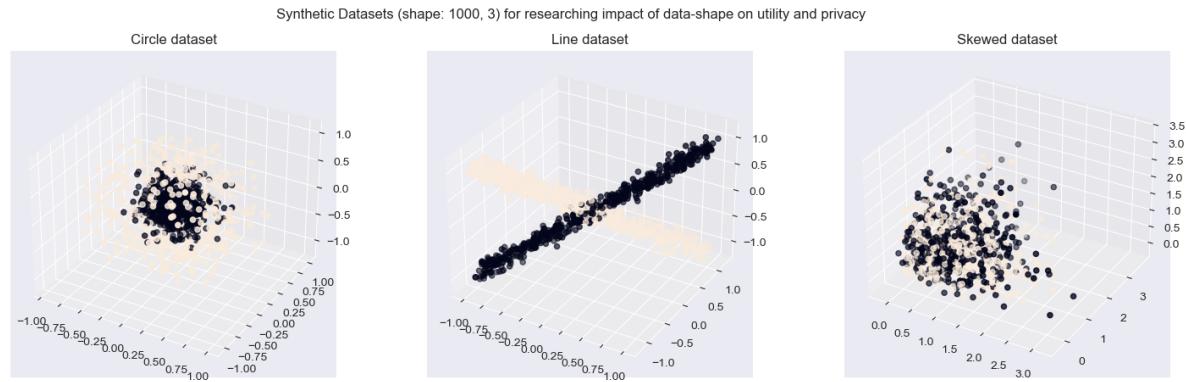


Figure 5.2: Synthetic datasets with 1000 samples and 3-dimensions

1. **Circle dataset:** Dataset with a circle shape.
2. **Line dataset:** Dataset is a line (shape is much like the seeds dataset).
3. **Skewed dataset:** Dataset samples are skewed/concentrated to one side (shape is much like the heart dataset).

#### 5.1.4. EXPERIMENT SETUP

We evaluate all privacy mechanisms by comparing them in utility and privacy. To reduce the measurement bias of results, we executed them ten times for multiple privacy budgets and reported the average for each [Huang et al., 2021].

1. The experiments run ten times, and we report the mean.
2. All experiments run for multiple epsilons: 0.5, 0.7, 1, 1.5, 2, 3.5, 5, 7, 9.

### 5.1.5. ENVIRONMENTAL SETUP

For running the experiments, we use 16GB RAM and an i7-10750H 2.6Ghz processor. The experiments are run using a Docker container which runs a pre-configured distribution of Linux Alpine. It includes a pre-installed Anaconda environment for Python <sup>3,4</sup>. We run the container using the dev-container feature for visual-studio code <sup>5</sup>. This allows us to create a reproducible experiment environment.

#### LIBRARIES & CODE VERSIONS

We use Python version 3.9.13 with Jupyter Notebook for creating a reproducible experimental environment. The packages for Python are:

- Scikit-learn: 1.0.\*
- Yellow-brick: 1.5
- Numpy: 1.24.\*
- Pandas: 2.0.\*
- Seaborn: 0.11.\*
- Matplotlib: 3.5.\*

## 5.2. DATA COLLECTION

This section explains what methods/ algorithms we use to collect the data necessary to answer the research questions.

### 5.2.1. HYPERPARAMETER SELECTION

For the three different algorithms: K-Means, AP, and OPTICS, we analyzed the most important decisions regarding hyperparameter selection (See section ??). This section gives a short list and explanation of the parameters we used throughout the experiments.

#### K-MEANS

Parameter	Description	Value	Dataset
K-value	Calculated based on an "elbow" plot.	k=4 (see Figure 1)	Seeds dataset
K-value	" "	k=4	Heart dataset
K-value	" "	k=5 (see Figure 2)	circle dataset
K-value	" "	k=4 (see Figure 3)	line dataset
K-value	" "	k=4 (see Figure 4)	skewed dataset

Table 5.1: K-Means hyperparameters for dataset 1 - 3

The above table shows the relevant hyperparameters for the K-Means clustering algorithm. The hyperparameter differs for each dataset, so we must display it for each. This

<sup>3</sup><https://github.com/devcontainers/images/tree/main/src/anaconda>

<sup>4</sup>tag: mcr.microsoft.com/devcontainers/anaconda:0-3

<sup>5</sup><https://code.visualstudio.com/docs/devcontainers/containers>

hyperparameter is the k-value used to determine the number of clusters. The "value" column shows the value we used for the experiments, with the corresponding "elbow" plot in the referenced figure.

### AGGLOMERATIVE CLUSTERING

Parameter	Description	Value	Dataset(s)
Number of clusters	This is decided based on elbow plots with the SC metric <a href="#">2.2.2</a>	n_clusters=2	Heart-dataset (2 & 3 dimensions) Seeds-dataset (7 dimensions) Circle-dataset (3 dimensions) Line-dataset (2 dimensions)
Number of clusters	""	n_clusters=3	Heart-dataset (9 dimensions) Seeds-dataset (2 & 3 dimensions) Skewed dataset (2 dimensions)
Number of clusters	""	n_clusters=4	Skewed dataset (2 dimensions)
Number of clusters	""	n_clusters=5	Circle dataset (2 dimensions) Line-dataset (3 dimensions)

Table 5.2: Agglomerative clustering hyperparameters for the datasets

The above table shows the relevant hyperparameters for the [Agglomerative Clustering \(AG\)](#) clustering algorithm per dataset. Because the hyperparameters "n\_clusters" rely heavily on the number of dimensions, each dataset's dimensions were provided. For the "n\_clusters," we used the [SC](#) in corporation with the elbow method to determine the number of clusters. In addition, there is the "linkage," which is "ward" for all configurations (See Section [2.2.2](#)).

## OPTICS

Parameter	Description	Value	Dataset
Minimum points	Decided using the formula $minPts = n * 2$ , where $n$ is the number of features (2.2.2)	$minPts = 4$	Seeds dataset Heart dataset Circle dataset Line dataset Skewed dataset (2-dimensions)
Minimum points	""	$minPts = 6$	Seeds dataset Heart dataset Circle dataset Line dataset Skewed dataset (3-dimensions)
Minimum points	""	$minPts = 14$	Seeds dataset (7-dimensions)
Minimum points	""	$minPts = 20$	Heart dataset (10-dimensions)

Table 5.3: OPTICS hyperparameters for datasets 1 - 3

The above table shows the relevant hyperparameters for the **OPTICS** clustering algorithm for each dataset. Because the minimum points depend on the number of dimensions, we must display them for each dataset (See Section: 2.2.2). We omitted the synthetic datasets for the same reason specified for **AP**. The first two rows display the minimum points for the datasets that have 2 and 3-dimensional data. The last two rows display the minimum points for the datasets that have 7 and 10-dimensional data.

## 5.2.2. UTILITY AND PRIVACY EVALUATION

### UTILITY

To measure cluster utility, both internal and external validation methods are used:

1. **External validation:** The external validation is measured by comparing the labels of non-private trained clustering algorithms with those trained using a privacy mechanism. The outcome is between 0 and 1, where 1 indicates the highest similarity (thus the best result). AMI (Adjusted Mutual Information) is used to assess the external validity of the clustering algorithms (See Section 2.2.3).
2. **Internal validation:** The internal validation measures the intrinsic properties of the clustering algorithms. The outcome is a value between -1 and 1 for the SC evaluation. Where -1 indicates incorrect and 1 dense clustering (See Section 2.2.3).

Because the clustering algorithms rely on Euclidean distance, we need to apply some data standardization. For this purpose, we use standard scaling provided by the Scikit-learn package<sup>6</sup>.

### PRIVACY

Privacy is hard to quantify, but we can measure the privacy loss/gain by calculating the Euclidean distance between the non-perturbed and perturbed data. In addition, we evaluate privacy by simulating a membership inference attack and calculating the adversary advantage.

1. **Privacy distance:** The first measurement we evaluate is the distance between the dataset's non-private and private variants. This metric gives us a sense of how much extra distance (and privacy) the privacy mechanisms offer compared to the non-private variant. To this end, the average Euclidean distance is measured and reported per epsilon. A higher distance indicates more privacy.
2. **Adversarial advantage:** We used the Membership Inference Attack (MIA) that Shokri et al. proposed with the implementation that was provided by Adversarial Robustness Toolkit (ART) [Nicolae et al., 2019]. An earlier study also explored this attack to evaluate differential privacy. Similar to this attack, we train a classifier with perturbed data and evaluate it using non-perturbed data (test-data / shadow-data) [Zhao et al., 2020]. We consider a semi-supervised setup (Figure 4.2), where we train a classifier with the perturbed data and evaluate it using non-perturbed data (test-data / shadow-data). For the classifier, we use a RandomForestClassifier [Rigaki and Garcia, 2021]. Finally, we evaluate the adversary advantage (percentage) using  $TPR - FPR$  [Yeom et al., 2018]. Figure: 5.3 provides a visual setup of the experiment.

Although the adversary advantage is a proven method for measuring MI attacks, we also include the TPR as an essential measurement. This metric is most interesting to us because this displays the percentage of the actual leaked labels (see Section 4.3).

---

<sup>6</sup><https://scikit-learn.org/stable/modules/preprocessing.html>

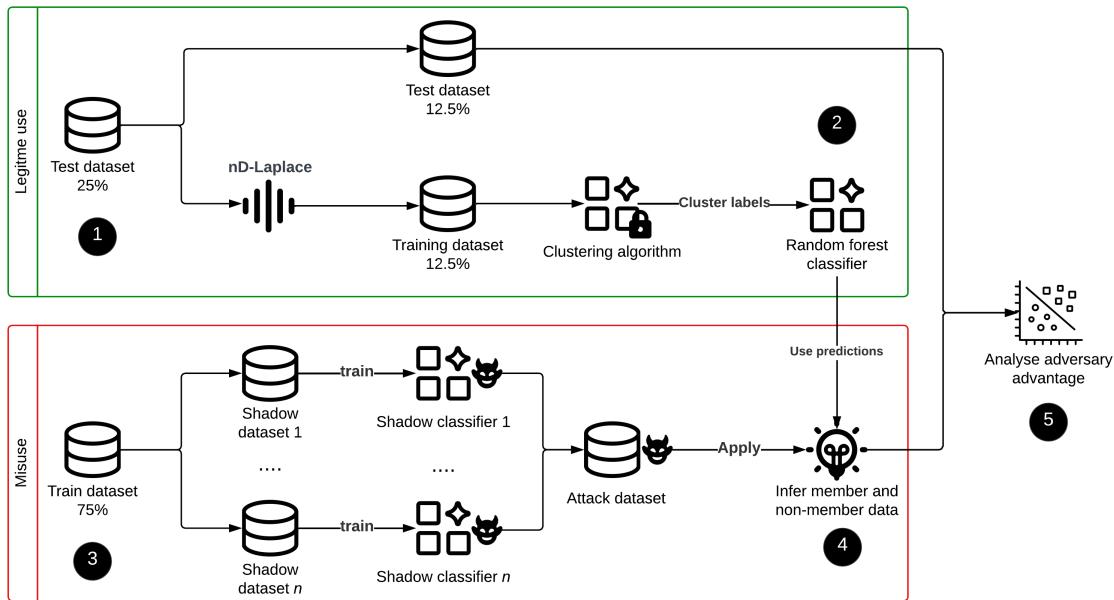


Figure 5.3: Member inference attack using shadow models. The green swim lane illustrates the normal setup, and the red swim lane projects the adversary steps.

The figure above visually represents the **MIA** attack we set up using semi-supervised learning.

- The green swim line presents the setup for legitimate system usage. Here, we use a general setup for a supervised machine learning setup. First, the data is split into test/train data, where only the training data is private.
- Then, the K-Means clustering algorithm is privately trained, and the labels are used with the RandomForestClassifier.
- The red swim line presents the adversary's steps. This approach uses Shokri et al. attack setup using shadow models [Shokri et al., 2017]. The random shadow datasets and classifiers are combined into an attack dataset.
- Then, the attacker tries to infer the training data from the attack dataset using the classifier's predictions.
- Finally, we calculate the adversary advantage using the inferred labels and the original labels.

The attacker could infer many original labels if the adversary advantage or **TPR** is high. In this case, a lower score is better than a higher score.

### 5.2.3. RESEARCH ROADMAP

This section explains the experiments' setup and the result reporting order.

- Cluster utility: The mechanisms kD-Laplace and Piecewise are compared using the external and internal validation methods. We report the **AMI** and **SC** in the results for both mechanisms and each dataset. For this experiment, we only consider real-world datasets.

2. Mechanism utility: Both mechanisms are compared against each other by evaluating the **AMI** and **SC**. In this research, we also conduct experiments to assess the impact of the number of dimensions on the utility of each dataset.
3. Mechanism privacy: We do the same for utility, but now with privacy. For this purpose, we evaluate the adversary advantage and privacy distance.
4. Mechanism comparison: Now we established the performance of kD-Laplace and Piecewise, we compare the different variants of kD-Laplace.

### 5.3. DATA ANALYSIS

For the data analysis, we utilize visualization libraries in Python, primarily Matplotlib, in combination with Seaborn.

- Line/bar plot: These plots visualize comparisons between the different privacy mechanisms. The different epsilon values are displayed on the x-axis, and the utility/privacy metric is displayed on the y-axis.
- Heatmap plot: These plots visualize the interaction between dimensions and epsilons. This way, we can display two categorical values (dimension/epsilon) and one continuous value.

# 6

## RESULTS

This chapter aims to present the results to the reader. As indicated in the methodology, we tested four mechanisms (three variants of kd-Laplace and Piecewise) for their utility and privacy.

1. Cluster utility: Visualise the utility of kD-Laplace and Piecewise for 2/3/n-dimensional data. The results are displayed using a line diagram, using **AMI** and **SC**.
2. Mechanism utility: Visualise the utility of kD-Laplace and Piecewise using a heatmap plot. The results for each dimension and privacy budget are provided for only the **AMI** score.
3. Mechanism privacy: Visualise the privacy of kD-Laplace and Piecewise using a heatmap plot. The results for each dimension and privacy budget are provided for only the adversary advantage and privacy distance.
4. Mechanism comparison: Visualise the comparison between the variants of kD-Laplace and Piecewise. The results for each mechanism & variant are provided using a bar-plot for the **AMI** score (utility) and the adversary advantage (privacy).

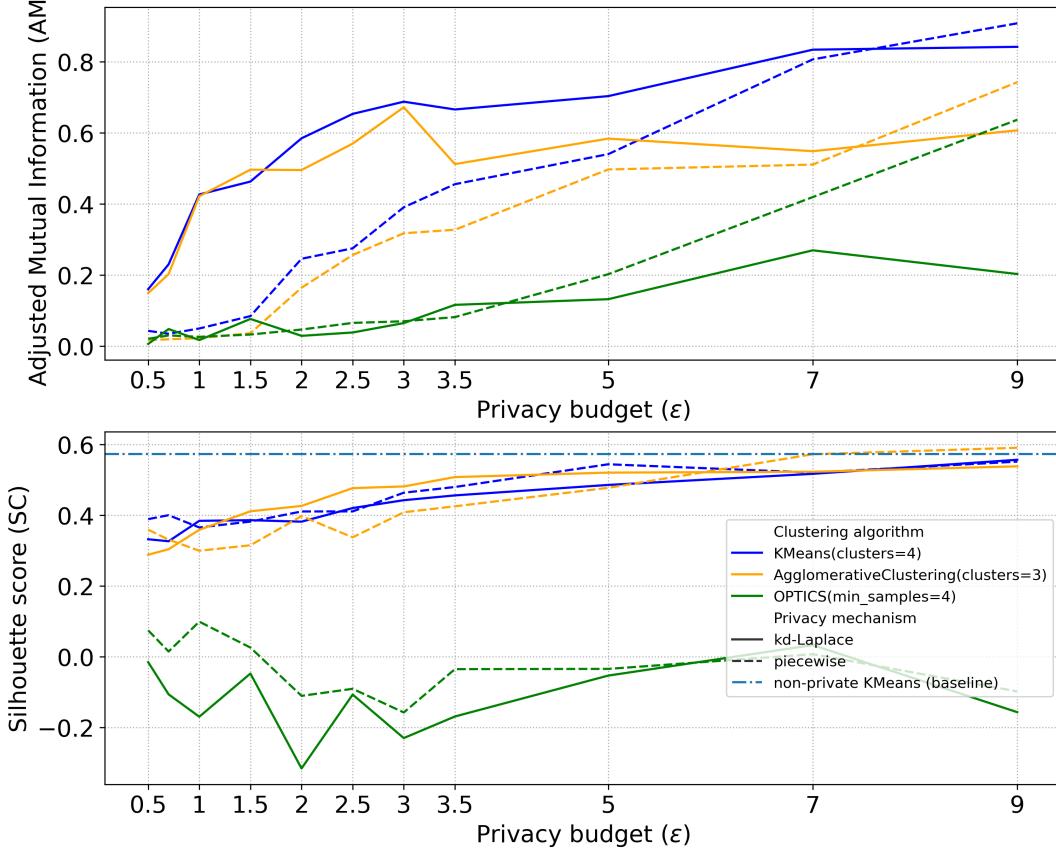
All results are reported for each dataset (See methodology: [5.1.3](#)) separately.

### 6.1. CLUSTER UTILITY

The results below display the difference in external and internal utility for the three clustering algorithms using the kd-Laplace and Piecewise mechanisms. The x-axis shows the privacy budget, and the y-axis shows the **AMI**. Please refer to the plots in the appendix for the same scores for kD-Laplace variants (See [.4](#)).

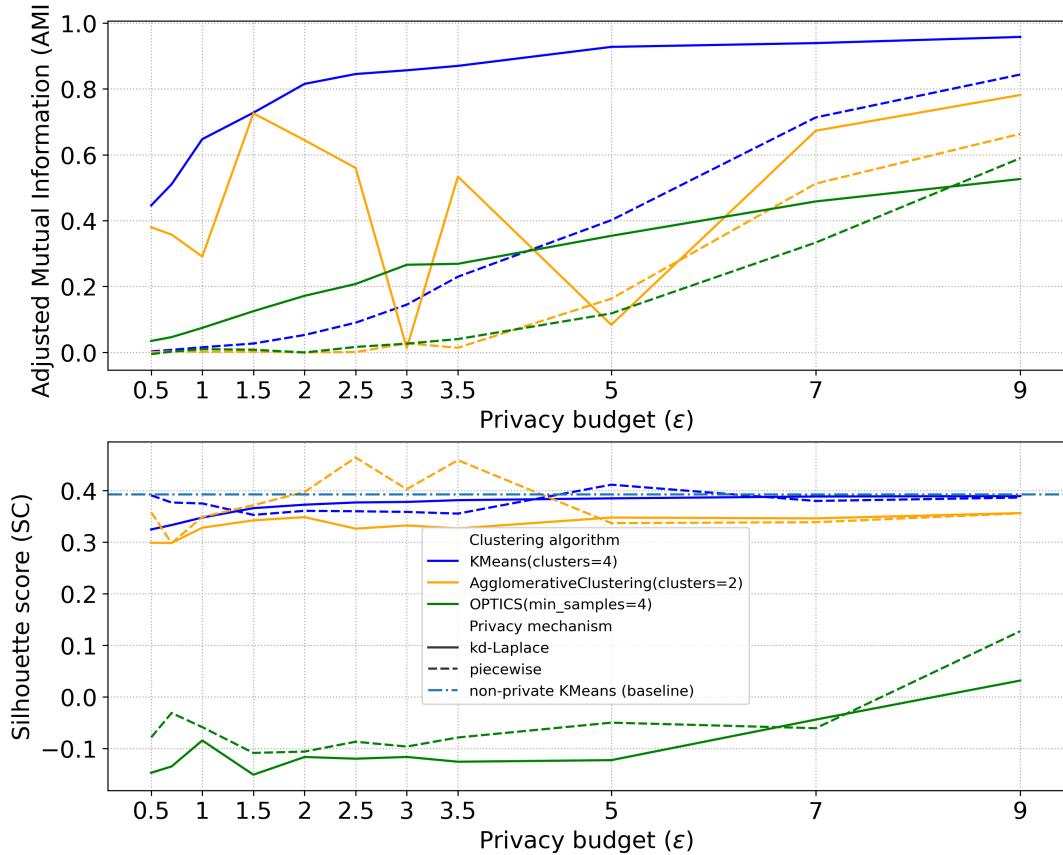
### 6.1.1. 2-DIMENSIONAL DATA

Figure 6.1: AMI (top) and SC (bottom) for the kD-Laplace and Piecewise mechanisms for the 2-dimensional data seeds-dataset



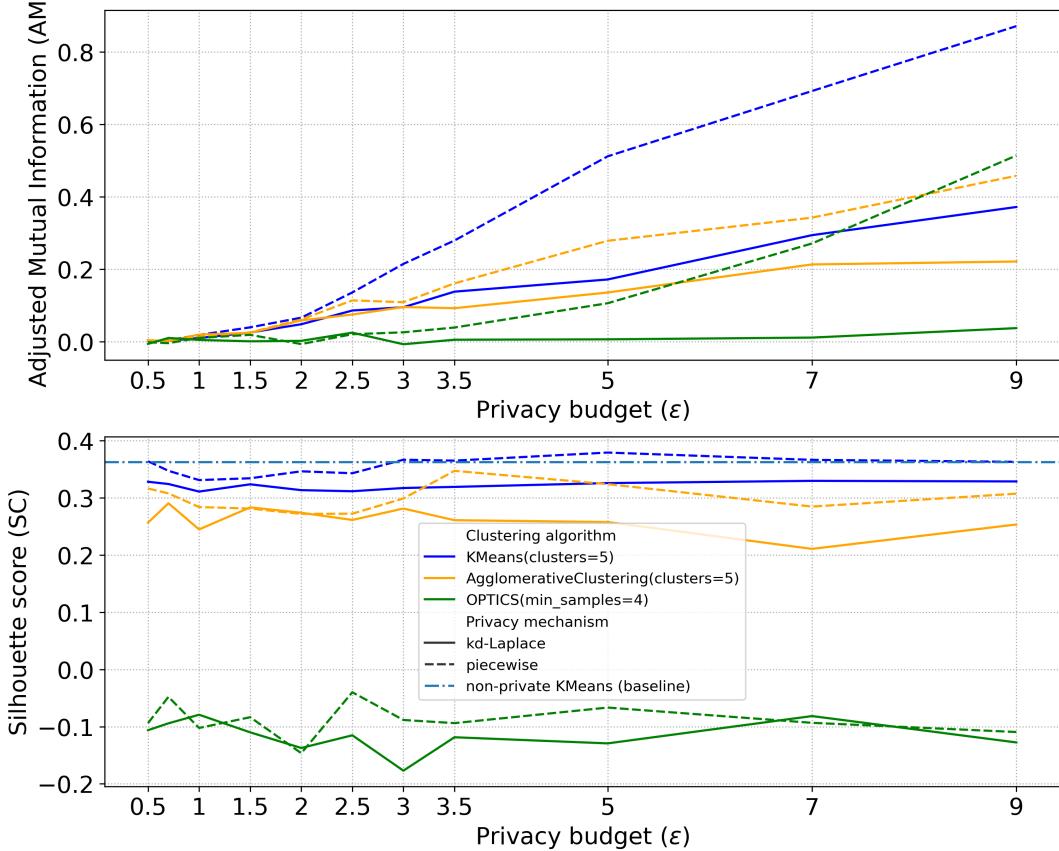
Piecewise shows a higher **AMI** for privacy budgets (epsilon) 7 and 9. In contrast, the kd-Laplace mechanism scores higher for all the other epsilons (0.1 onward 7). K-Means achieves the best score for both mechanisms with a slight margin compared to AG, but they score equal for **SC**. For both the mechanisms, OPTICS underperforms heavily but still scores the same as AP for Piecewise.

Figure 6.2: AMI (top) and SC (bottom) for the kD-Laplace mechanism for the 2-dimensional data heart-dataset



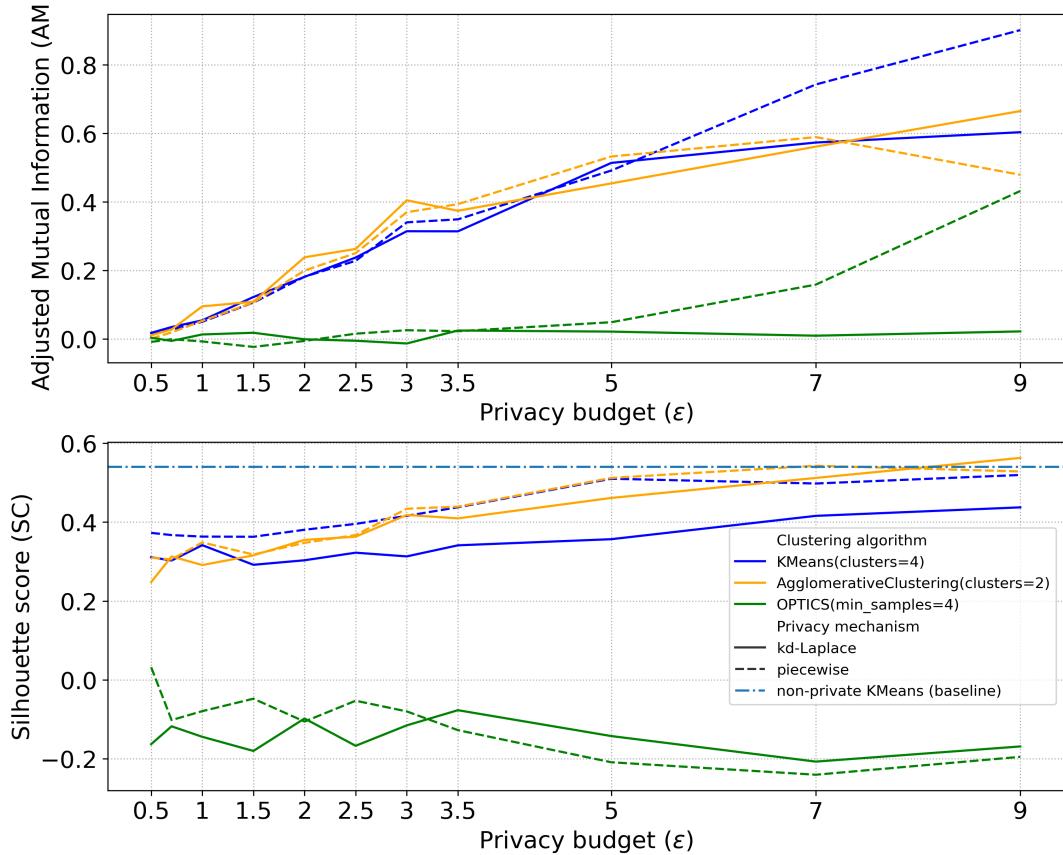
K-Means achieves a high **AMI** score of 0.95+ using kD-Laplace from privacy budgets 5 to 9, compared to Piecewise scoring between 0.4 to 0.8. Both mechanisms score high for **AG** and K-Means, reaching the baseline value, but **OPTICS** performs poorly. **AG** fluctuates more with kD-Laplace. For **SC**, both kD-Laplace and Piecewise reach the baseline for **AG** and K-Means, while **OPTICS** performs poorly.

Figure 6.3: **AMI** (top) and **SC** (bottom) for the kD-Laplace mechanism for the 2-dimensional data circle-dataset



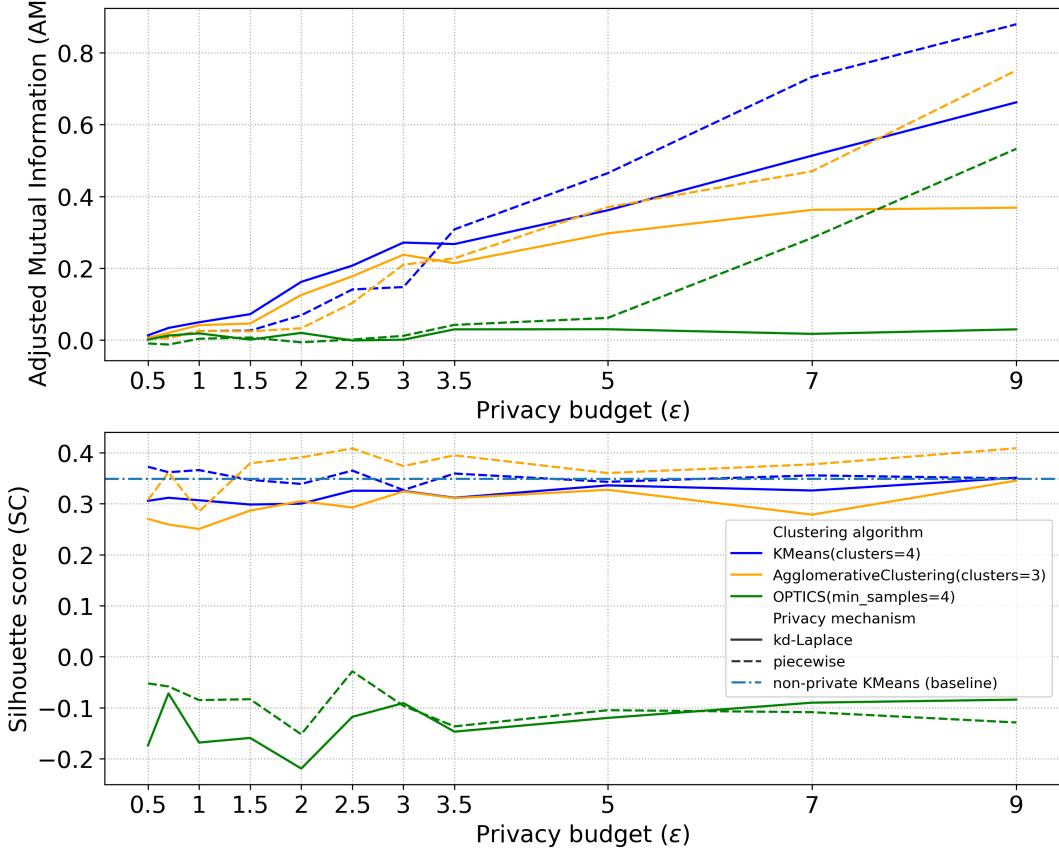
For both mechanisms, K-Means scores the privacy budgets from 2.5. Compared to Piecewise, the kD-Laplace mechanism scores the worst overall for almost all privacy budgets. The Piecewise algorithm scores 0.84 as the highest **AMI** for K-Means, while kd-Laplace scores 0.36. The difference is a little less for **SC**, but overall Piecewise is 0.10 higher. K-Means is for Piecewise on the baseline value for privacy budgets 3 to 9.

Figure 6.4: AMI (top) and SC (bottom) for the kD-Laplace and Piecewise mechanisms for the 2-dimensional data line-dataset



The kD-Laplace algorithm performs best with a privacy budget of 5 for K-Means (AMI 0.55 - 0.6). AG performs slightly better with a budget of 9 (AMI 0.63). Piecewise outperforms other methods, scoring 0.83 AMI for K-Means at budget 9 and better for different budgets, except for AG. For SC, the mechanisms have similar scores, with Piecewise at baseline from budget 5, kD-Laplace at 7, and AG beyond baseline from budget 9.

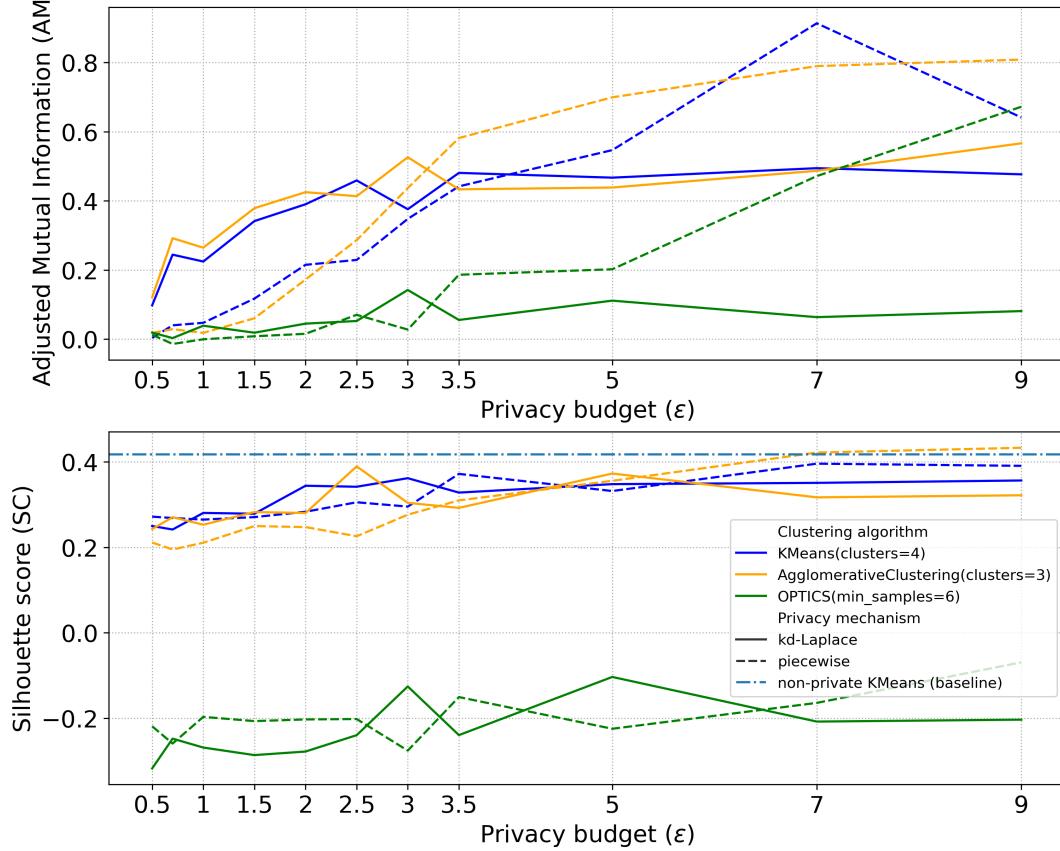
Figure 6.5: **AMI (top)** and **SC (bottom)** for the kD-Laplace and Piecewise mechanisms for the 2-dimensional data skewed-dataset



kD-Laplace performs best with K-Means, reaching 0.63 **AMI** at privacy budget 9. **AG** follows, scoring 0.38 **AMI** at budgets 7 and 9. Piecewise outperforms kD-Laplace for all clustering algorithms, scoring 0.82 **AMI** for K-Means at budget 9. In terms of **SC**, Piecewise also outperforms both K-Means and **AG**, scoring above the baseline (0.4), while kD-Laplace is slightly worse (around 0.3).

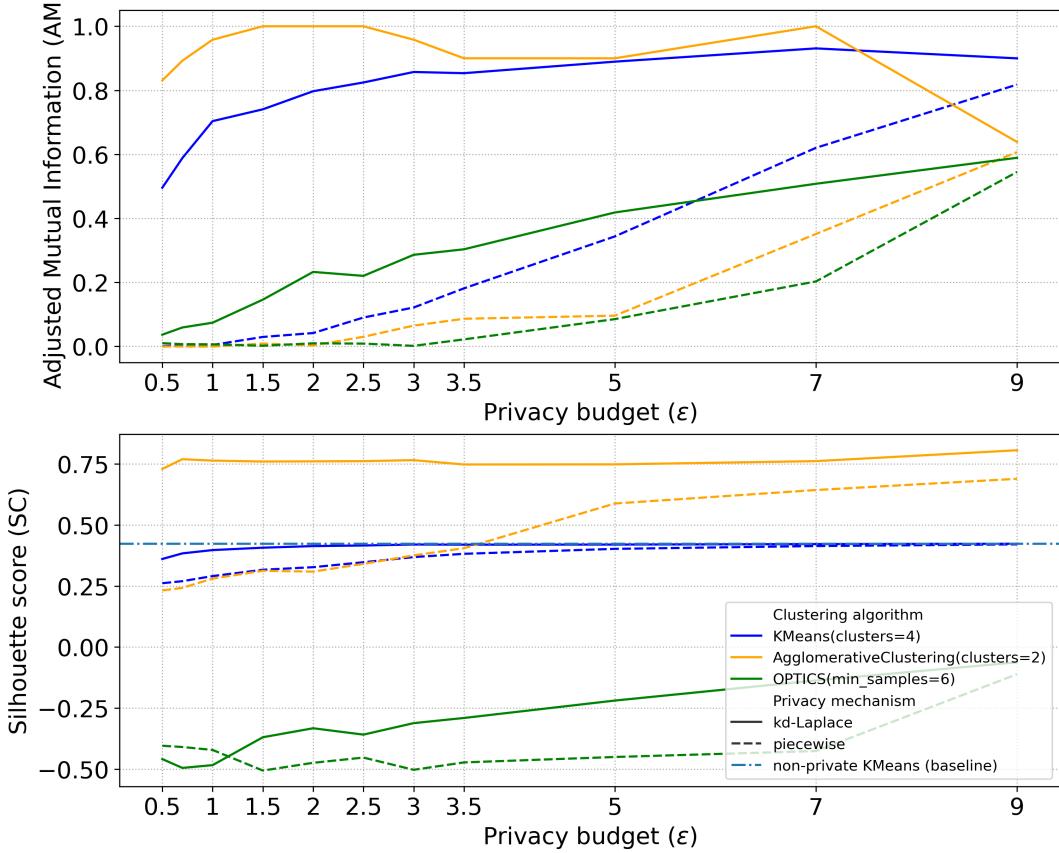
### 6.1.2. 3-DIMENSIONAL DATA

Figure 6.6: AMI (top) and SC (bottom) for the kD-Laplace and Piecewise mechanisms for the 3-dimensional data seeds-dataset



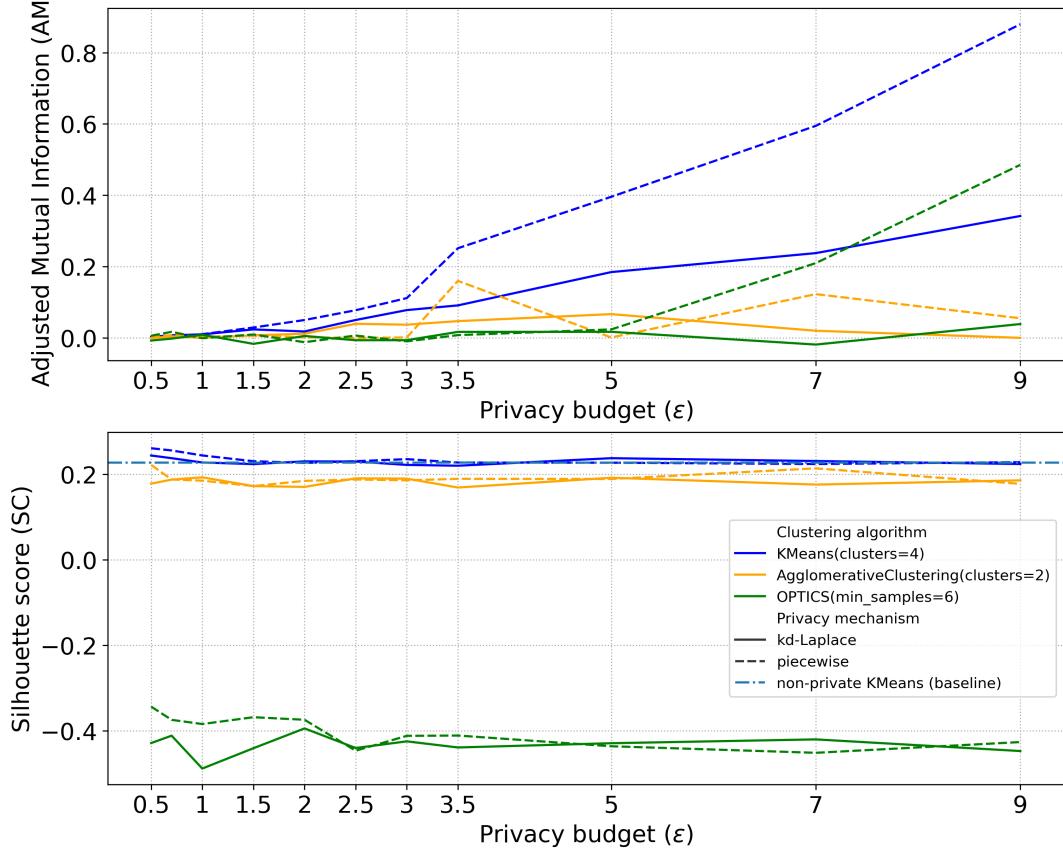
Piecewise performs better for the **AMI** for privacy budgets 3.5 to 9 but worse for the lower epsilons (0.1 - 3.5). K-Means scores 0.5 from epsilon 2 to 9 for kD-Laplace. For both mechanisms, **AG** scores the same as K-Means for both **AMI** and **SC**. **OPTICS** scores low (< 0.1) for kD-Laplace, but for Piecewise, it scores 0.4 and 0.6 for epsilon 7 and 9, respectively.

Figure 6.7: **AMI (top) and SC (bottom) for the kD-Laplace and Piecewise mechanisms for the 3-dimensional heart-dataset**



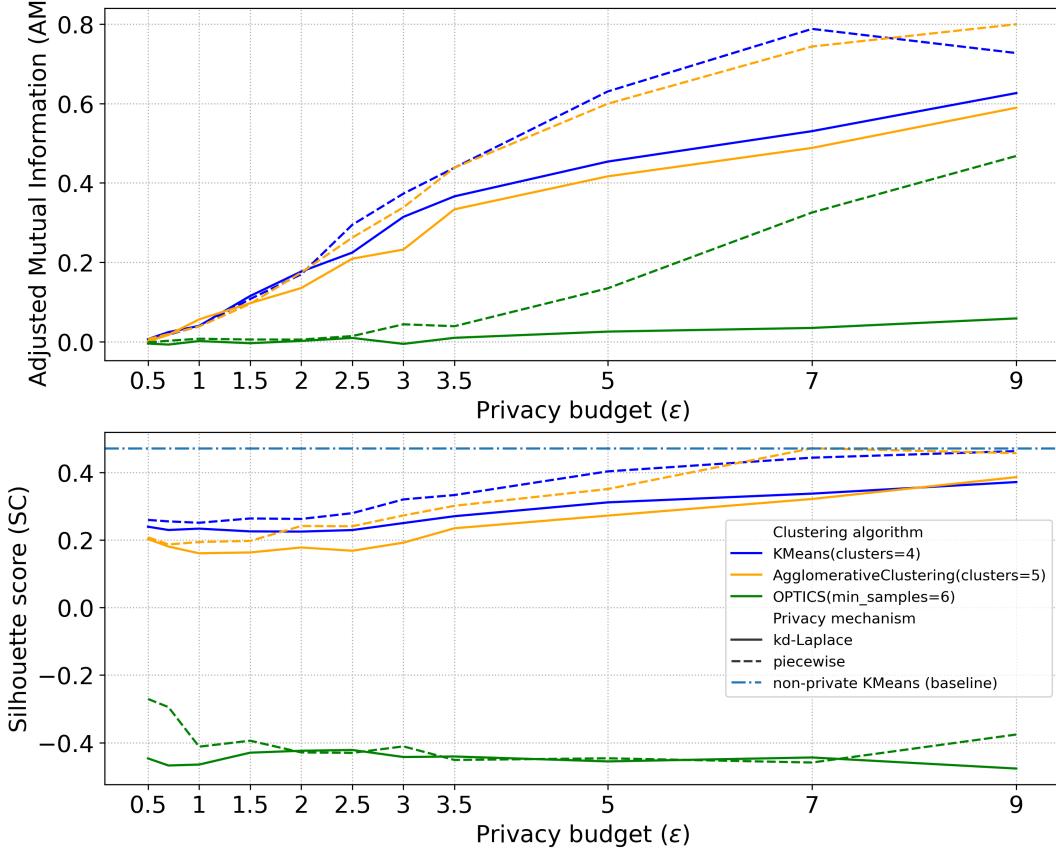
KD-Laplace performs well (+/- 0.90 AMI) for K-Means at budgets 7 and 9. AG scores around +/- 1.0 AMI for epsilon values of 1.5 to 2.5, staying above 0.8 for most budgets. OPTICS scores better (0.6 AMI at budget 9). Piecewise scores below 0.6 AMI for budgets < 7, reaching a maximum of 0.8 AMI for K-Means. For SC, both Piecewise and kD-Laplace reach the baseline for AG and K-Means. KD-Laplace surpasses the baseline by 0.25 for AG.

Figure 6.8: AMI (top) and SC (bottom) for the kD-Laplace and Piecewise mechanisms for the 3-dimensional data circle-dataset



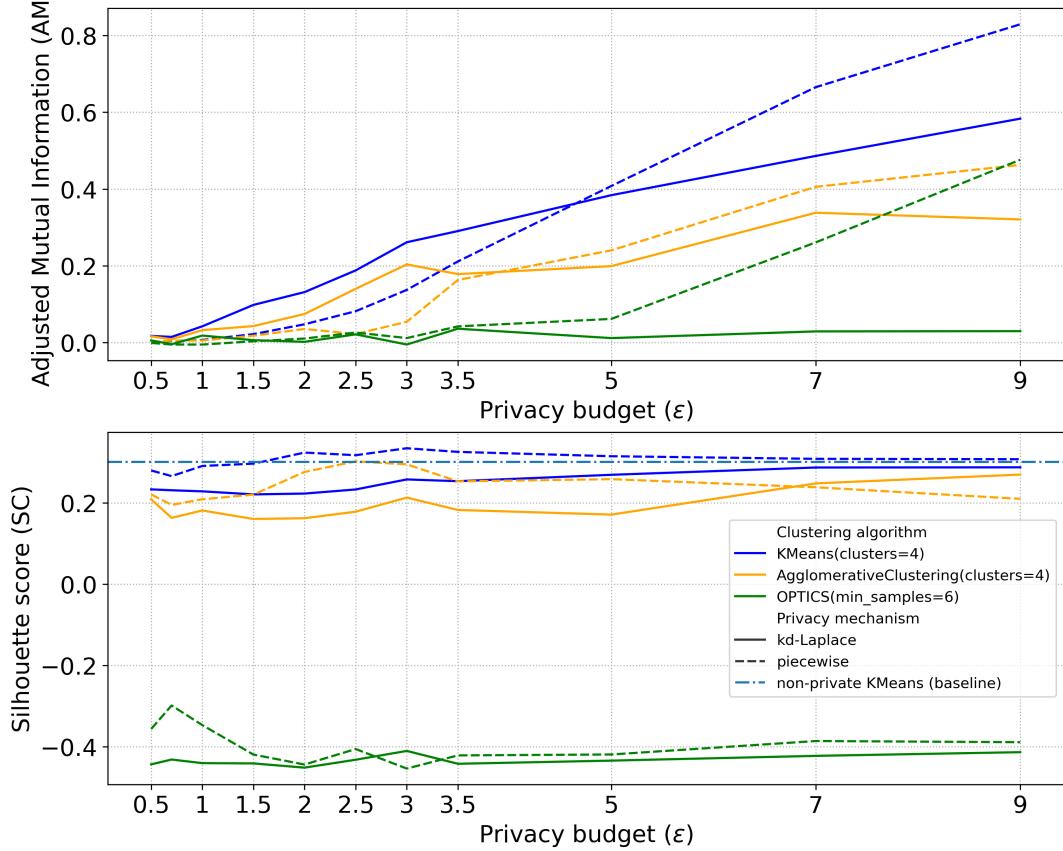
kD-Laplace performs best with K-Means, scoring a maximum of 0.33 AMI. The other algorithms perform poorly (< 0.1 AMI). Piecewise also shows the best results with K-Means, scoring 0.82 AMI, while the other algorithms perform weakly (< 0.2 AMI). OPTICS performs relatively well with 0.49 AMI at privacy budget 9. Both mechanisms have similar SC scores (around 0.2), which is at the baseline level. OPTICS performs below 0 SC for all privacy budgets.

Figure 6.9: AMI (top) and SC (bottom) for the kD-Laplace and Piecewise mechanisms for the 3-dimensional data line-dataset



K-Means performs best with kD-Laplace, scoring 0.6 AMI from privacy budget 3.5 to 9. AG follows the same pattern but with around 0.2 AMI lower. OPTICS scores poorly (< 0.1 AMI). Piecewise shows a similar trend as kD-Laplace, with slightly higher scores. OPTICS improves to 0.43 AMI at privacy budget 9. Both mechanisms are similar for SC, with Piecewise slightly higher, while OPTICS remains relatively lower.

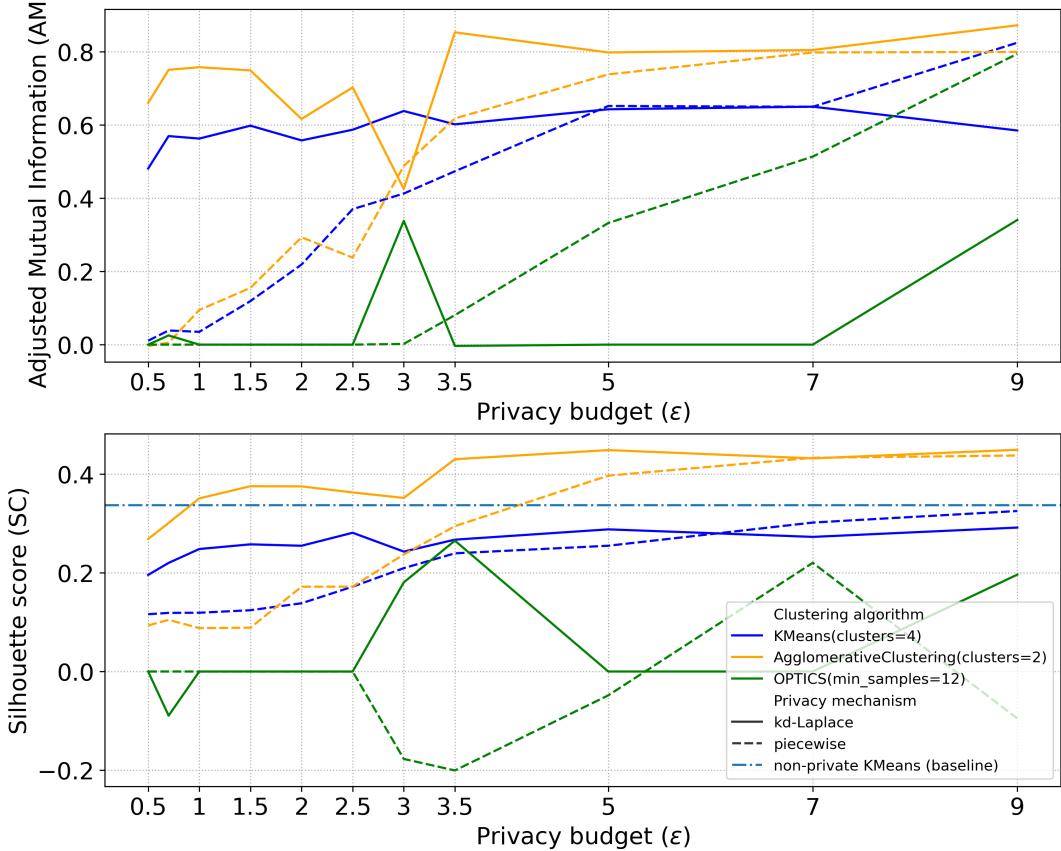
Figure 6.10: AMI (top) and SC (bottom) for the kD-Laplace and Piecewise mechanisms for the 3-dimensional data skewed-dataset



kD-Laplace scores best with the K-Means algorithm, 0.59 AMI for privacy budget 9. In contrast, AG is +/- 0.3 for privacy budgets 7 and 9. OPTICS scores worst with a maximum value of 0.02 AMI. The Piecewise mechanism scores much better with a AMI of 0.8 for K-Means. The same trend is observed for AG and OPTICS as with kD-Laplace. Both score worse, but OPTICS scores equal to AG for privacy budget 9. Both mechanisms score equal for the SC, but Piecewise is slightly higher and scores equal to the baseline.

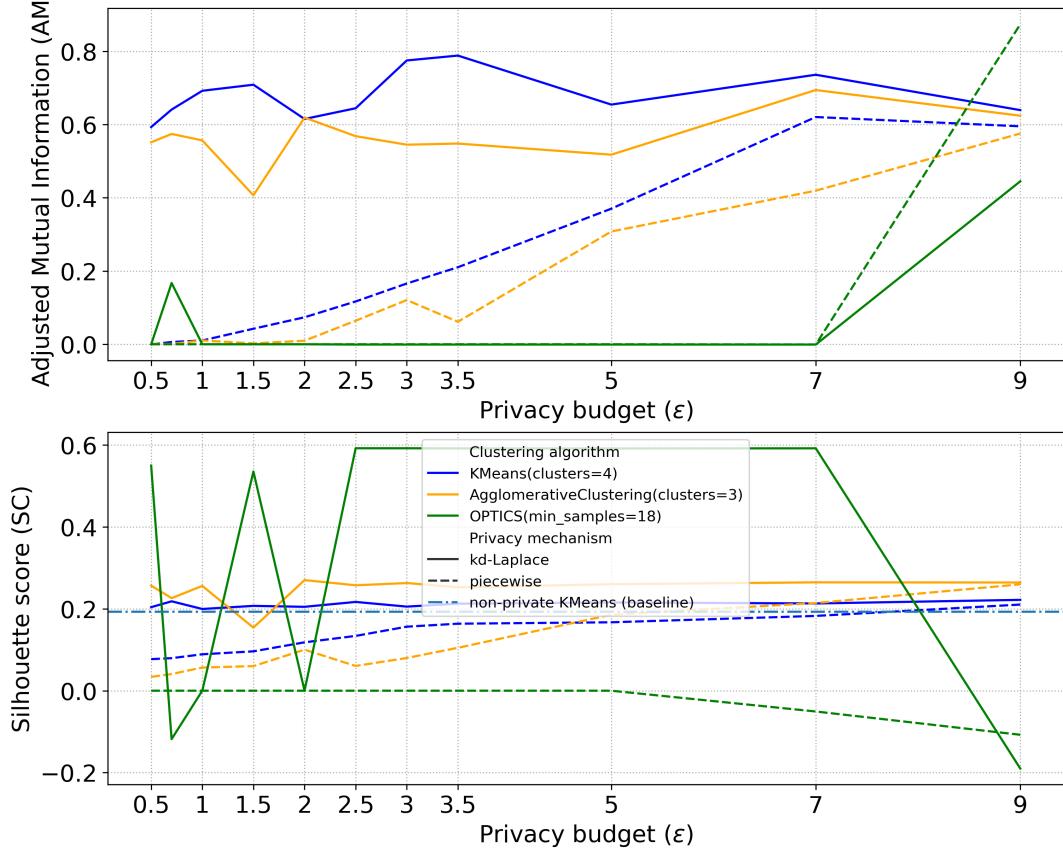
### 6.1.3. N-DIMENSIONAL DATA

Figure 6.11: AMI (top) and SC (bottom) for the kD-Laplace and Piecewise mechanisms for the n-dimensional data seeds-dataset



kD-Laplace scores 0.83 **AMI** for **AG** at privacy budget 9, while K-Means scores 0.5 to 0.6 **AMI** for all budgets. Piecewise shows a similar trend, with all algorithms scoring nearly equal at budget 9. **SC** follows a similar trend as **AMI**, with peaks at epsilon 3 and 3.5 for kD-Laplace with **OPTICS**. Overall, **OPTICS** scores lower for both mechanisms, with an upward trend for Piecewise but not for kD-Laplace.

Figure 6.12: AMI (top) and SC (bottom) for the kD-Laplace and Piecewise mechanisms for the n-dimensional data heart-dataset



kD-Laplace shows consistently high scores for K-Means and AG. Piecewise scores worse for K-Means (0.6) at budgets 7 and 9, while AG scores slightly lower. Remarkably, OPTICS achieves 0.85 AMI for privacy budgets 9 but zero for other budgets. For SC, kD-Laplace performs well for K-Means and AG, while OPTICS fluctuates from -0.1 to 0.6. Similarly, for Piecewise, OPTICS scores zero for most budgets but 0.1 for budgets 7 and 9.

## **6.2. MECHANISM UTILITY**

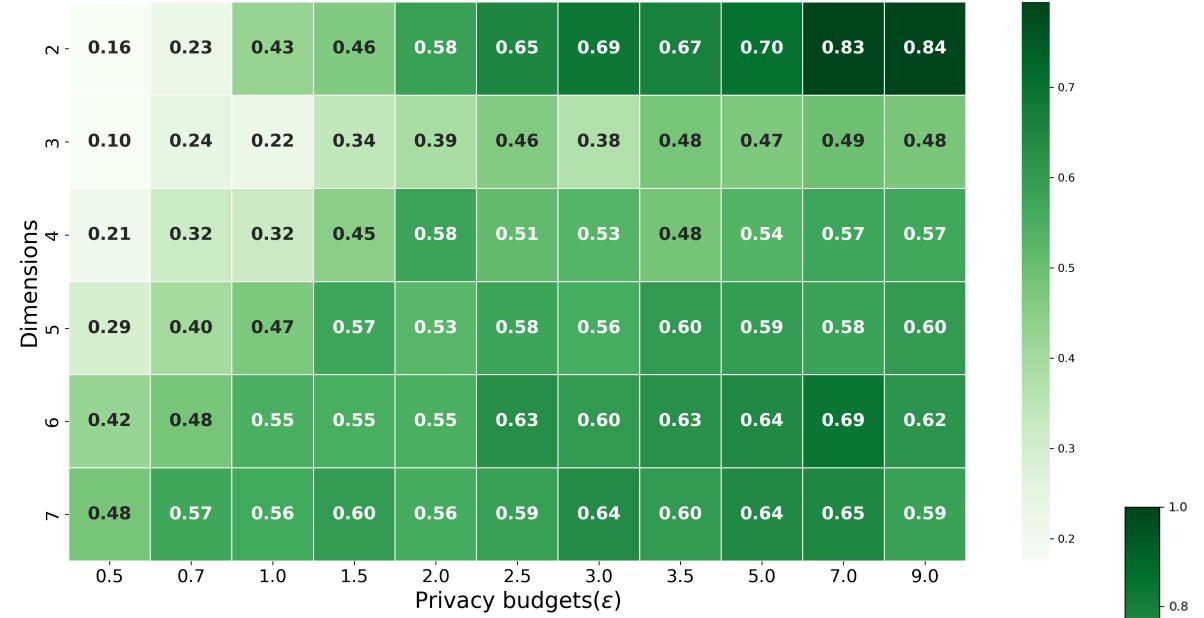
The kd-Laplace and Piecewise mechanisms are compared using a heatmap in the sections below. We used a heatmap to show both the privacy budget (epsilon) and the number of dimensions. As explained earlier in the research design, only the K-Means algorithm was used to calculate the AMI. The x-axis displays the privacy budget, and the y-axis shows the number of dimensions.

Please refer to the plots in the appendix for:

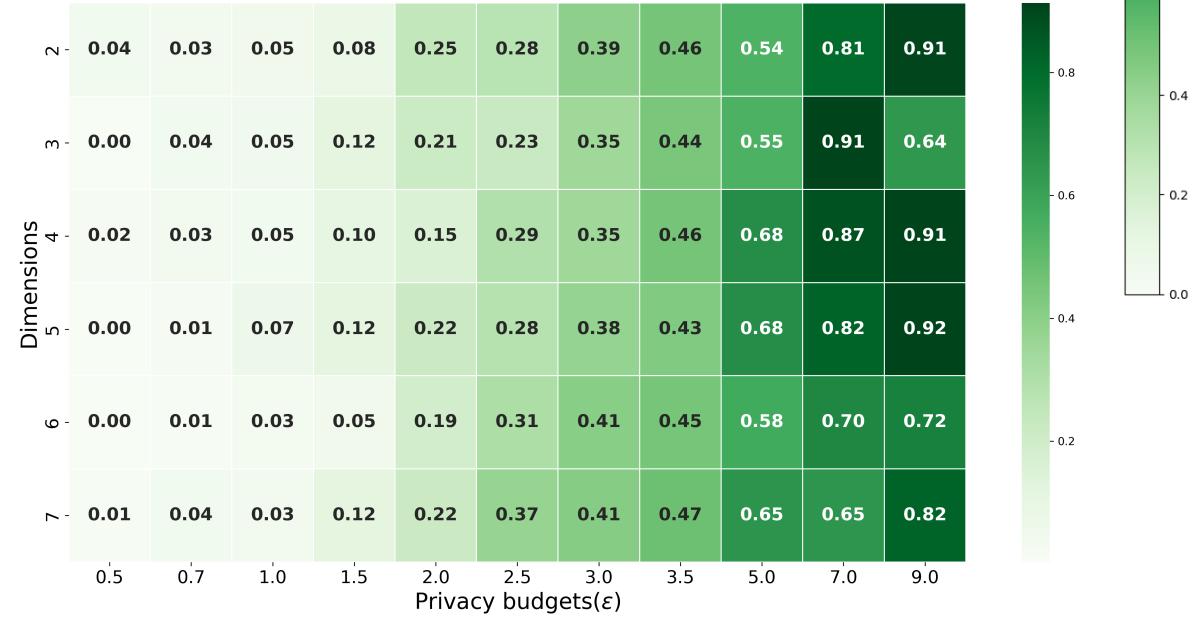
1. Privacy distance: [.5](#).
2. True Positive Rate: [.5](#).

### 6.2.1. SEEDS-DATASET

(a) Adjusted Mutual Information comparison for the kd-Laplace mechanism



(b) Adjusted Mutual Information comparison for the Piecewise mechanism



The two plots compare external validation (AMI) between the two privacy mechanisms. The kd-Laplace mechanism scores better for all privacy budgets (epsilon) for K-Means, except for privacy budgets 7 and 9. For epsilon 5, Piecewise is slightly better, but kD-Laplace scores better for 2 - 5 dimensions. KD-Laplace scores low for 2 - 5 dimensions and epsilon 0.5 and 0.7. However, Piecewise scores < 0.10 AMI for epsilons < 3 for all dimensions.

## 6.2.2. HEART-DATASET

(a) Adjusted Mutual Information comparison for the kd-Laplace mechanism



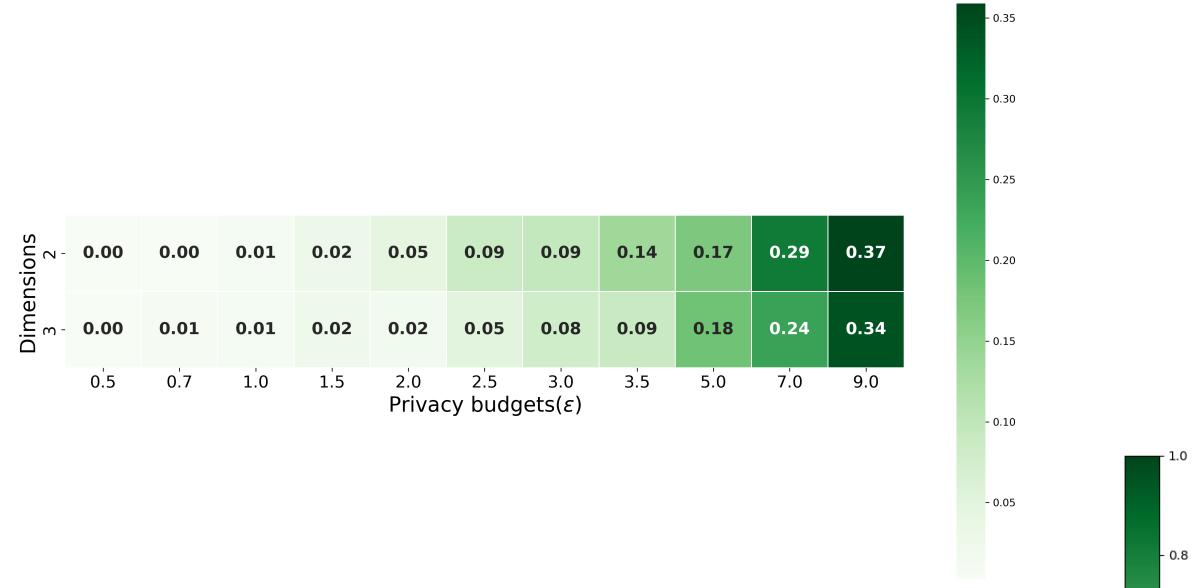
(b) Adjusted Mutual Information comparison for the Piecewise mechanism



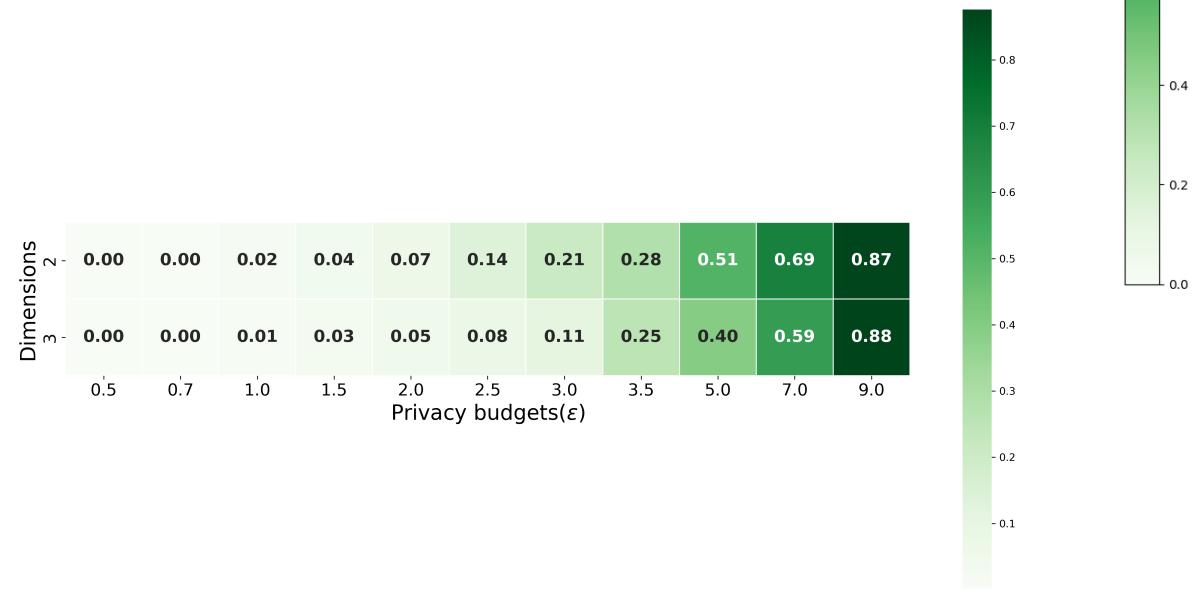
For kD-Laplace, all scores are higher than 0.23 AMI. Dimensions from 4 to 10 with a privacy budget of 0.5 have the lowest scores, but then the AMI increases. Piecewise scores a maximum of 0.24 AMI for privacy budgets 0.5 to 3.5, with most values below 0.10. The number of dimensions does not seem to have an influence here. The values for privacy budgets 7 and 9 are approximately 0.70 and 0.80 AMI, respectively.

### 6.2.3. CIRCLE-DATASET

(a) Adjusted Mutual Information comparison for the kd-Laplace mechanism



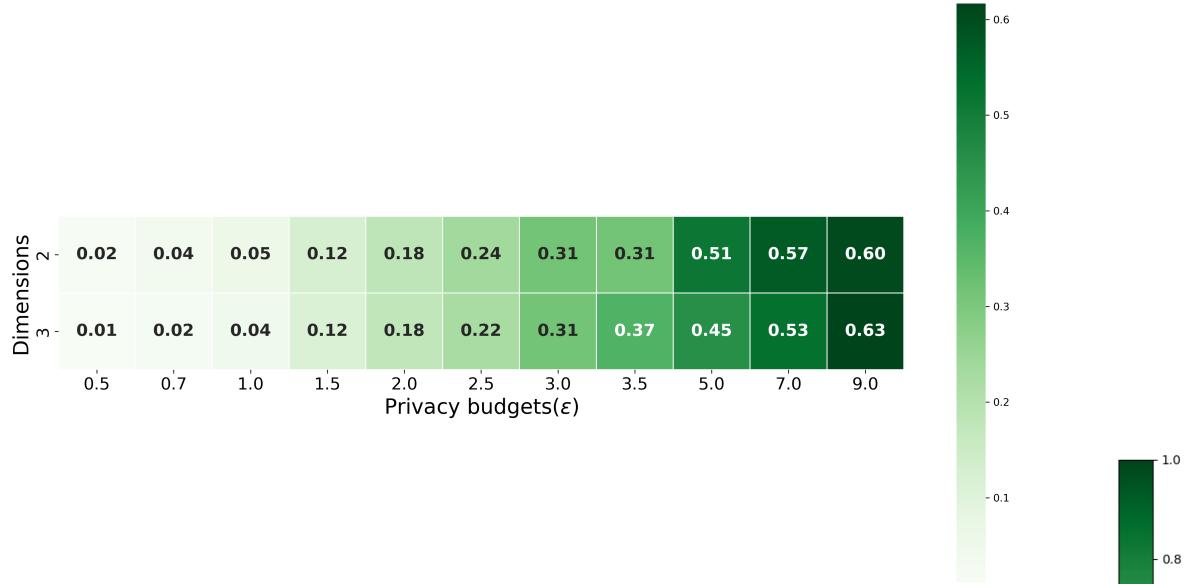
(b) Adjusted Mutual Information comparison for the Piecewise mechanism



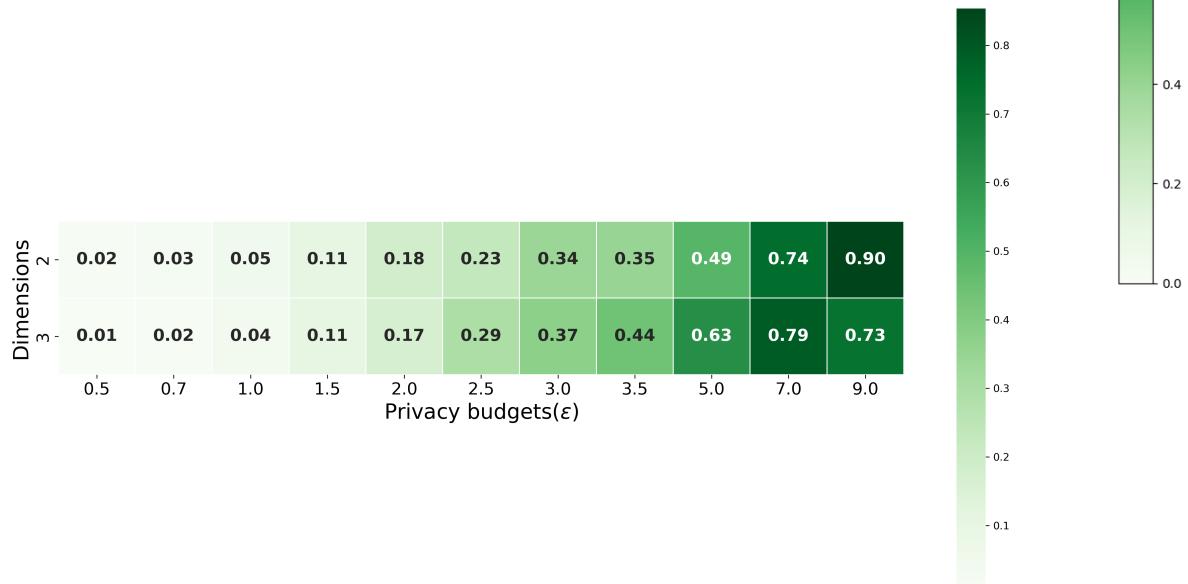
kD-Laplace achieves the highest value of 0.37 AMI for privacy budget 9. However, the score remains low for the other privacy budgets, regardless of the dimension. For Piecewise, the highest score is significantly higher, reaching a maximum of 0.88 AMI for privacy budget 9. The dimension does not seem to have much impact here, either.

## 6.2.4. LINE-DATASET

(a) Adjusted Mutual Information comparison for the kd-Laplace mechanism



(b) Adjusted Mutual Information comparison for the Piecewise mechanism



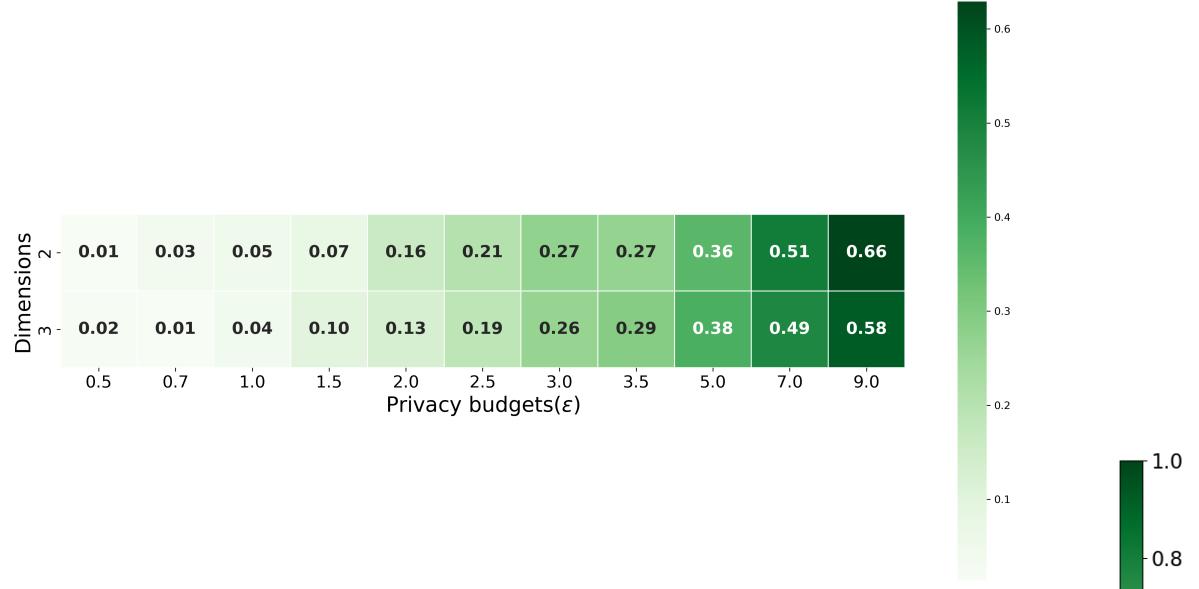
Up to privacy budget 3.5, the scores remain relatively similar for both mechanisms, but after that, they become better for Piecewise.

For kd-Laplace, the highest score is 0.63 AMI for privacy budget 9. For privacy budgets 5 and 7, the score is around 0.5 - 0.55. After that, the scores are no higher than 0.37.

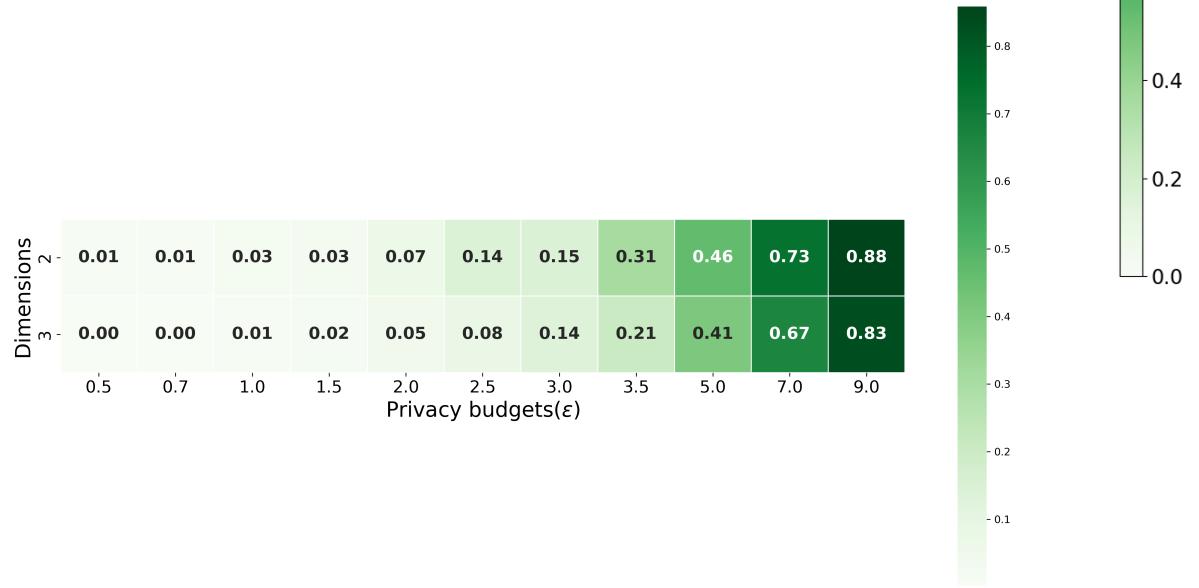
Piecewise performs better, with a maximum of 0.90 AMI for privacy budget 9 for 2 dimensions. For both mechanisms, the dimensions do not significantly impact the score.

### 6.2.5. SKEWED-DATASET

(a) Adjusted Mutual Information comparison for the kd-Laplace mechanism



(b) Adjusted Mutual Information comparison for the Piecewise mechanism



The kd-Laplace mechanism performs slightly better from privacy budgets 2 to 5. However, Piecewise performs significantly better after that, achieving around 0.73 and 0.85 AMI for privacy budgets 7 and 9, respectively.

In contrast, kd-Laplace scores much lower in this range. Furthermore, it appears that the dimensions do not significantly impact the performance of the mechanisms in this context.

## 6.3. PRIVACY

### 6.3.1. SEEDS DATASET

(a) Heatmap showing adversary advantage for the kD-Laplace mechanism, per privacy budget & dimension for seeds-dataset.



(b) Heatmap showing adversary advantage for the Piecewise mechanism, per privacy budget & dimension for seeds-dataset.



The heatmaps show the adversary advantage of kD-Laplace and Piecewise mechanisms. For kD-Laplace, the advantage slightly decreases with higher privacy budgets. Dimension 4 has a relatively high advantage for budgets 0.5 and 0.7. Piecewise generally has a higher advantage, especially for 2 dimensions (around +/- 0.5). Overall, dimensions don't have a

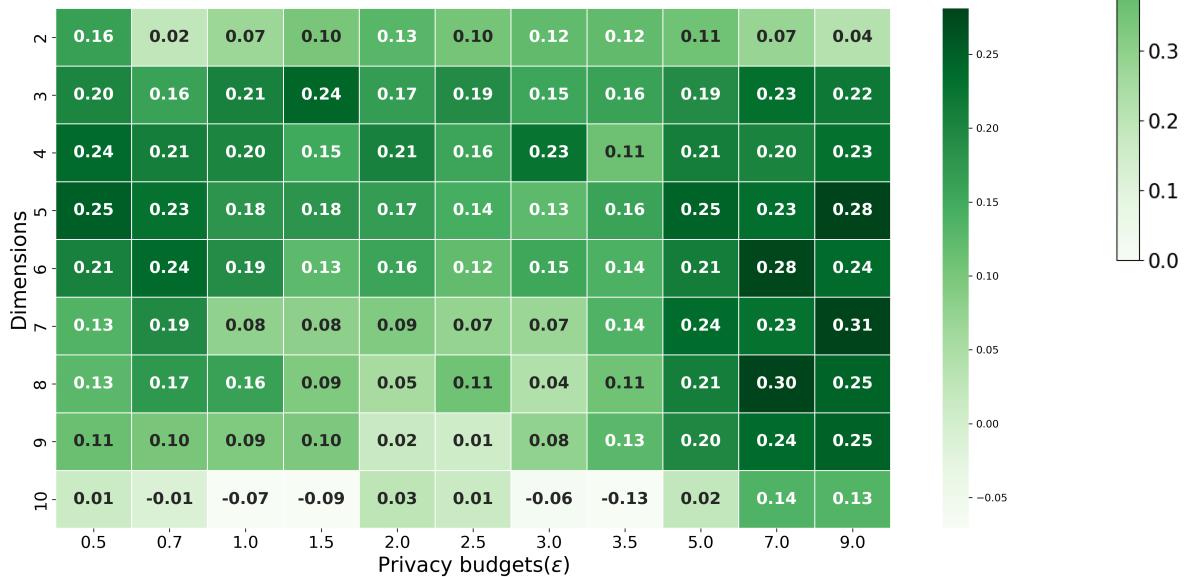
significant impact on the advantage for both mechanisms.

### 6.3.2. HEART DATASET

(a) Heatmap showing adversary advantage for the kD-Laplace mechanism, per privacy budget & dimension for heart-dataset.

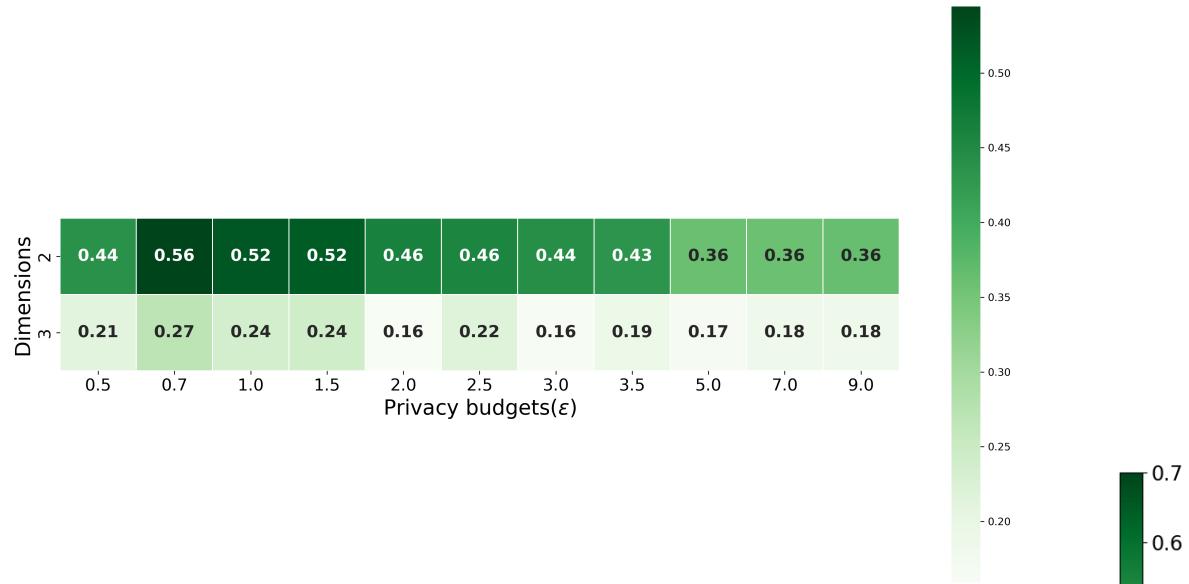


(b) Heatmap showing adversary advantage for the Piecewise mechanism, per privacy budget & dimension for heart-dataset.

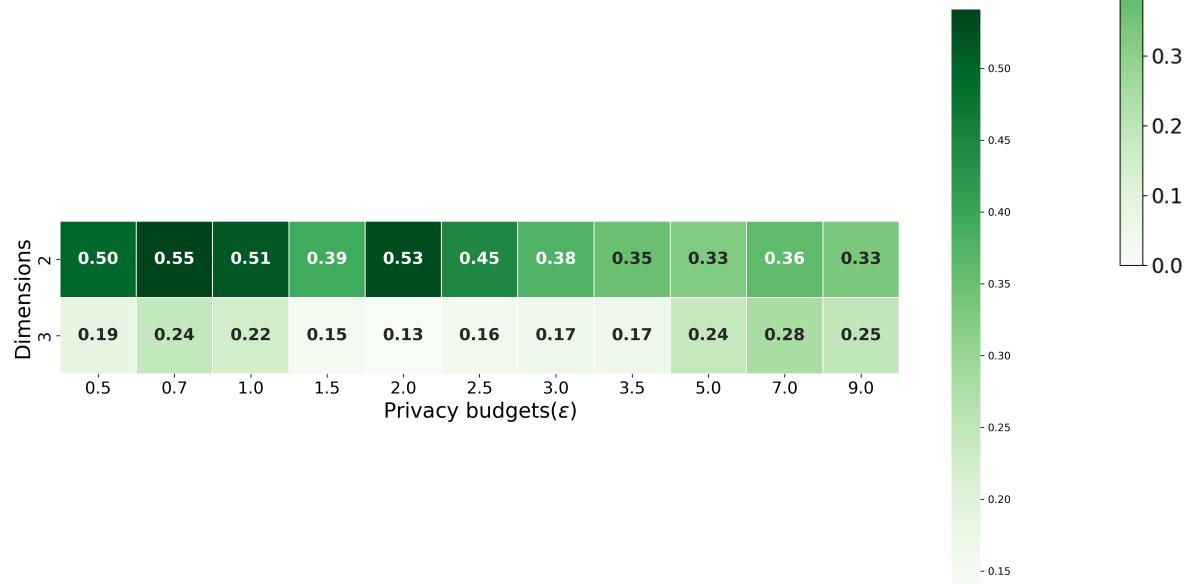


### 6.3.3. CIRCLE DATASET

(a) Heatmap showing adversary advantage for the kD-Laplace mechanism, per privacy budget & dimension for circle-dataset.

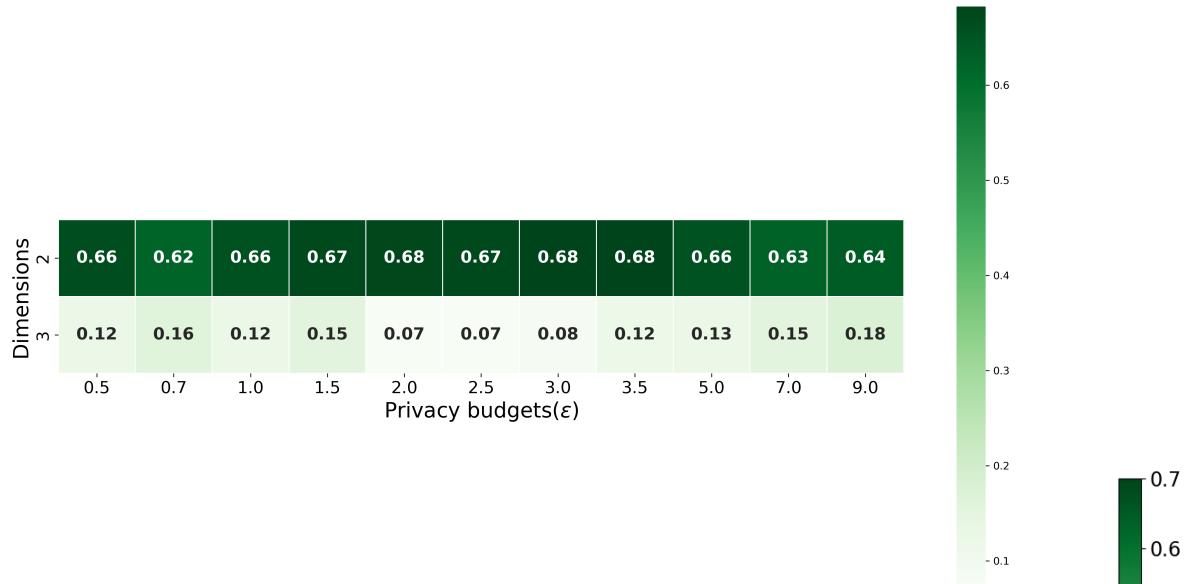


(b) Heatmap showing adversary advantage for the Piecewise mechanism, per privacy budget & dimension for seeds-dataset.

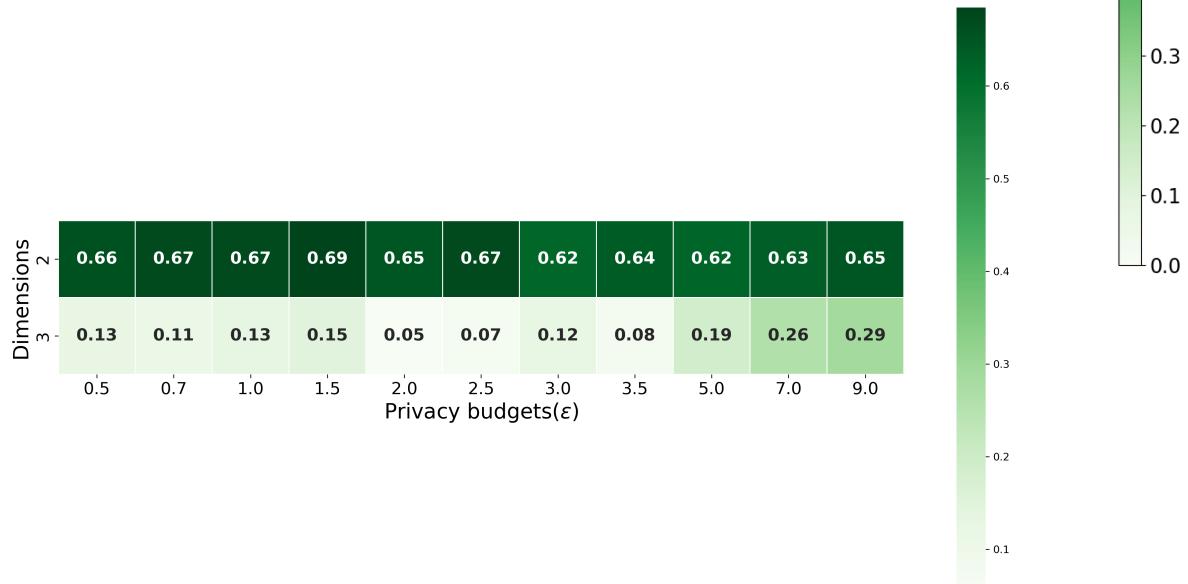


### 6.3.4. LINE DATASET

(a) Heatmap showing adversary advantage for the kD-Laplace mechanism, per privacy budget & dimension for seeds-dataset.



(b) Heatmap showing adversary advantage for the Piecewise mechanism, per privacy budget & dimension for seeds-dataset.



### 6.3.5. SKEWED DATASET

(a) Heatmap showing adversary advantage for the kD-Laplace mechanism, per privacy budget & dimension for seeds-dataset.



(b) Heatmap showing adversary advantage for the Piecewise mechanism, per privacy budget & dimension for seeds-dataset.



### 6.3.6. MECHANISM COMPARISON

In this section, we compare the different mechanisms for each dataset. For this purpose, we also include all the different variants of kd-Laplace to see if there is a difference between them. So, instead of comparing the mechanisms based on the number of dimensions, we compare them on the average scores for all dimensions per mechanism. We are most interested in the utility and performance, and to compare them, we only show the **AMI** and adversary advantage scores.

Also needs privacy distance comparison

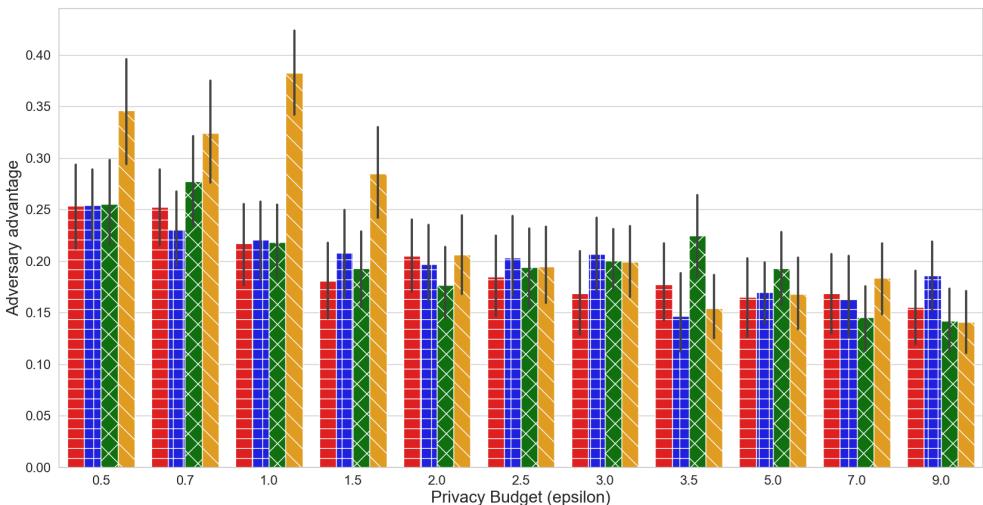
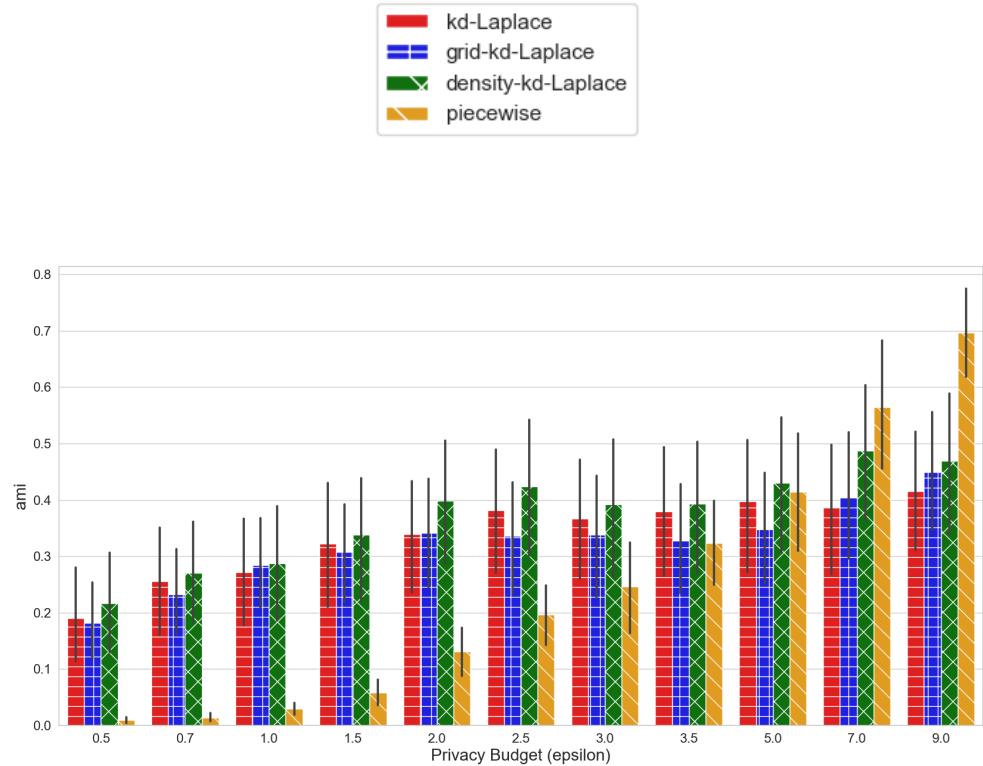


Figure 6.23: Average AMI (top) and Adversary Advantage (bottom) comparison for each mechanism for seeds-dataset (8 dimensions).

Add interpretation

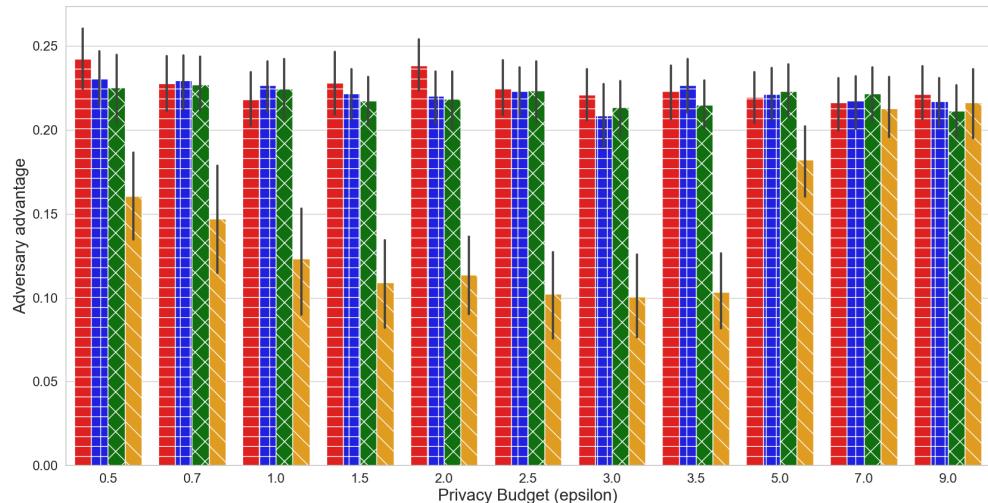
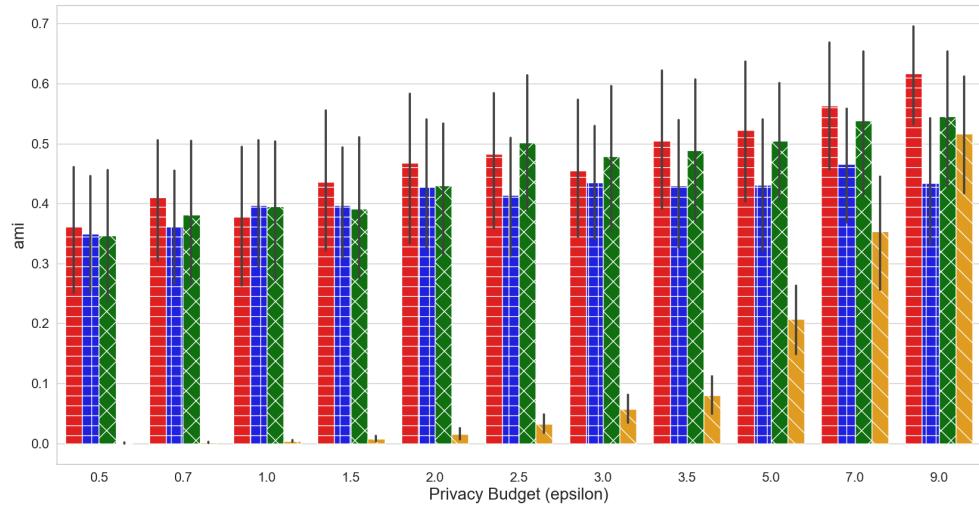


Figure 6.24: Average AMI (top) and Adversary Advantage (bottom) comparison for each mechanism for heart-dataset (10 dimensions).

Add interpretation

### 6.3.7. SHAPE DATASETS

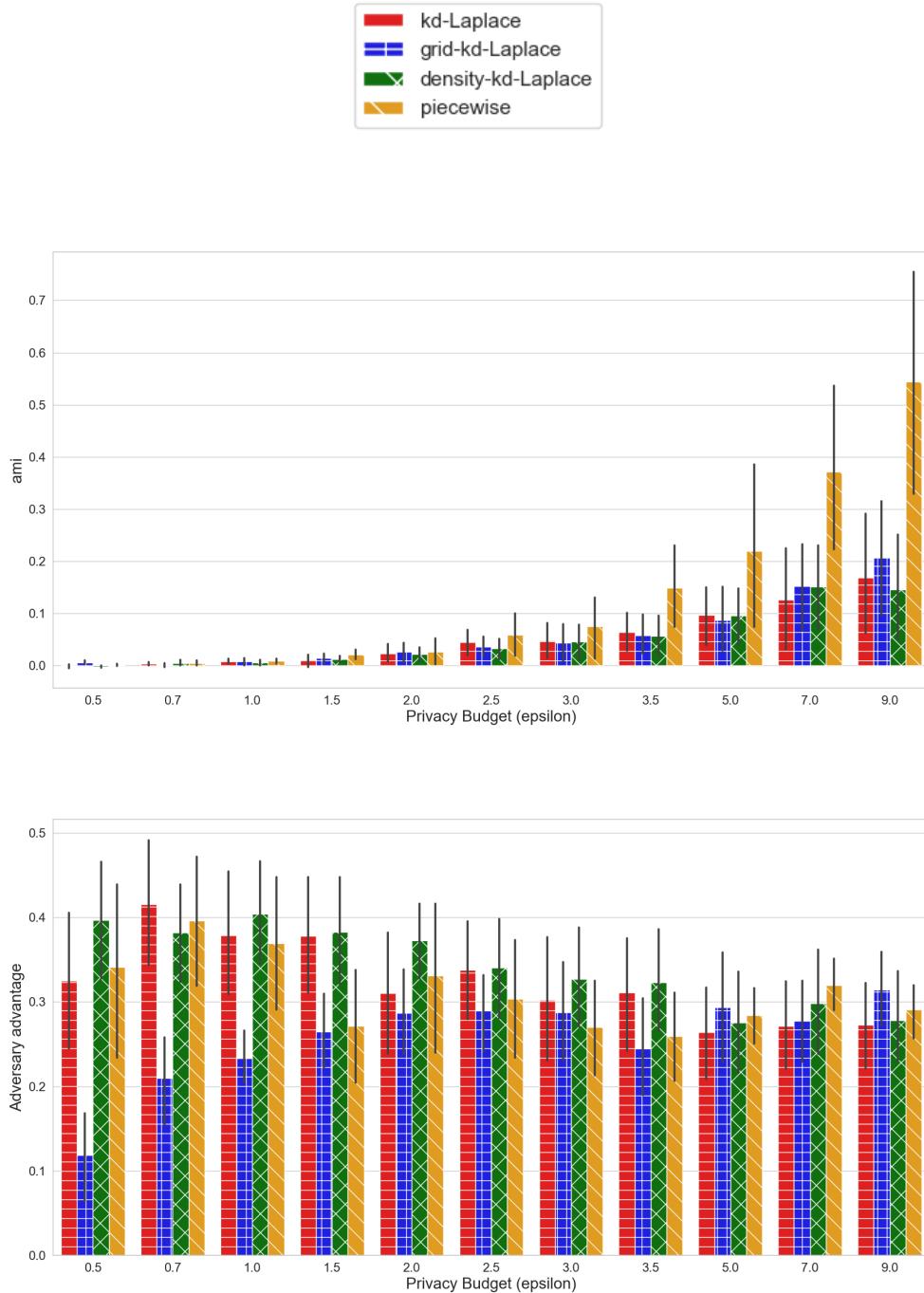


Figure 6.25: Average AMI (top) and Adversary Advantage (bottom) comparison for each mechanism for circle-dataset (3 dimensions).

Add interpretation

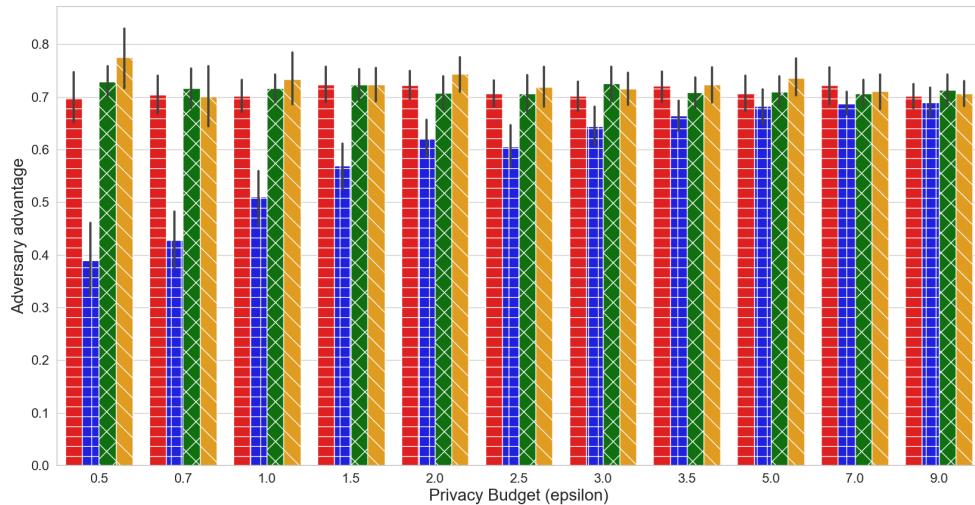
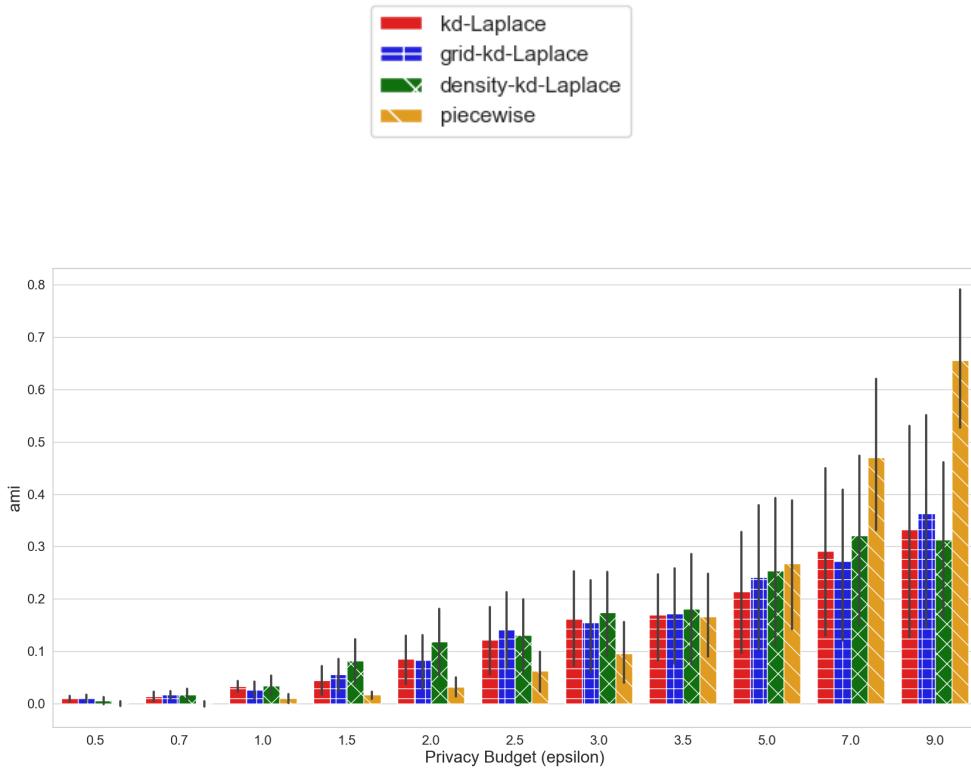


Figure 6.26: Average AMI (top) and Adversary Advantage (bottom) comparison for each mechanism for skewed-dataset (3 dimensions).

Add interpretation

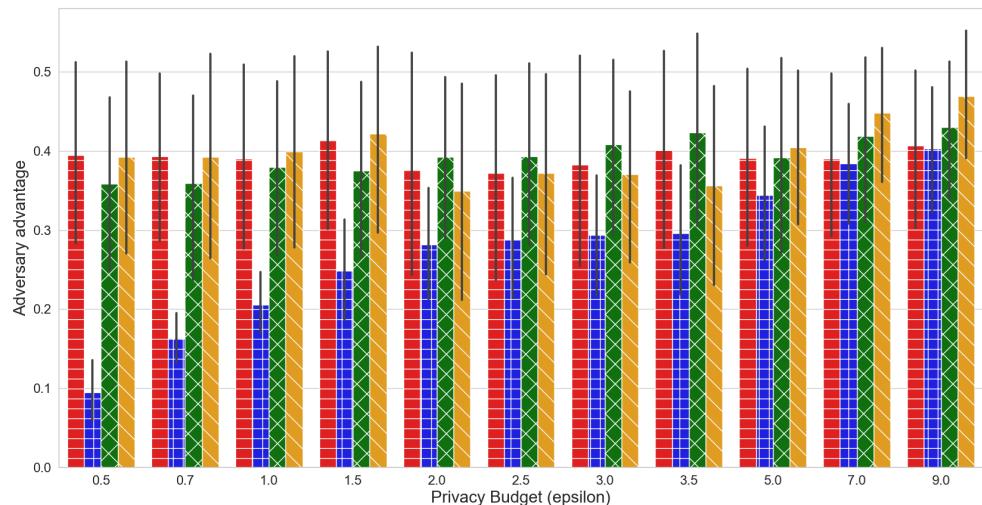
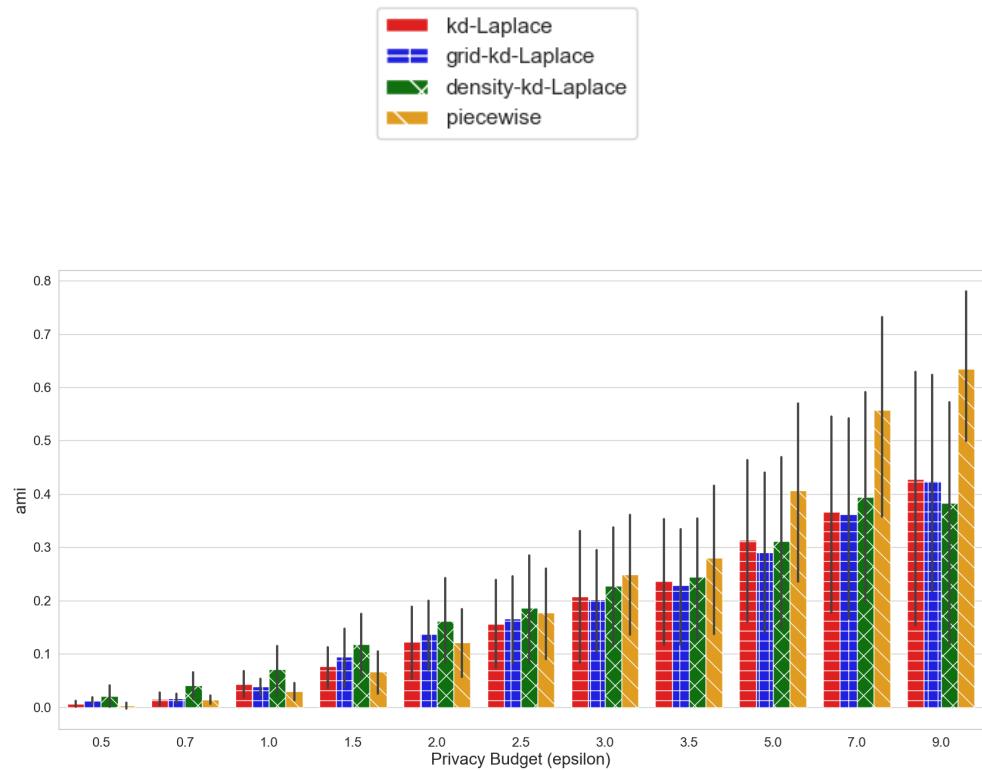


Figure 6.27: Average AMI (top) and Adversary Advantage (bottom) comparison for each mechanism for line-dataset (3 dimensions).

Add interpretation

# 7

## DISCUSSION

1. Why does kd-Laplace perform poorly on shape datasets?
2. Why does OPTICS perform so poorly for clustering with both privacy mechanisms?
3. Why does the Heart dataset perform better than the Seeds dataset?

# 8

## CONCLUSION

This thesis has explored the application of the kD-Laplace algorithm in training privacy-preserving clustering algorithms on distributed k-dimensional data. The research was guided by three research questions, which have been addressed as follows:

**RQ1: How can 2D-Laplace be used to protect the privacy of 2-dimensional data employed for training clustering algorithms?**

Implementing kD-Laplace on n-dimensional data has shown promising results, with the mechanism balancing data utility and privacy protection. For the Heart and Seeds datasets, the kD-Laplace mechanism scores better in **AMI** and adversary advantage for lower privacy budgets (0.1 - 3/5) between kD-Laplace and Piecewise. The Piecewise mechanism scores better in data utility for higher privacy budgets (5 - 9).

**RQ2: How can 3D-Laplace be used to protect the privacy of 3-dimensional data employed for training clustering algorithms?**

### Results

**RQ3: How can nD-Laplace be used to protect the privacy of n-dimensional data employed for training clustering algorithms?**

This research question consists of three hypotheses, which are addressed as follows:

- H1: Adding remapping based on density improves utility without sacrificing privacy:

### In progress

- H2: The privacy leakage (adversary advantage) and utility increases for more dimensions.

### In progress

But especially for the higher privacy budgets.

- H3: The shape of the data negatively impacts the kd-Laplace mechanism in terms of privacy and utility. The kD-Laplace mechanism scores worse for **AMI** for the shape datasets in comparison to the real-world datasets (Heart and Seeds). Piecewise, conversely, scored more consistently and better in **AMI** for the shape datasets. Regarding

the adversary advantage, kD-Laplace scores higher than Piecewise. Upon closer examination, we observe that the True Positive Rate (TPR) is also higher for kD-Laplace. For both the Circle and Line datasets, kD-Laplace significantly outperforms Piecewise. Therefore, the shape of the data negatively impacts the kd-Laplace mechanism in terms of privacy and utility.

#### Still in progress

Our investigation into the utility and privacy differences among the three kd-Laplace variants and the Piecewise mechanism has revealed that the kd-Laplace mechanism can effectively protect the privacy of n-dimensional data. The kd-Laplace mechanism generally scores better in data utility and privacy than the Piecewise mechanism for lower privacy budgets (0.1 - 3/5). The findings of this research have significant implications for the field of privacy-preserving data analysis. This research approach presents a fresh take on a solution for protecting the privacy of n-dimensional data used in clustering algorithms. Also, this thesis presents a new take on evaluating privacy mechanisms using real-world data and attacks.

However, further research is needed to address the limitations identified in this study. In particular, the impact of the shape of the data on the mechanism's effectiveness, also for dimensions higher than 2-dimensions. In addition, it is also worth exploring the impact of the number of dimensions on the privacy and utility of the mechanism.

In conclusion, this thesis has contributed to understanding how the kd-Laplace algorithm can be applied in training privacy-preserving clustering algorithms on distributed k-dimensional data. The findings provide a foundation for future research in this area, with the potential to advance the field of privacy-preserving data analysis.

# BIBLIOGRAPHY

EUR-Lex - 32018R1725 - EN - EUR-Lex. <https://eur-lex.europa.eu/eli/reg/2018/1725/oj>, a. 1

Global big data industry market size 2011-2027. <https://www.statista.com/statistics/254266/global-big-data-market-forecast/>, b. 1

UCI Machine Learning Repository: Cardiotocography Data Set. <https://archive.ics.uci.edu/ml/datasets/cardiotocography>, c. 59

Mohiuddin Ahmed, Raihan Seraj, and Syed Mohammed Shamsul Islam. The k-means algorithm: A comprehensive survey and performance evaluation. *Electronics*, 9(8):1295, 2020. ISSN 2079-9292. 12

Muhammad Aitsam. Differential Privacy Made Easy, December 2021. 6

Miguel E. Andrés, Nicolás Emilio Bordenabe, Konstantinos Chatzikokolakis, and Catuscia Palamidessi. Geo-indistinguishability: Differential privacy for location-based systems. *CoRR*, abs/1212.1984, 2012. 2, 8, 9, 28, 29, 30, 31, 32, 38, 39, 42, 44, 47

Mihai Ankerst, Markus M Breunig, Hans-Peter Kriegel, and Jörg Sander. OPTICS: Ordering Points To Identify the Clustering Structure. 11, 12

Samah Baraheem and Zhongmei Yao. A Survey on Differential Privacy with Machine Learning and Future Outlook, November 2022. 17

Raef Bassily and Adam Smith. Local, Private, Efficient Protocols for Succinct Histograms. In *Proceedings of the Forty-Seventh Annual ACM Symposium on Theory of Computing*, pages 127–135, June 2015. doi: 10.1145/2746539.2746632. 19

Jon Louis Bentley. Multidimensional binary search trees used for associative searching. *Communications of the ACM*, 18(9):509–517, September 1975. ISSN 0001-0782. doi: 10.1145/361002.361007. 43

Avrim Blum, Cynthia Dwork, Frank McSherry, and Kobbi Nissim. *Practical Privacy: The SulQ Framework*. June 2005. doi: 10.1145/1065167.1065184. 20

Beyza Bozdemir, Sébastien Canard, Orhan Ermis, Helen Möllering, Melek Önen, and Thomas Schneider. Privacy-preserving Density-based Clustering. 22, 23

Hanbo Cai, Jinyan Wang, Xiaohong Liu, and Xianxian Li. DP-AP: Differential Privacy-Preserving Affinity Propagation Clustering. In *2020 IEEE 14th International Conference on Big Data Science and Engineering (BigDataSE)*, pages 73–79, Guangzhou, China, December 2020. IEEE. ISBN 978-1-66540-396-2. doi: 10.1109/BigDataSE50710.2020.00018. 21, 23

Tadeusz Caliński and Jerzy Harabasz. A dendrite method for cluster analysis. *Communications in Statistics-theory and Methods*, 3(1):1–27, 1974. ISSN 0090-3272. [14](#), [ix](#)

Nicholas Carlini, Steve Chien, Milad Nasr, Shuang Song, Andreas Terzis, and Florian Tramer. Membership Inference Attacks From First Principles, April 2022. [55](#)

Konstantinos Chatzikokolakis, Catuscia Palamidessi, and Marco Stronati. Constructing elastic distinguishability metrics for location privacy. *Proceedings on Privacy Enhancing Technologies*, 2015(2):156–170, June 2015. ISSN 2299-0984. doi: 10.1515/popets-2015-0023. [9](#)

Konstantinos Chatzikokolakis, Ehab Elsalamouny, and Catuscia Palamidessi. Efficient utility improvement for location privacy. *Proceedings on Privacy Enhancing Technologies*, 2017(4):308–328, 2017. [42](#), [46](#), [47](#)

Jianbo Chen, Michael I. Jordan, and Martin J. Wainwright. HopSkipJumpAttack: A Query-Efficient Decision-Based Attack, April 2020.

Christopher A. Choquette-Choo, Florian Tramer, Nicholas Carlini, and Nicolas Papernot. Label-Only Membership Inference Attacks, December 2021. [55](#)

Toon Van Craenendonck and Hendrik Blockeel. Using Internal Validity Measures to Compare Clustering Algorithms. [14](#)

David L Davies and Donald W Bouldin. A cluster separation measure. *IEEE transactions on pattern analysis and machine intelligence*, (2):224–227, 1979. ISSN 0162-8828. [15](#)

Irit Dinur and Kobbi Nissim. Revealing information while preserving privacy. In *Proceedings of the Twenty-Second ACM SIGMOD-SIGACT-SIGART Symposium on Principles of Database Systems*, pages 202–210, San Diego California, June 2003. ACM. ISBN 978-1-58113-670-8. doi: 10.1145/773153.773173. [54](#)

Freedom of Information Act (FOIA) Division. The Privacy Act. <https://www.hhs.gov/foia/privacy/index.html>, April 2007. [1](#)

Salim Dridi. *Unsupervised Learning - A Systematic Literature Review*. December 2021. doi: 10.13140/RG.2.2.16963.12323. [1](#)

John Duchi, Martin Wainwright, and Michael Jordan. Minimax Optimal Procedures for Locally Private Estimation, November 2017. [19](#), [24](#), [26](#)

John C. Duchi, Michael I. Jordan, and Martin J. Wainwright. Privacy Aware Learning, October 2013. [19](#), [26](#)

Cynthia Dwork. Differential privacy. In *Automata, Languages and Programming: 33rd International Colloquium, ICALP 2006, Venice, Italy, July 10-14, 2006, Proceedings, Part II* 33, pages 1–12. Springer, 2006. ISBN 3-540-35907-9. [1](#), [6](#), [9](#), [10](#), [19](#)

Cynthia Dwork, Adam Smith, Thomas Steinke, and Jonathan Ullman. Exposed! A Survey of Attacks on Private Data. *Annual Review of Statistics and Its Application*, 4(1):61–84, March 2017. ISSN 2326-8298, 2326-831X. doi: 10.1146/annurev-statistics-060116-054123. [54](#)

Martin Ester, Hans-Peter Kriegel, and Xiaowei Xu. A Density-Based Algorithm for Discovering Clusters in Large Spatial Databases with Noise. 11

Adil Fahad, Najlaa Alshatri, Zahir Tari, Abdullah Alamri, Ibrahim Khalil, Albert Y. Zomaya, Sebti Foufou, and Abdelaziz Bouras. A Survey of Clustering Algorithms for Big Data: Taxonomy and Empirical Analysis. *IEEE Transactions on Emerging Topics in Computing*, 2(3):267–279, September 2014. ISSN 2168-6750. doi: 10.1109/TETC.2014.2330519. 11

Natasha Fernandes, Mark Dras, and Annabelle McIver. Generalised Differential Privacy for Text Document Processing, February 2019. 2, 39

Pasi Fränti, Mohammad Rezaei, and Qinpei Zhao. Centroid index: Cluster level similarity measure. *Pattern Recognition*, 47(9):3034–3045, September 2014. ISSN 00313203. doi: 10.1016/j.patcog.2014.03.017. 14

Matt Fredrikson, Somesh Jha, and Thomas Ristenpart. Model Inversion Attacks that Exploit Confidence Information and Basic Countermeasures. In *Proceedings of the 22nd ACM SIGSAC Conference on Computer and Communications Security*, pages 1322–1333, Denver Colorado USA, October 2015. ACM. ISBN 978-1-4503-3832-5. doi: 10.1145/2810103.2813677. 54

Arik Friedman and Assaf Schuster. Data mining with differential privacy. In *Proceedings of the 16th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 493–502, Washington DC USA, July 2010. ACM. ISBN 978-1-4503-0055-1. doi: 10.1145/1835804.1835868. 6

Quan Geng and Pramod Viswanath. The Optimal Mechanism in Differential Privacy, October 2013. 24

Quan Geng, Peter Kairouz, Sewoong Oh, and Pramod Viswanath. The staircase mechanism in differential privacy. *IEEE Journal of Selected Topics in Signal Processing*, 9(7):1176–1184, October 2015. ISSN 1932-4553, 1941-0484. doi: 10.1109/JSTSP.2015.2425831. 19, 24

Muneeb Ul Hassan, Mubashir Husain Rehmani, and Jinjun Chen. Differential Privacy Techniques for Cyber Physical Systems: A Survey, September 2019. 17

Marwan Hassani and Thomas Seidl. Using internal evaluation measures to validate the quality of diverse stream clustering algorithms. *Vietnam Journal of Computer Science*, 4(3):171–183, August 2017. ISSN 2196-8896. doi: 10.1007/s40595-016-0086-9. 14, 15

Hongsheng Hu, Zoran Salcic, Lichao Sun, Gillian Dobbie, Philip S. Yu, and Xuyun Zhang. Membership Inference Attacks on Machine Learning: A Survey, February 2022. 52, 55

D. Huang, X. Yao, S. An, and S. Ren. Private distributed K-means clustering on interval data. In *2021 IEEE International Performance, Computing, and Communications Conference (IPCCC)*, pages 1–9, Los Alamitos, CA, USA, October 2021. IEEE Computer Society. doi: 10.1109/IPCCC51483.2021.9679364. 2, 21, 23, 26, 60

Lawrence Hubert and Phipps Arabie. Comparing partitions. *Journal of classification*, 2:193–218, 1985. ISSN 0176-4268. 14, ix

Bargav Jayaraman and David Evans. Evaluating Differentially Private Machine Learning in Practice. [52](#), [54](#)

Bargav Jayaraman and David Evans. Are Attribute Inference Attacks Just Imputation? In *Proceedings of the 2022 ACM SIGSAC Conference on Computer and Communications Security*, pages 1569–1582, Los Angeles CA USA, November 2022. ACM. ISBN 978-1-4503-9450-5. doi: [10.1145/3548606.3560663](#). [54](#)

Marija Jegorova, Chaitanya Kaul, Charlie Mayor, Alison Q. O’Neil, Alexander Weir, Roderick Murray-Smith, and Sotirios A. Tsaftaris. Survey: Leakage and Privacy at Inference Time, September 2022. [54](#)

Zhanglong Ji, Zachary C. Lipton, and Charles Elkan. Differential Privacy and Machine Learning: A Survey and Review, December 2014. [17](#)

Matthew Joseph, Jieming Mao, Seth Neel, and Aaron Roth. The Role of Interactivity in Local Differential Privacy. In *2019 IEEE 60th Annual Symposium on Foundations of Computer Science (FOCS)*, pages 94–105, Baltimore, MD, USA, November 2019. IEEE. ISBN 978-1-72814-952-3. doi: [10.1109/FOCS.2019.00015](#). [8](#)

Haim Kaplan and Uri Stemmer. Differentially Private k-Means with Constant Multiplicative Error, July 2018. [20](#)

Shiva Prasad Kasiviswanathan, Homin K. Lee, Kobbi Nissim, Sofya Raskhodnikova, and Adam Smith. What Can We Learn Privately?, February 2010. [1](#)

Trupti M Kodinariya and Prashant R Makwana. Review on determining number of Cluster in K-Means Clustering. *International Journal*, 1(6):90–95, 2013. [13](#)

Jussi Lehtonen. The Lambert W function in ecological and evolutionary models. *Methods in Ecology and Evolution*, 7(9):1110–1118, 2016. ISSN 2041-210X. doi: [10.1111/2041-210X.12568](#). [29](#)

Zheng Li and Yang Zhang. Membership Leakage in Label-Only Exposures, September 2021.

Bo Liu, Ming Ding, Sina Shaham, Wenny Rahayu, Farhad Farokhi, and Zihuai Lin. When Machine Learning Meets Privacy: A Survey and Outlook. *ACM Computing Surveys*, 54(2):1–36, March 2022. ISSN 0360-0300, 1557-7341. doi: [10.1145/3436755](#). [1](#)

Jinfei Liu, Joshua Huang, Jun Luo, and Li Xiong. Privacy preserving distributed DBSCAN clustering. *Transactions on Data Privacy*, 6, March 2012. doi: [10.1145/2320765.2320819](#). [11](#)

Yanchi Liu, Zhongmou Li, Hui Xiong, Xuedong Gao, and Junjie Wu. Understanding of internal clustering validation measures. In *2010 IEEE International Conference on Data Mining*, pages 911–916. IEEE, 2010. ISBN 1-4244-9131-2. [15](#)

S. Lloyd. Least squares quantization in PCM. *IEEE Transactions on Information Theory*, 28(2):129–137, 1982. doi: [10.1109/TIT.1982.1056489](#). [11](#)

George Marsaglia. Choosing a point from the surface of a sphere. *The Annals of Mathematical Statistics*, 43(2):645–646, 1972. ISSN 0003-4851.

Xianrui Meng, Dimitrios Papadopoulos, Alina Oprea, and Nikos Triandopoulos. Private Hierarchical Clustering and Efficient Approximation. In *Proceedings of the 2021 on Cloud Computing Security Workshop*, pages 3–20, November 2021. doi: 10.1145/3474123.3486760. [12](#)

Minghui Min, Liang Xiao, Jiahao Ding, Hongliang Zhang, Shiycin Li, Miao Pan, and Zhu Han. 3D Geo-Indistinguishability for Indoor Location-Based Services. *IEEE Transactions on Wireless Communications*, 21(7):4682–4694, 2022. doi: 10.1109/TWC.2021.3132464. [2](#), [28](#), [30](#), [33](#), [34](#), [35](#), [36](#), [39](#), [44](#), [47](#)

Thôong T. Nguyễn, Xiaokui Xiao, Yin Yang, Siu Cheung Hui, Hyejin Shin, and Junbum Shin. Collecting and Analyzing Data from Smart Device Users with Local Differential Privacy, June 2016. [19](#), [21](#), [24](#), [26](#)

Maria-Irina Nicolae, Mathieu Sinn, Minh Ngoc Tran, Beat Buesser, Ambrish Rawat, Martin Wistuba, Valentina Zantedeschi, Nathalie Baracaldo, Bryant Chen, Heiko Ludwig, Ian M. Molloy, and Ben Edwards. Adversarial Robustness Toolbox v1.0.0, November 2019. [54](#), [64](#)

Kobbi Nissim and Uri Stemmer. Clustering Algorithms for the Centralized and Local Models. In *Proceedings of Algorithmic Learning Theory*, pages 619–653. PMLR, April 2018. [20](#), [23](#)

Kobbi Nissim, Sofya Raskhodnikova, and Adam Smith. Smooth sensitivity and sampling in private data analysis. In *Proceedings of the Thirty-Ninth Annual ACM Symposium on Theory of Computing*, pages 75–84, San Diego California USA, June 2007. ACM. ISBN 978-1-59593-631-8. doi: 10.1145/1250790.1250803. [10](#), [20](#), [23](#)

Yuefeng Peng, Bo Zhao, and Hui Liu. Unsupervised Membership Inference Attacks Against Machine Learning Models. [53](#)

William M Rand. Objective criteria for the evaluation of clustering methods. *Journal of the American Statistical association*, 66(336):846–850, 1971. ISSN 0162-1459. [14](#)

Maria Rigaki and Sebastian Garcia. A Survey of Privacy Attacks in Machine Learning, April 2021. [52](#), [53](#), [54](#), [55](#), [64](#)

Peter J Rousseeuw. Silhouettes: A graphical aid to the interpretation and validation of cluster analysis. *Journal of computational and applied mathematics*, 20:53–65, 1987. ISSN 0377-0427. [14](#), [15](#)

Maurice Roux. A comparative study of divisive hierarchical clustering algorithms, September 2015. [12](#), [13](#)

Jörg Sander, Martin Ester, Hans-Peter Kriegel, and Xiaowei Xu. Density-Based Clustering in Spatial Databases: The Algorithm GDBSCAN and Its Applications. *Data Mining and Knowledge Discovery*, 2(2):169–194, June 1998. ISSN 1573-756X. doi: 10.1023/A:1009745219419. [14](#)

Danny Matthew SAPUTRA, Daniel SAPUTRA, and Liniyanti D OSWARI. Effect of distance metrics in determining k-value in k-means clustering using elbow and silhouette method. In *Sriwijaya International Conference on Information Technology and Its Applications (SICONIAN 2019)*, pages 341–346. Atlantis Press, 2020. ISBN 94-6252-963-9. [13](#)

Erich Schubert, Jörg Sander, Martin Ester, Hans Peter Kriegel, and Xiaowei Xu. DBSCAN Revisited, Revisited: Why and How You Should (Still) Use DBSCAN. *ACM Transactions on Database Systems*, 42(3):1–21, September 2017. ISSN 0362-5915, 1557-4644. doi: 10.1145/3068335. [14](#)

Sreedhar Kumar Seetharaman, Madheswaran Muthusamy, Vinutha A, Manjunatha H, and Charan K V. A Brief Survey of Unsupervised Agglomerative Hierarchical Clustering Schemes. *International Journal of Engineering and Technology*, 8:29–37, January 2019. doi: 10.14419/ijet.v8i1.13971. [13](#)

Reza Shokri, Marco Stronati, Congzheng Song, and Vitaly Shmatikov. Membership Inference Attacks against Machine Learning Models, March 2017. [53](#), [65](#)

Jordi Soria-Comas and Josep Domingo-Ferrer. Optimal data-independent noise for differential privacy. *Information Sciences*, 250:200–214, November 2013. ISSN 0020-0255. doi: 10.1016/j.ins.2013.07.004. [19](#)

Uri Stemmer. Locally private k-means clustering. *The Journal of Machine Learning Research*, 22(1):7964–7993, 2021. ISSN 1532-4435. [20](#), [23](#)

Alexander Strehl and Joydeep Ghosh. Cluster ensembles—a knowledge reuse framework for combining multiple partitions. *Journal of machine learning research*, 3(Dec):583–617, 2002. [14](#)

Dong Su, Jianneng Cao, Ninghui Li, Elisa Bertino, and Hongxia Jin. Differentially Private \$k\\$-Means Clustering, April 2015. [20](#), [21](#), [23](#)

Lin Sun, Jun Zhao, and Xiaojun Ye. Distributed Clustering in the Anonymized Space with Local Differential Privacy, June 2019. [2](#), [17](#), [23](#), [26](#)

Lin Sun, Guolou Ping, and Xiaojun Ye. PrivBV: Distance-aware encoding for distributed data with local differential privacy. *Tsinghua Science and Technology*, 27(2):412–421, April 2022. ISSN 1007-0214. doi: 10.26599/TST.2021.9010027. [23](#)

Latanya Sweeney. K-ANONYMITY: A MODEL FOR PROTECTING PRIVACY. *International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems*, 10(05):557–570, October 2002. ISSN 0218-4885, 1793-6411. doi: 10.1142/S0218488502001648. [1](#)

Robert Tibshirani, Guenther Walther, and Trevor Hastie. Estimating the number of clusters in a data set via the gap statistic. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 63(2):411–423, 2001. ISSN 1369-7412. [13](#)

Nguyen Xuan Vinh, Julien Epps, and James Bailey. Information Theoretic Measures for Clusterings Comparison: Variants, Properties, Normalization and Correction for Chance. [14](#), [15](#), [16](#)

Silke Wagner and Dorothea Wagner. Comparing Clusterings - An Overview. [14](#)

Ning Wang, Xiaokui Xiao, Yin Yang, Jun Zhao, Siu Cheung Hui, Hyejin Shin, Junbum Shin, and Ge Yu. Collecting and Analyzing Multidimensional Data with Local Differential Privacy, June 2019. [19](#), [24](#), [26](#), [27](#)

Teng Wang, Xuefeng Zhang, Jingyu Feng, and Xinyu Yang. A Comprehensive Survey on Local Differential Privacy toward Data Statistics and Analysis. *Sensors*, 20(24), 2020. ISSN 1424-8220. doi: 10.3390/s20247030. [8](#), [49](#)

Tianhao Wang, Jeremiah Blocki, Ninghui Li, and Somesh Jha. Locally Differentially Private Protocols for Frequency Estimation. [19](#), [21](#)

Matthijs J. Warrens and Hanneke van der Hoef. Understanding the Adjusted Rand Index and Other Partition Comparison Indices Based on Counting Object Pairs. *Journal of Classification*, 39(3):487–509, November 2022. ISSN 1432-1343. doi: 10.1007/s00357-022-09413-z. [14](#)

Washington. K-D Trees, February 2002. [43](#), [44](#)

Chang Xia, Jingyu Hua, Wei Tong, and Sheng Zhong. Distributed K-Means clustering guaranteeing local differential privacy. *Computers & Security*, 90:101699, 2020. ISSN 0167-4048. [2](#), [21](#), [23](#), [26](#)

Xingxing Xiong, Shubo Liu, Dan Li, Zhaojun Cai, and Xiaoguang Niu. A Comprehensive Survey on Local Differential Privacy. *Security and Communication Networks*, 2020:8829523, October 2020a. ISSN 1939-0114. doi: 10.1155/2020/8829523. [7](#), [9](#), [10](#)

Xingxing Xiong, Shubo Liu, Dan Li, Zhaojun Cai, and Xiaoguang Niu. A Comprehensive Survey on Local Differential Privacy. *Security and Communication Networks*, 2020:8829523, October 2020b. ISSN 1939-0114. doi: 10.1155/2020/8829523. [8](#), [49](#)

Dongkuan Xu and Yingjie Tian. A Comprehensive Survey of Clustering Algorithms. *Annals of Data Science*, 2(2):165–193, June 2015. ISSN 2198-5812. doi: 10.1007/s40745-015-0040-1. [11](#)

Mengmeng Yang, Lingjuan Lyu, Jun Zhao, Tianqing Zhu, and Kwok-Yan Lam. Local Differential Privacy and Its Applications: A Comprehensive Survey, August 2020. [1](#)

Samuel Yeom, Irene Giacomelli, Matt Fredrikson, and Somesh Jha. Privacy Risk in Machine Learning: Analyzing the Connection to Overfitting. In *2018 IEEE 31st Computer Security Foundations Symposium (CSF)*, pages 268–282, Oxford, July 2018. IEEE. ISBN 978-1-5386-6680-7. doi: 10.1109/CSF2018.00027. [53](#), [54](#), [55](#), [64](#)

Chunhui Yuan and Haitao Yang. Research on K-value selection method of K-means clustering algorithm. *J*, 2(2):226–235, 2019. ISSN 2571-8800. [11](#), [13](#)

Liujiu Yuan, Shaobo Zhang, Gengming Zhu, and Karim Alinani. Privacy-preserving mechanism for mixed data clustering with local differential privacy. *Concurrency and Computation: Practice and Experience*, July 2021. ISSN 1532-0626, 1532-0634. doi: 10.1002/cpe.6503. [21](#), [23](#), [26](#)

Benjamin Zi Hao Zhao, Mohamed Ali Kaafar, and Nicolas Kourtellis. Not one but many Tradeoffs: Privacy Vs. Utility in Differentially Private Machine Learning. In *Proceedings of the 2020 ACM SIGSAC Conference on Cloud Computing Security Workshop*, pages 15–26, November 2020. doi: 10.1145/3411495.3421352. [64](#)

Benjamin Zi Hao Zhao, Aviral Agrawal, Catisha Coburn, Hassan Jameel Asghar, Raghav Bhaskar, Mohamed Ali Kaafar, Darren Webb, and Peter Dickinson. On the (In)Feasibility of Attribute Inference Attacks on Machine Learning Models, March 2021. [54](#)

## GLOSSARY

**Adjusted Mutual Information** Comparable with **Adjusted Rand Index** this algorithm is modified to account to chance. This means it accounts for a higher MI for a higher amount of clusters between two clustering algorithms. Therefore, the calculations are strongly influenced by that of **Adjusted Rand Index** [?]. . 13

**Adjusted Rand Index** The Rand Index is improved and adjusted for chance [Hubert and Arabie, 1985]. This algorithm takes also into consideration the number of clusters and can be used to also compare different clustering algorithms [?]. ix, 13

**Calinski-Harabasz Index** This is a way to measure the similarity of clusters [Caliński and Harabasz, 1974]. It tells how well the clusters are separated from each other and how well the points are grouped.. 13

**Mutual Information** This metric can be used to explain the amount of information about a random variable if compared to another random variable. Therefore, it can also be used to compare two cluster similarities.. ix, 13

**Normalized Mutual Information** The normalized version is a scaled version of **Mutual Information** to always be a value between 0 (no correlation) and 1 (perfect correlation). This version of **Mutual Information** is not adjusted and therefore highly influenced by cluster amount [?]. So it suffers the same issue as with **Mutual Information**.. 13

## ACRONYMS

**AEE** Average Estimated Error. *Glossary:* Average Estimation Error

**AG** Agglomerative Clustering. 57, 63, 64, 66–69, 71–74

**AMI** Adjusted Mutual Information. 13, 14, 21, 60, 62–74, 76–80, 86, 93, *Glossary:* Adjusted Mutual Information

**AP** Affinity Propagation. 10, 14, 20, 22, 56, 58

**ARI** Adjusted Rank Index. 13, 14, 21, *Glossary:* Adjusted Rand Index

**BIRCH** Balanced Iterative Reducing and Clustering using Hierarchies. 11

**CDF** Cumulative Distribution Function. 27

**CHI** Calinski-Harabasz Index. 13, *Glossary:* Calinski-Harabasz Index

**DBSCAN** Density-based spatial clustering of applications with noise. 10–13, 20, 21

**DP** Differential Privacy. 2, 5, 6, 9, 14, 17, 18

**FPR** False Positive Rate. 50

**GI** Geo-indistinguishability. 2, 8, 11, 26, 29, 52

**LDP** Local Differential Privacy. 1, 2, 7, 9, 18, 20, 24, 29, 52, 53

**MI** Mutual Information. 13, 14, *Glossary: Mutual Information*

**MIA** Membership Inference Attack. 47–50, 60

**NMI** Normalized Mutual Information. 13, 14, *Glossary: Normalized Mutual Information*

**OPTICS** Ordering Points to Identify Clustering Structure. 10–13, 58, 64, 68–74

**PDF** Probability Density Function. 18, 24, 27

**PM** Piecewise Mechanism. 18

**SC** Silhouette Coefficient. 13, 14, 57, 62–74

**SVM** Support Vector Machine. 18

**TPR** True Positive Rate. 50, 60

# HYPERPARAMETERS

## .1. K-MEANS

For selecting the appropriate amount of clusters, we used an "elbow" plot in combination with the silhouette score.

1. Seeds dataset: 4 clusters (see figure: 1)
2. Heart dataset: TODO
3. Circle dataset: 5 clusters (see figure: 2)
4. Line dataset: 4 clusters (see figure: 3)
5. Skewed dataset: 4 clusters (see figure: 4)

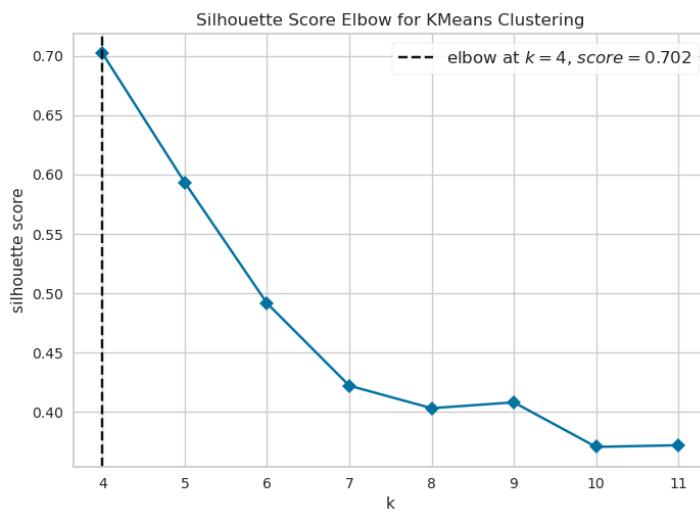


Figure 1: Selecting the  $k$  for K-Means for seeds dataset using the "elbow plot" using section 2.2.2

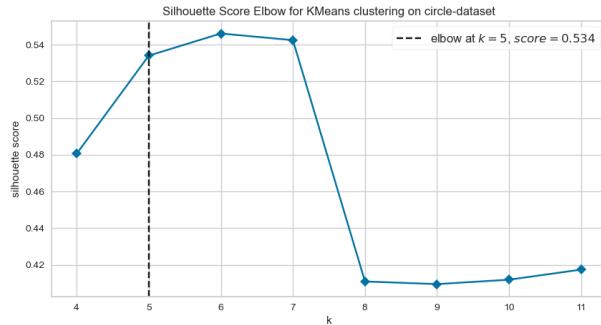


Figure 2: Selecting the  $k$  for K-Means for the circle dataset using the "elbow plot" using section 2.2.2

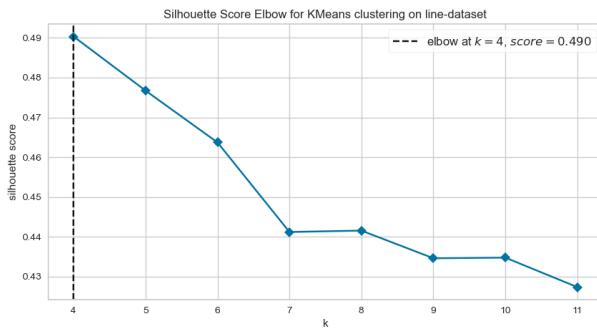


Figure 3: Selecting the  $k$  for K-Means for the line dataset using the "elbow plot" using section 2.2.2

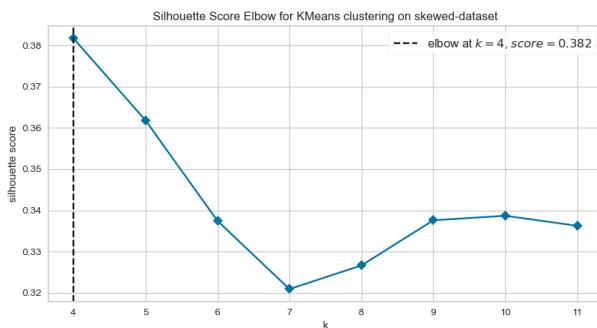


Figure 4: Selecting the  $k$  for K-Means for the skewed dataset using the "elbow plot" using section 2.2.2

## .2. AGGLOMERATIVE CLUSTERING

### .2.1. SEEDS DATASET

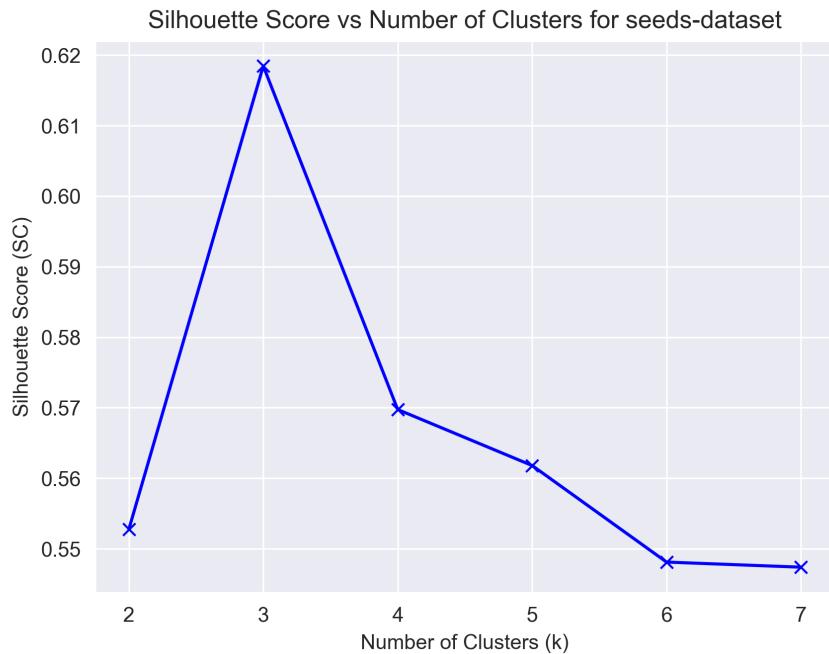


Figure 5: Selecting the  $k$  for Agglomerative clustering for seeds dataset (2 dimensions) using the "elbow plot."

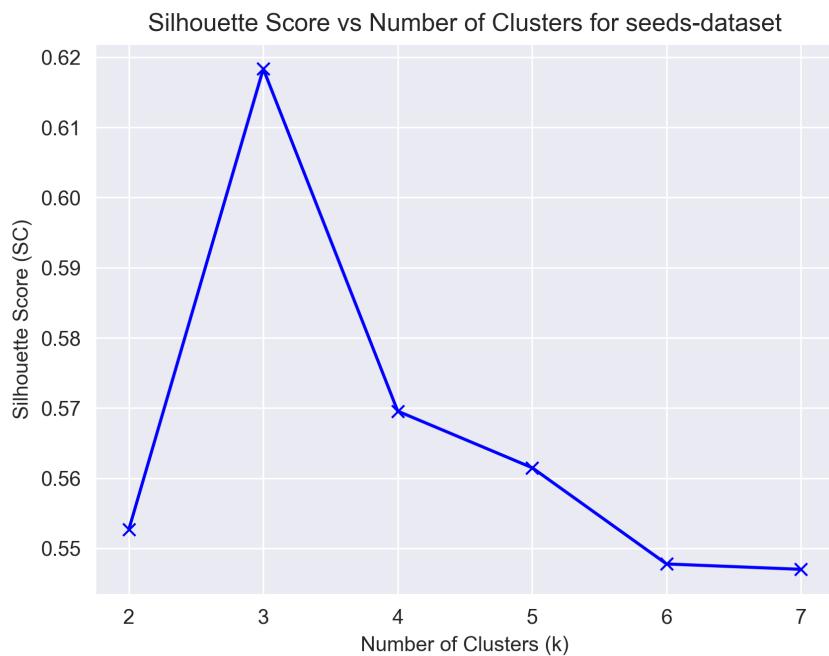


Figure 6: Selecting the  $k$  for Agglomerative clustering for seeds dataset (3 dimensions) using the "elbow plot."

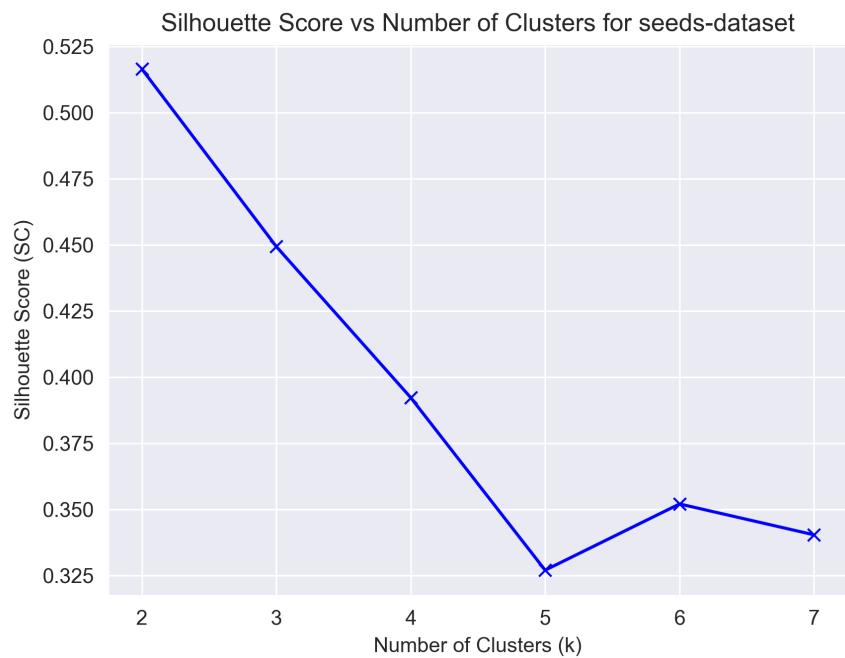


Figure 7: Selecting the  $k$  for Agglomerative clustering for seeds dataset (7 dimensions) using the "elbow plot."

## .2.2. HEART DATASET

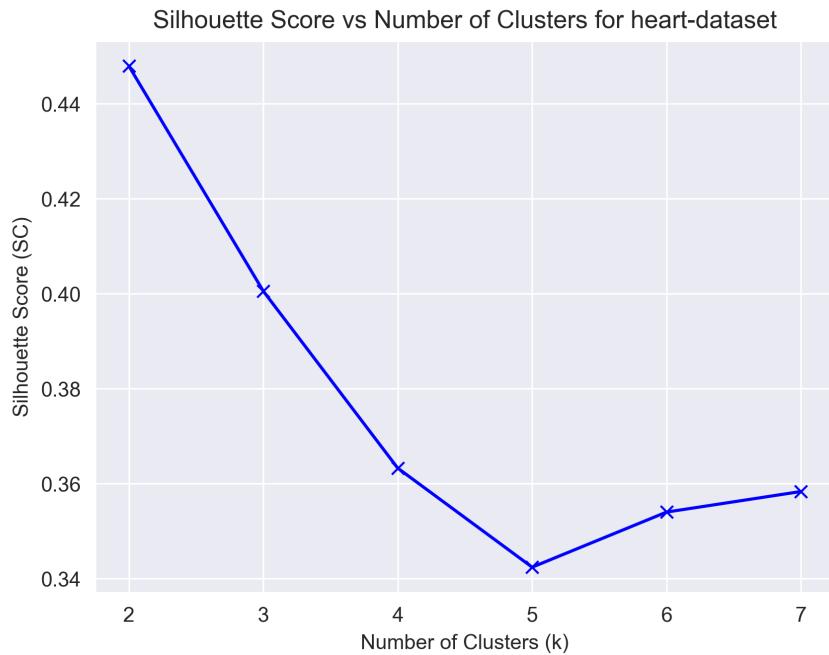


Figure 8: Selecting the  $k$  for Agglomerative clustering for Heart dataset (2 dimensions) using the "elbow plot."

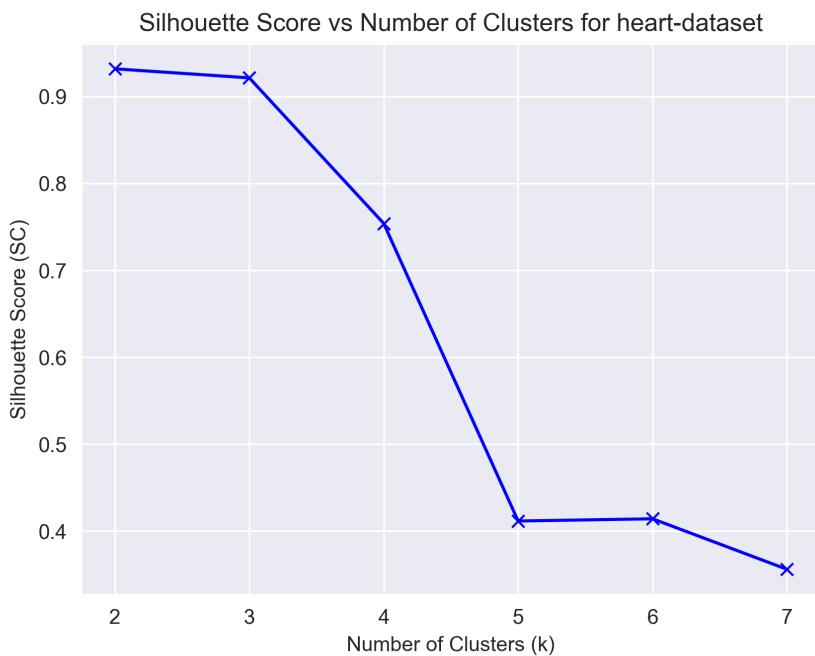


Figure 9: Selecting the  $k$  for Agglomerative clustering for Heart dataset (3 dimensions) using the "elbow plot."

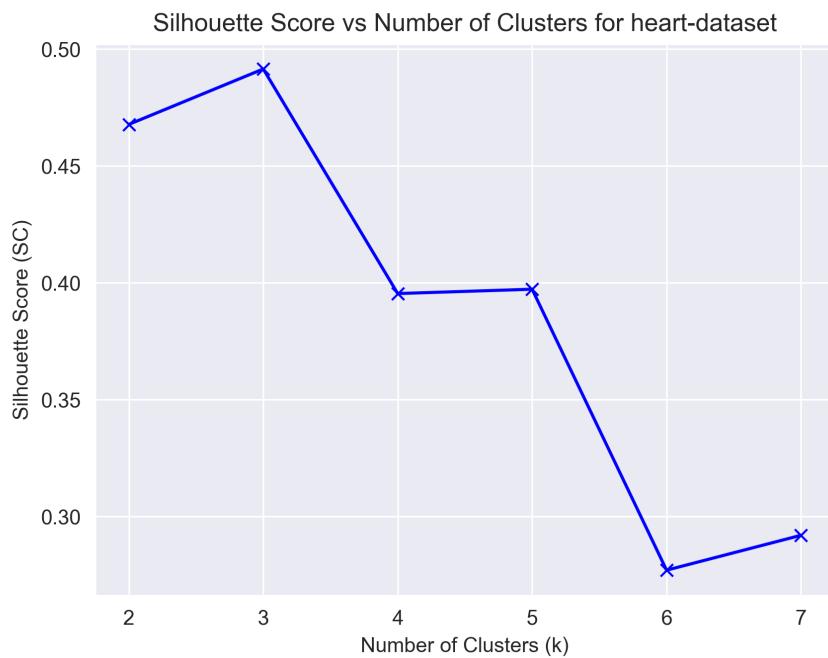


Figure 10: Selecting the  $k$  for Agglomerative clustering for Heart dataset (9 dimensions) using the "elbow plot."

## 2.3. CIRCLE DATASET

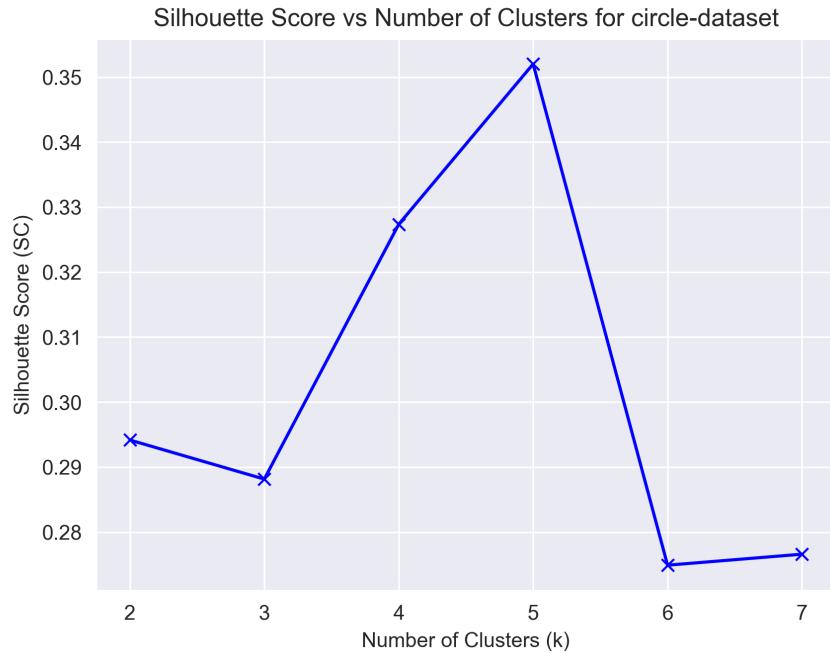


Figure 11: Selecting the  $k$  for Agglomerative clustering for circle dataset (2-dimensions) using the "elbow plot."

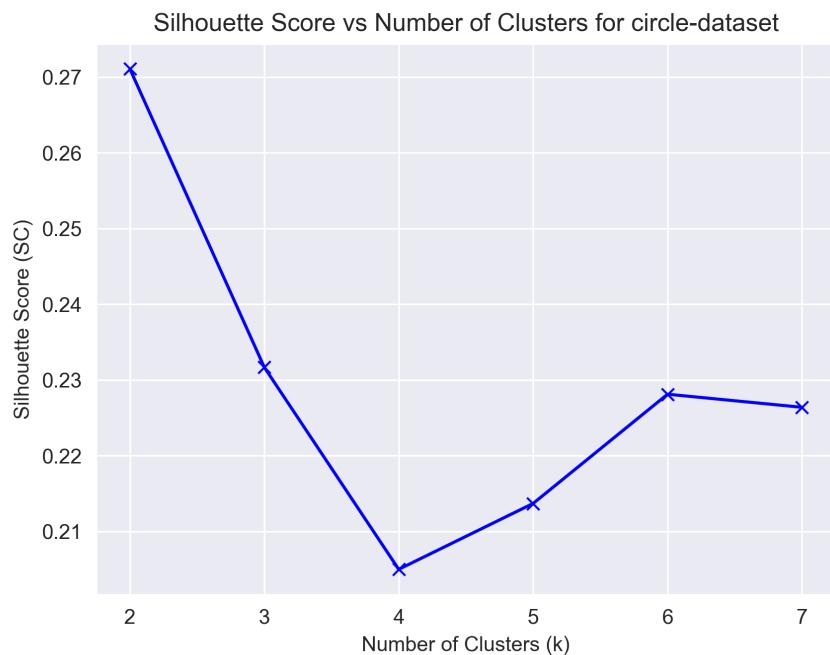


Figure 12: Selecting the  $k$  for Agglomerative clustering for circle dataset (3-dimensions) using the "elbow plot."

## 2.4. LINE DATASET

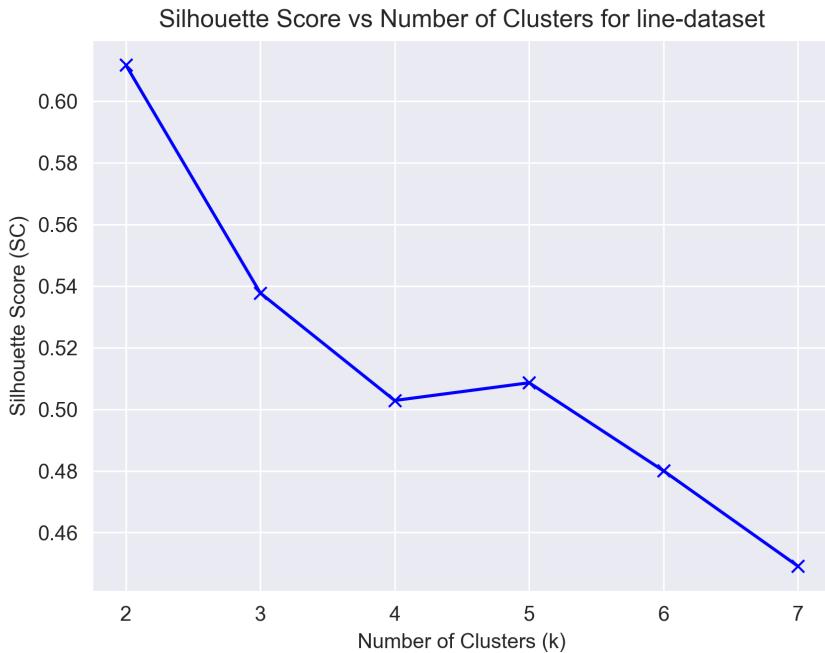


Figure 13: Selecting the  $k$  for Agglomerative clustering for line dataset (2-dimensions) using the "elbow plot."

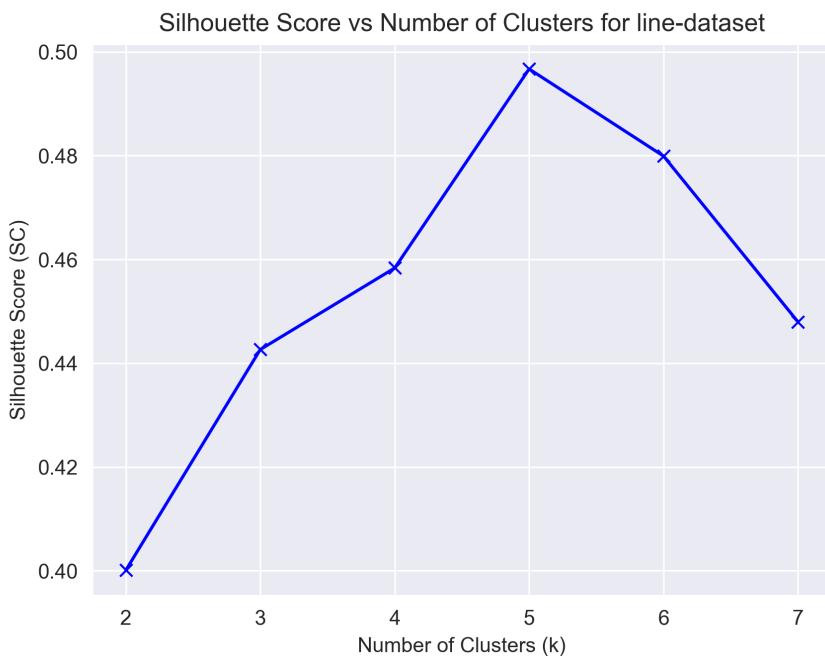


Figure 14: Selecting the  $k$  for Agglomerative clustering for line dataset (3-dimensions) using the "elbow plot."

## .2.5. SKEWED DATASET

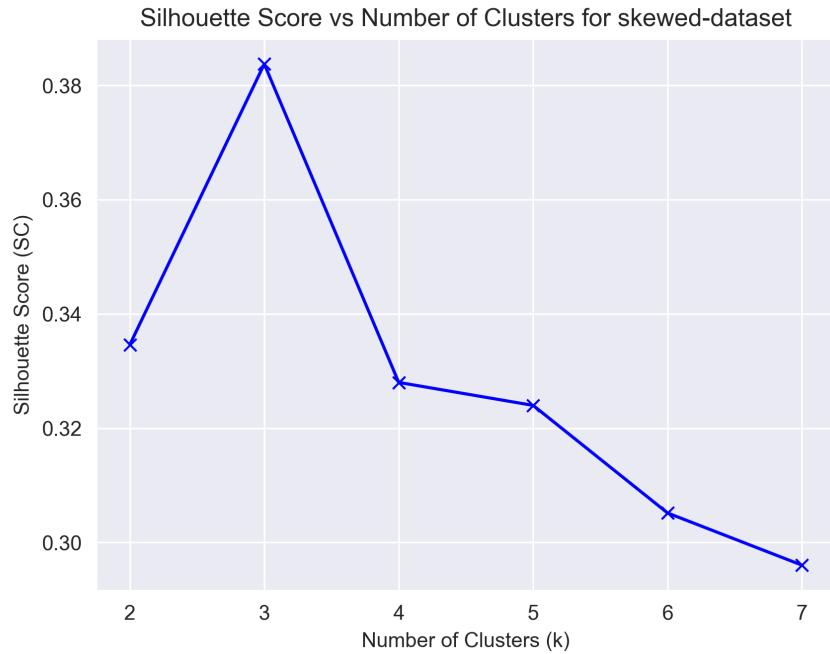


Figure 15: Selecting the  $k$  for Agglomerative clustering for skewed dataset (2-dimensions) using the "elbow plot."

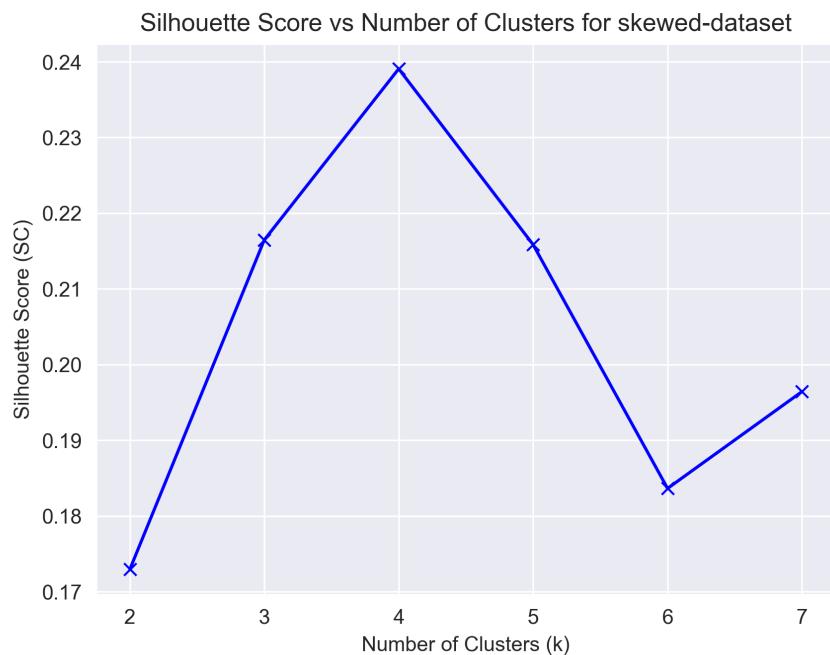


Figure 16: Selecting the  $k$  for Agglomerative clustering for skewed dataset (3-dimensions) using the "elbow plot."

# THEORY

## 3. BIG O NOTATION

The big O notation is a common way to describe the complexity of an algorithm. It is used to describe the worst-case scenario of an algorithm. For example, if an algorithm has a complexity of  $O(n)$ , it means that the algorithm will take at most  $n$  steps to complete (source for figure: [https://en.wikipedia.org/wiki/Big\\_O\\_notation](https://en.wikipedia.org/wiki/Big_O_notation)):

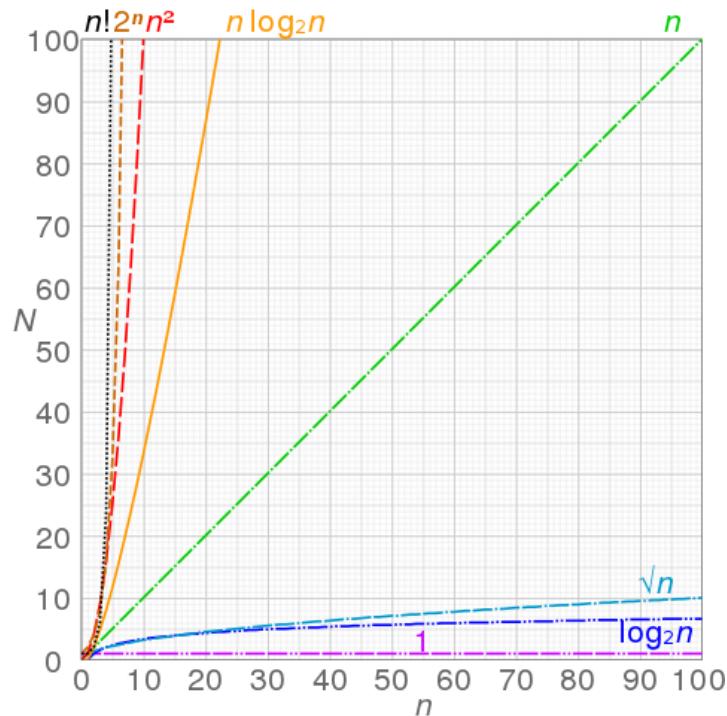


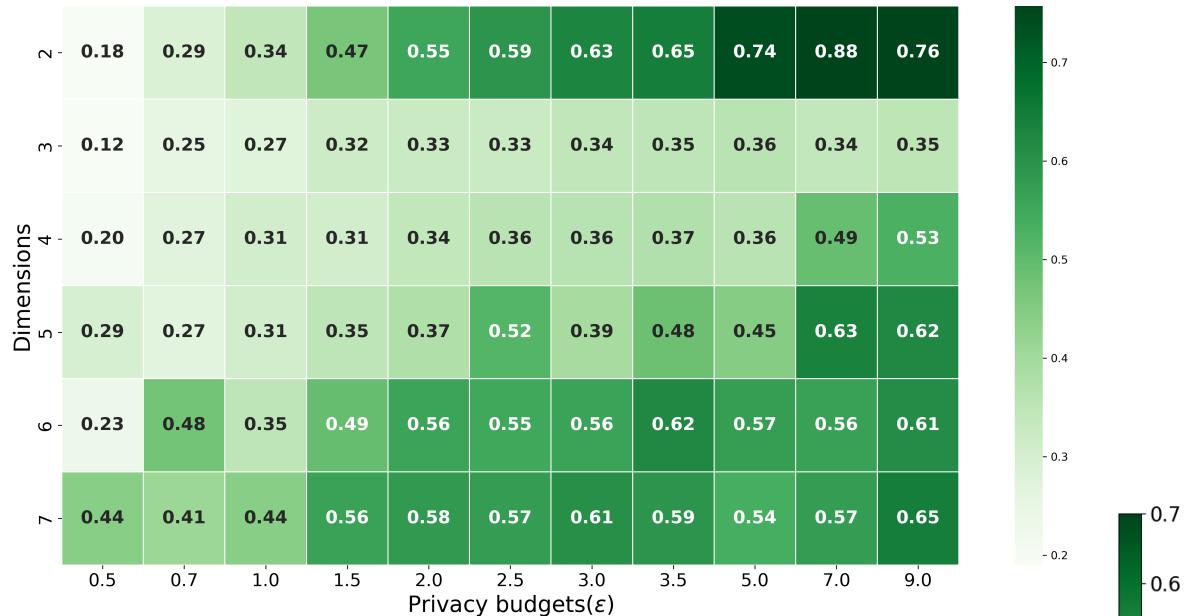
Figure 17: Graphical representation of the big O notation

# RESULTS

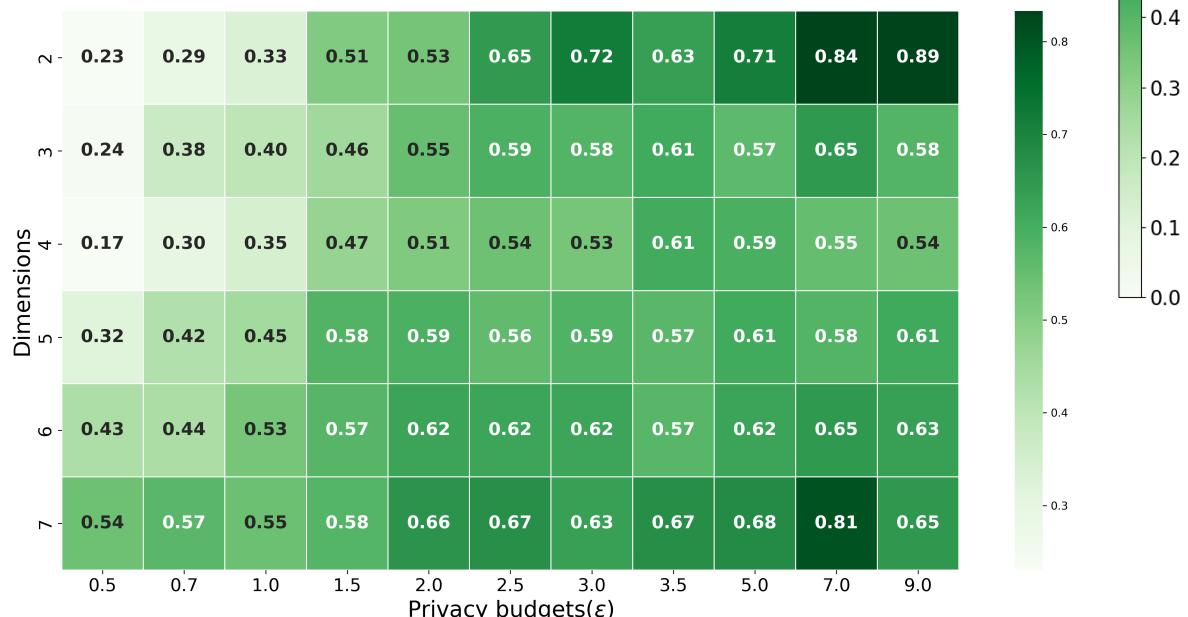
## .4. MECHANISM UTILITY

### .4.1. SEEDS DATASET

(a) Adjusted Mutual Information comparison for the grid-kD-Laplace mechanism



(b) Adjusted Mutual Information comparison for the density-kD-Laplace mechanism



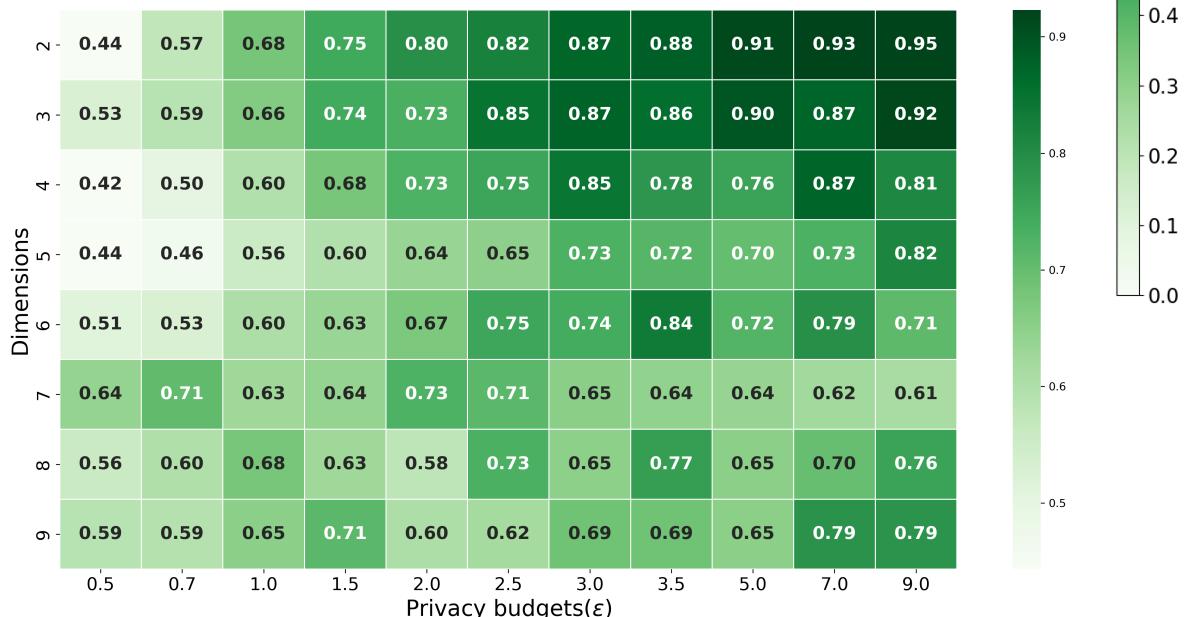
Add comments

## 4.2. HEART DATASET

(a) Adjusted Mutual Information comparison for the grid-kD-Laplace mechanism



(b) Adjusted Mutual Information comparison for the density-kD-Laplace mechanism



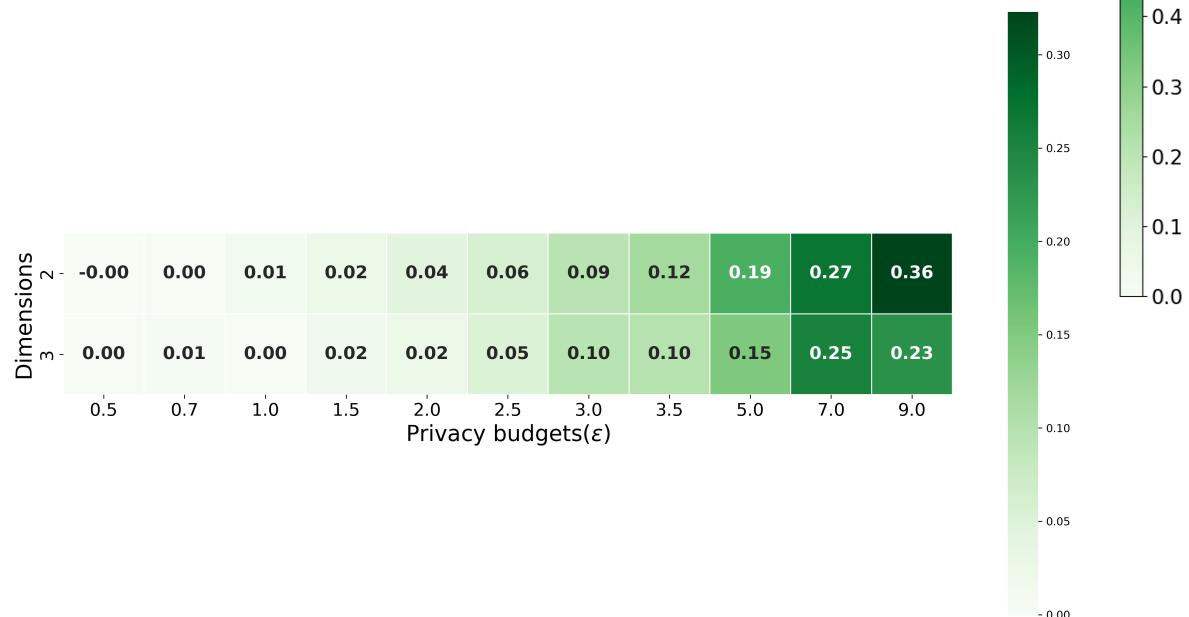
Add comments

### 4.3. CIRCLE DATASET

(a) Adjusted Mutual Information comparison for the grid-kD-Laplace mechanism



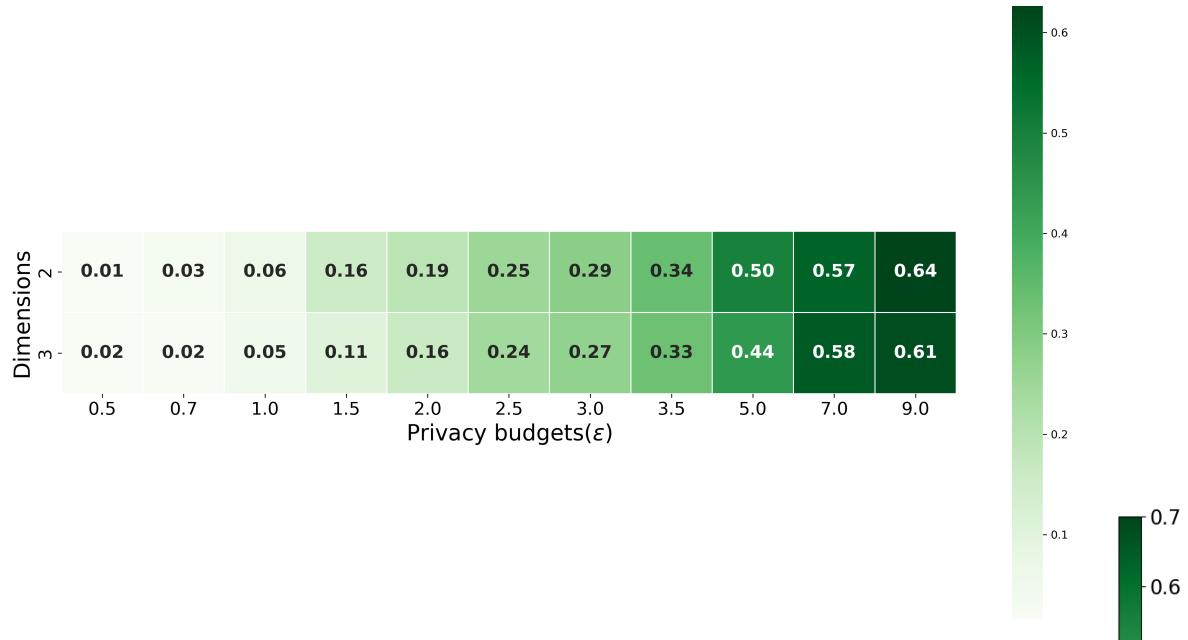
(b) Adjusted Mutual Information comparison for the density-kD-Laplace mechanism



Add comments

#### 4.4. LINE DATASET

(a) Adjusted Mutual Information comparison for the grid-kD-Laplace mechanism



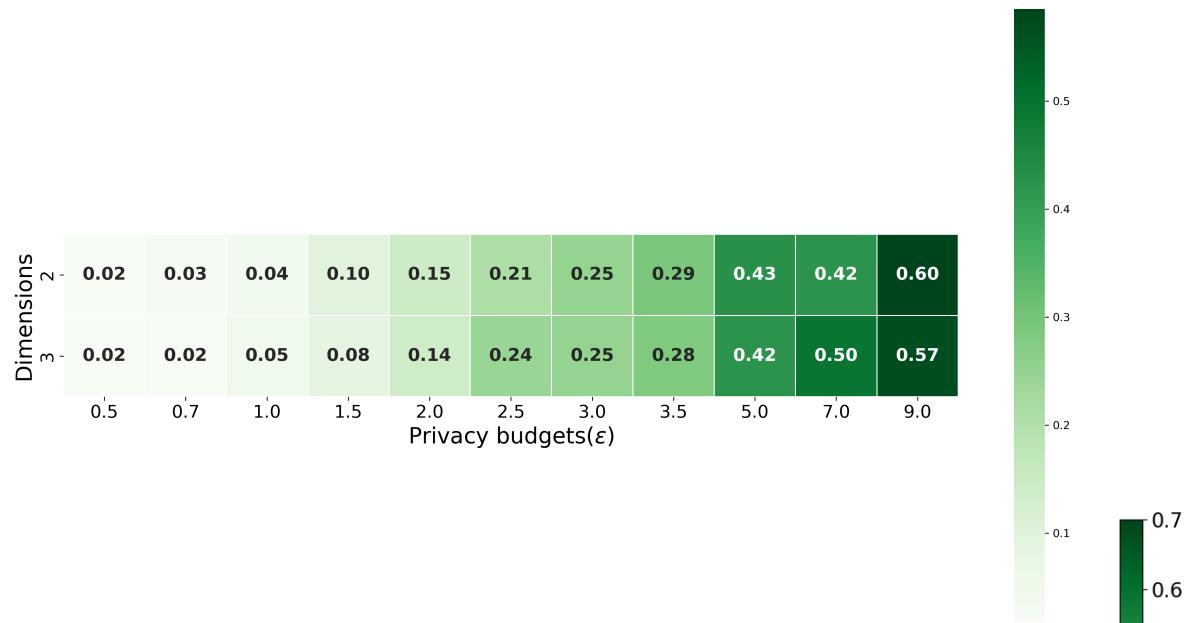
(b) Adjusted Mutual Information comparison for the density-kD-Laplace mechanism



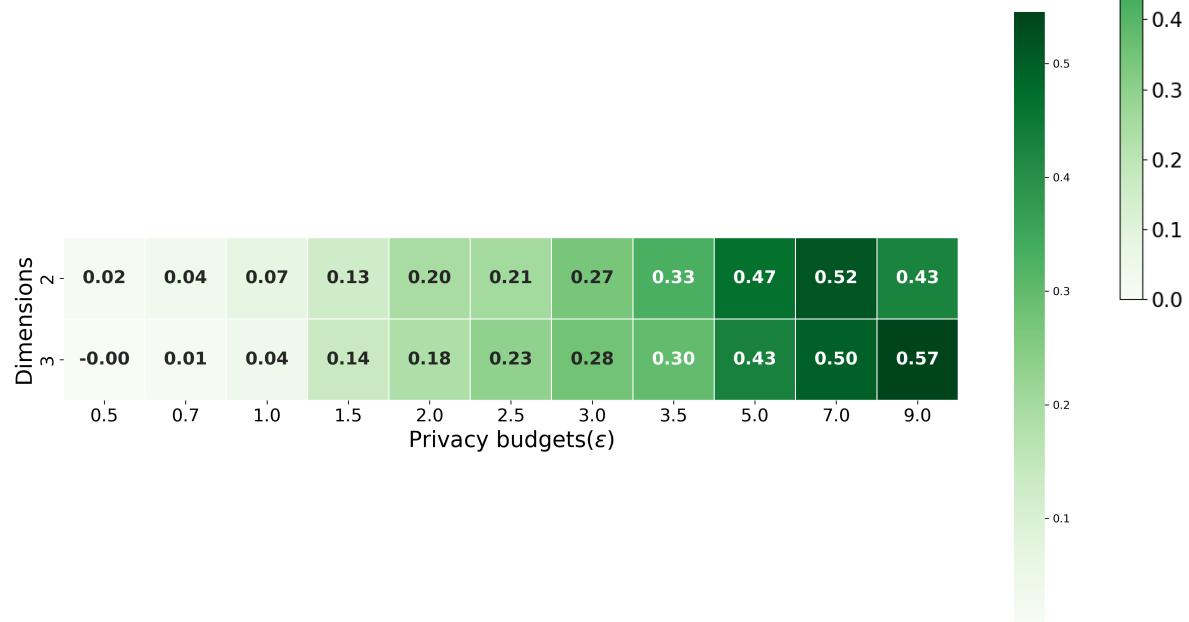
Add comments

#### 4.5. SKEWED DATASET

(a) Adjusted Mutual Information comparison for the grid-kD-Laplace mechanism



(b) Adjusted Mutual Information comparison for the density-kD-Laplace mechanism



Add comments

## .5. PRIVACY

### .5.1. SEEDS DATASET

(a) Heatmap showing disclosure risk for the kD-Laplace mechanism, per privacy budget & dimension for seeds-dataset.



(b) Heatmap showing disclosure risk for the Piecewise mechanism, per privacy budget & dimension for seeds-dataset.

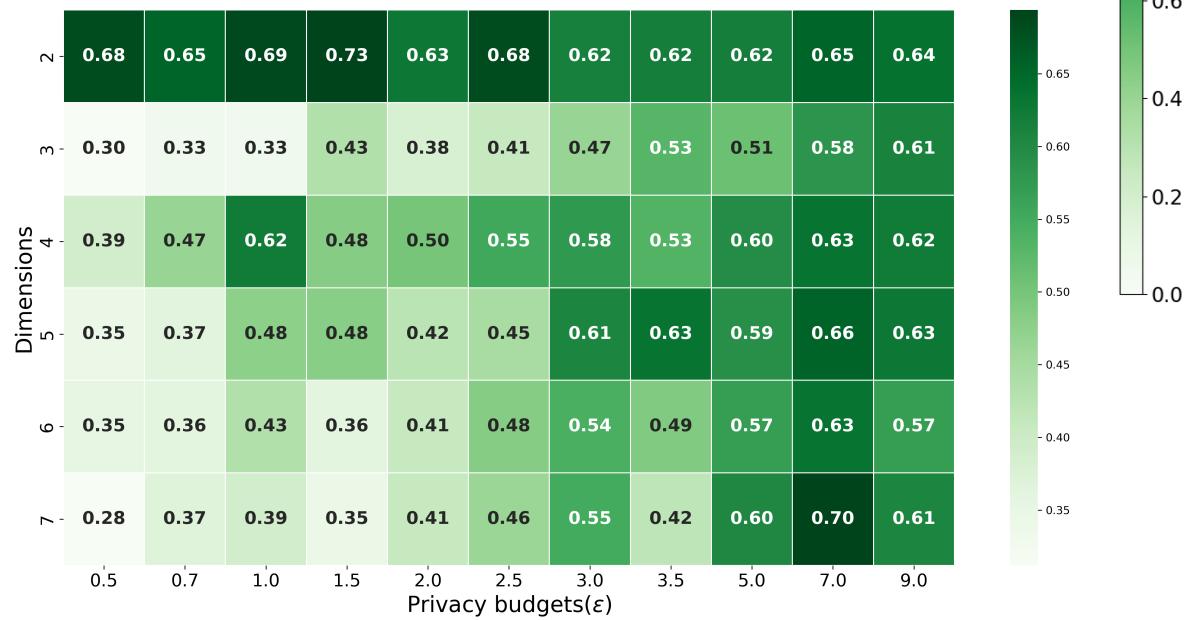


Add comments

(a) Heatmap TPR for the kD-Laplace mechanism, per privacy budget & dimension for seeds-dataset.



(b) Heatmap TPR for the Piecewise mechanism, per privacy budget & dimension for seeds-dataset.



Add comments

## 5.2. HEART DATASET

(a) Heatmap showing disclosure risk for the kD-Laplace mechanism, per privacy budget & dimension for heart-dataset.



(b) Heatmap showing disclosure risk for the Piecewise mechanism, per privacy budget & dimension for heart-dataset.

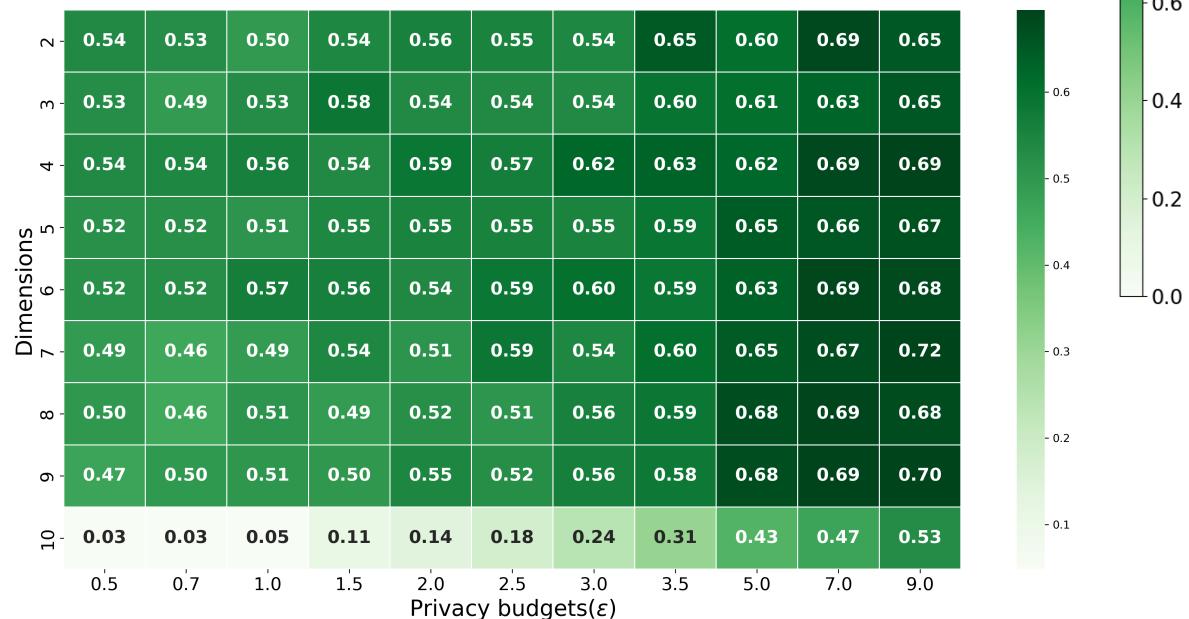


Add comments

(a) Heatmap TPR for the kD-Laplace mechanism, per privacy budget & dimension for heart-dataset.



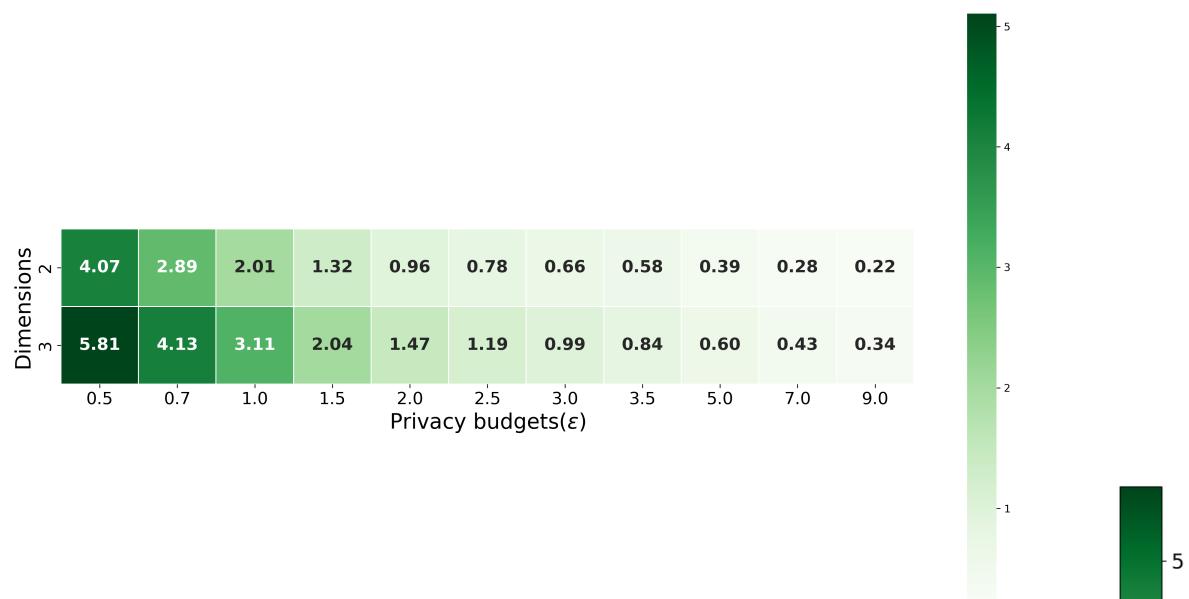
(b) Heatmap TPR for the Piecewise mechanism, per privacy budget & dimension for heart-dataset.



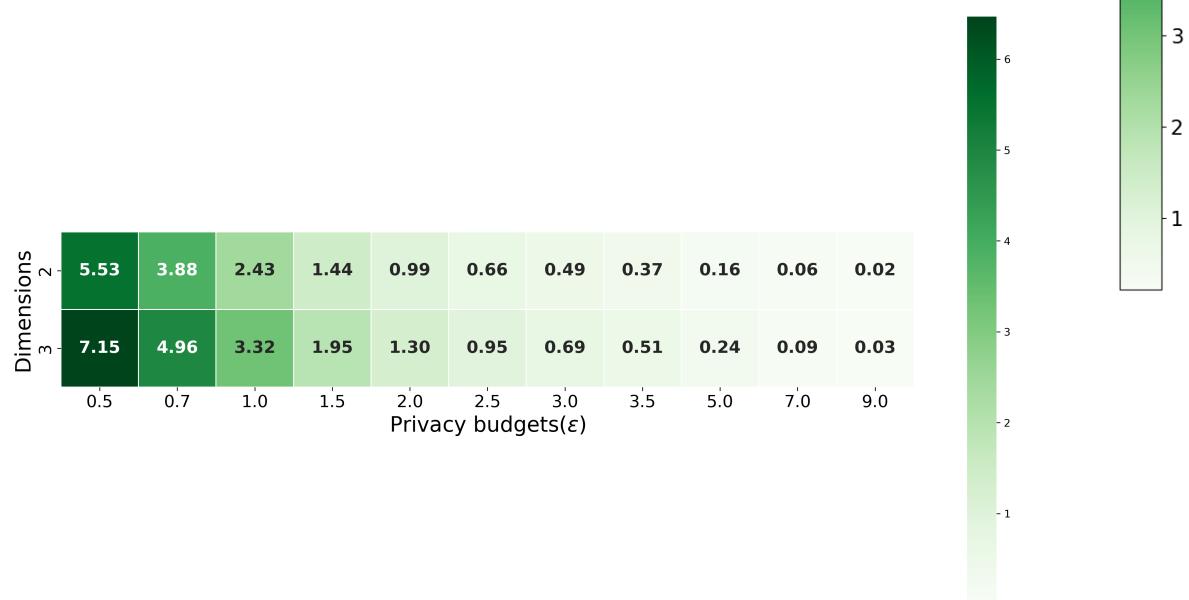
Add comments

### **.5.3. CIRCLE DATASET**

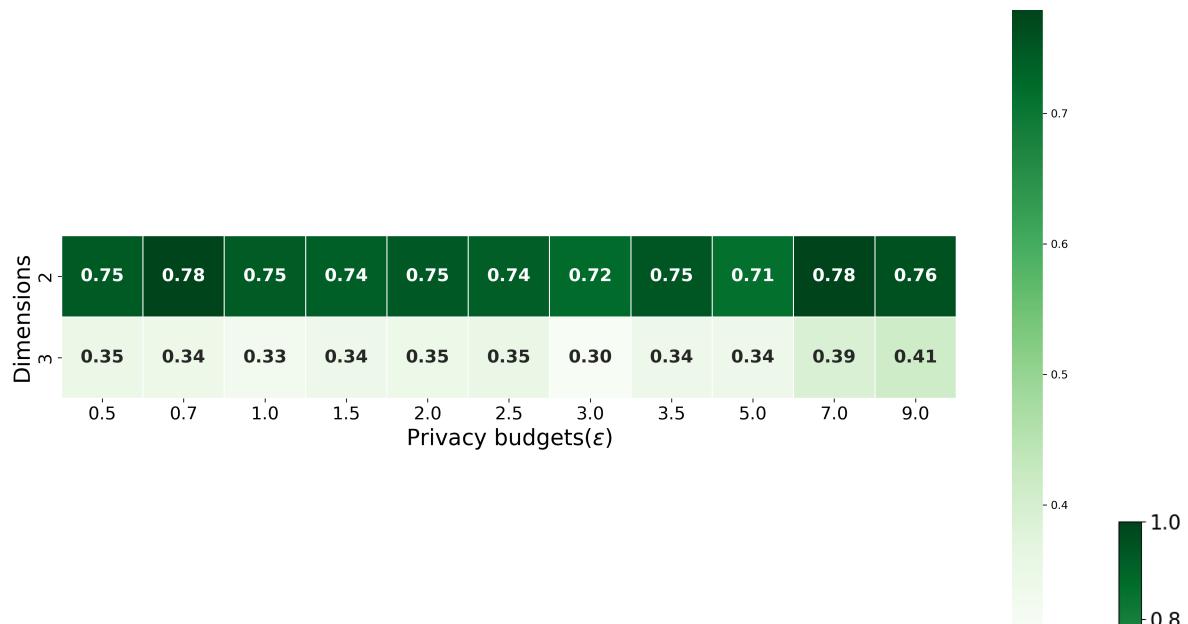
(a) Heatmap showing disclosure risk for the kD-Laplace mechanism, per privacy budget & dimension for circle-dataset.



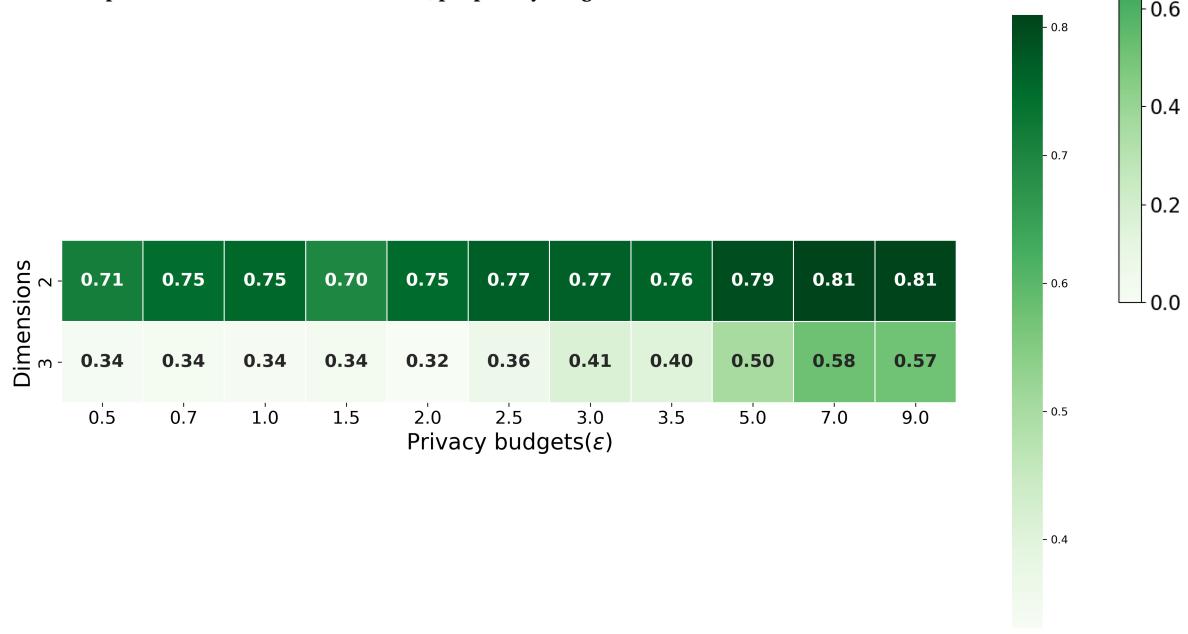
(b) Heatmap showing disclosure risk for the Piecewise mechanism, per privacy budget & dimension for circle-dataset.



(a) Heatmap TPR for the kD-Laplace mechanism, per privacy budget & dimension for Circle-dataset.



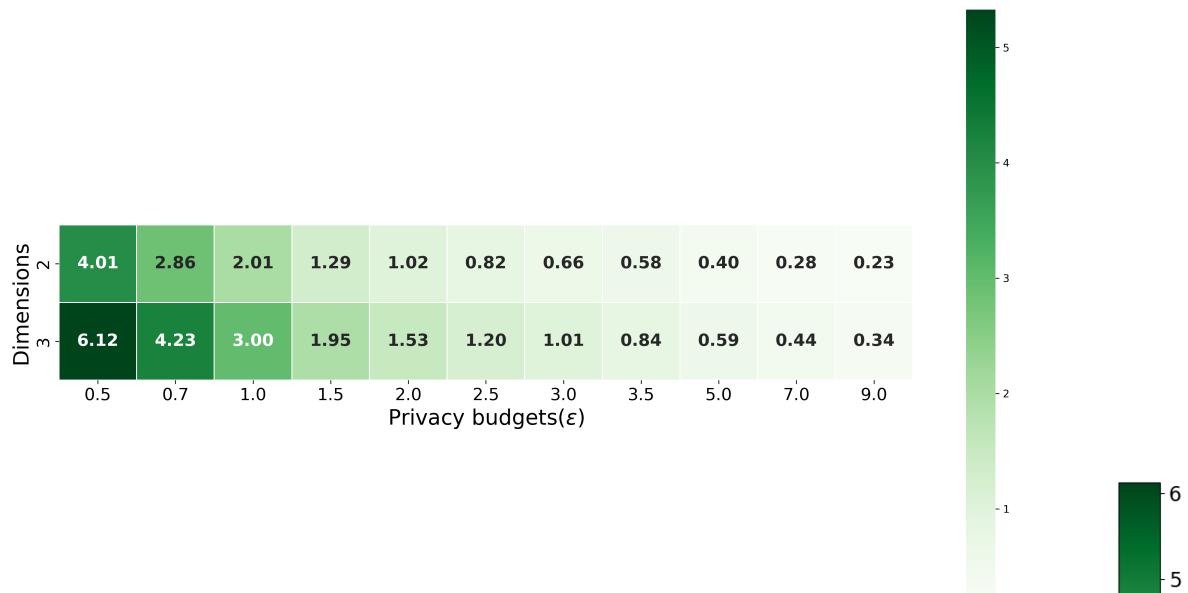
(b) Heatmap TPR for the Piecewise mechanism, per privacy budget & dimension for Circle-dataset.



Add comments

#### **.5.4. LINE DATASET**

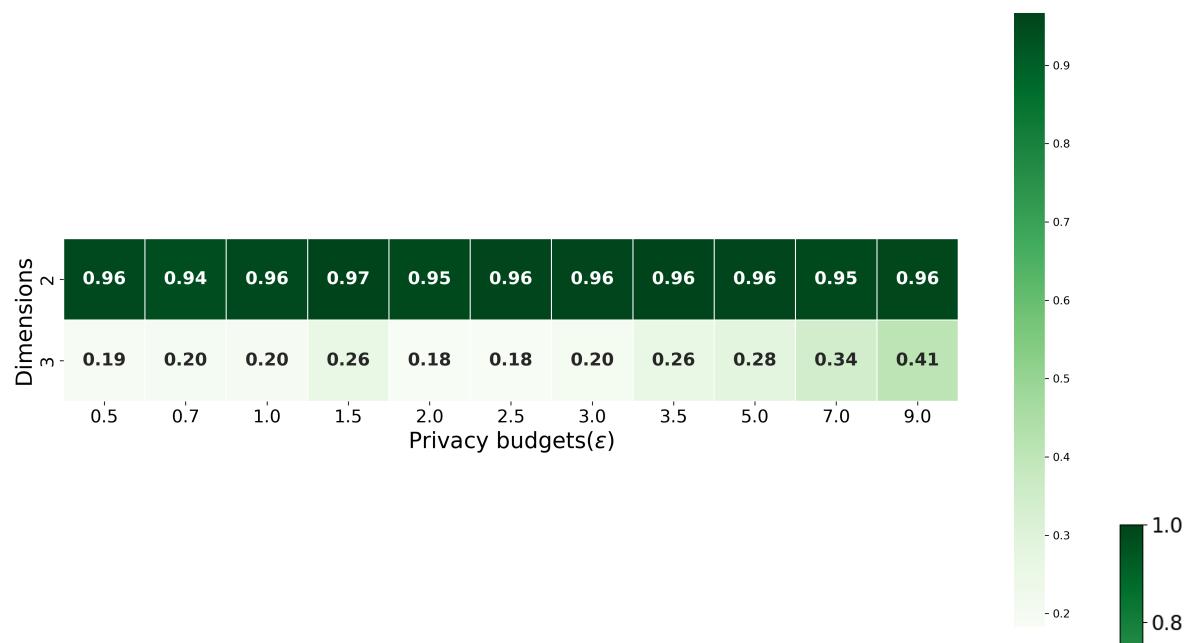
(a) Heatmap showing disclosure risk for the kD-Laplace mechanism, per privacy budget & dimension for line-dataset.



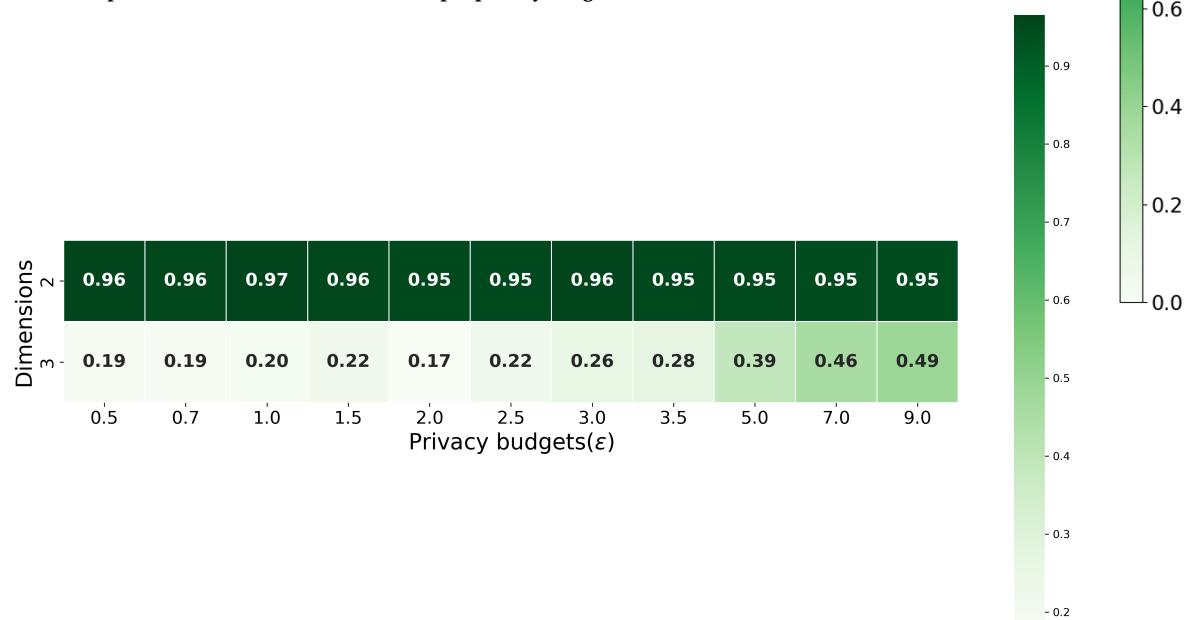
(b) Heatmap showing disclosure risk for the Piecewise mechanism, per privacy budget & dimension for line-dataset.



(a) Heatmap TPR for the kD-Laplace mechanism, per privacy budget & dimension for line-dataset.



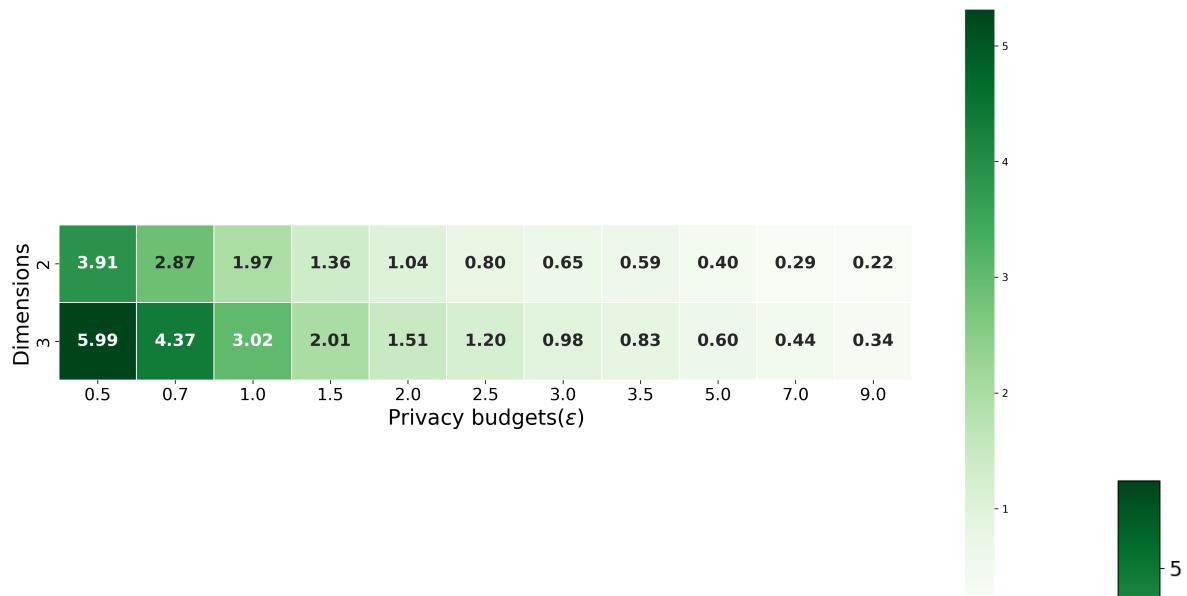
(b) Heatmap TPR for the Piecewise mechanism, per privacy budget & dimension for line-dataset.



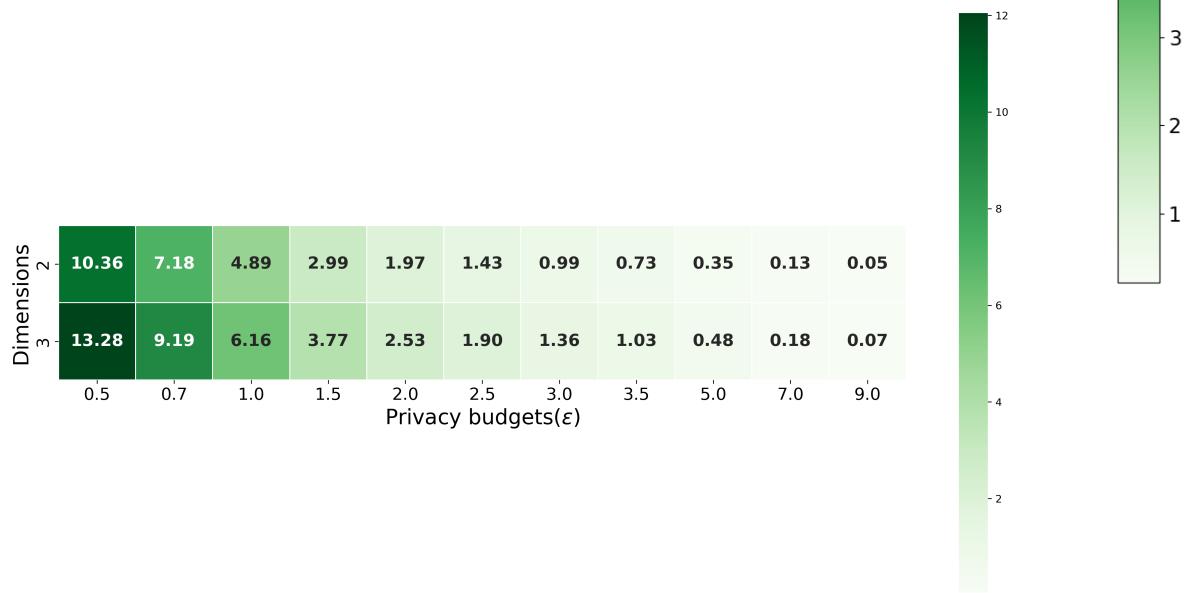
Add comments

### **.5.5. SKEWED DATASET**

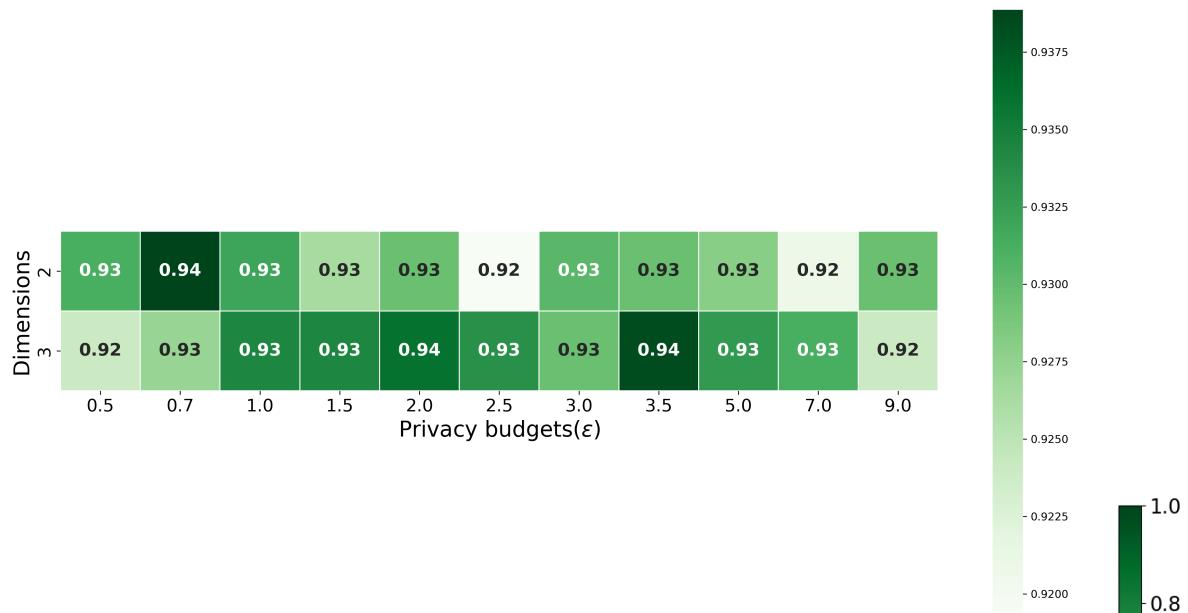
(a) Heatmap showing disclosure risk for the kD-Laplace mechanism, per privacy budget & dimension for skewed-dataset.



(b) Heatmap showing disclosure risk for the Piecewise mechanism, per privacy budget & dimension for skewed-dataset.



(a) Heatmap TPR for the kD-Laplace mechanism, per privacy budget & dimension for skewed-dataset.



(b) Heatmap TPR for the Piecewise mechanism, per privacy budget & dimension for skewed-dataset.



Add comments