

# An Algebraic Topological Perspective to Privacy: Numerical and Categorical Data\*

Alberto Speranzon

Shaunak D. Bopardikar<sup>†</sup>

June 20, 2018

## Abstract

In this paper, we cast the classic problem of achieving  $k$ -anonymity for a given database as a problem in algebraic topology. Using techniques from this field of mathematics, we propose a framework for  $k$ -anonymity that brings new insights and algorithms to anonymize a database. We begin by addressing the simpler case when the data lies in a metric space. This case is instrumental to introduce the main ideas and notation. Specifically, by mapping a database to the Euclidean space and by considering the distance between datapoints, we introduce a simplicial representation of the data and show how concepts from algebraic topology, such as the nerve complex and persistent homology, can be applied to efficiently obtain the entire spectrum of  $k$ -anonymity of the database for various values of  $k$  and levels of generalization. For this representation, we provide an analytic characterization of conditions under which a given representation of the dataset is  $k$ -anonymous. We introduce a weighted barcode diagram which, in this context, becomes a computational tool to tradeoff data anonymity with data loss expressed as level of generalization. Some simulations results are used to illustrate the main idea of the paper. We conclude the paper with a discussion on how to extend this method to address the general case of a mix of categorical and metric data.

## 1 Introduction

Recent times have seen a revolution in computing technologies. Third-party computing services, such as the Cloud, have been creating new paradigms for both, data storage and computation. Such technologies require one to repeatedly revisit the basic question of “How do we protect the data from a privacy perspective?”. Although this problem originated in database theory and computer science applications, it has recently expanded to domains such as systems and control. In control systems, there is always an information flow between sensors/actuators and controllers/plants and even between controllers in the case of distributed or cloud-based control. In critical cyber-physical systems such as power or transportation networks, where information about users and utility companies are exchanged or information about multiple users is fused, privacy has become a critical element to be considered [1]. Clearly, as the world is becoming more connected and loops are increasingly being closed over millions of sensors [2], privacy is becoming a top priority in the control community.

In this paper, we consider static data collected within a database that, in the context of cyber-physical systems, could represent log/monitoring data that needs to be analyzed offline to determine

---

\*This work was supported by United Technologies Research Center.

<sup>†</sup>Alberto Speranzon was with United Technologies Research Center at the time when this work was developed. He is now with Honeywell Aerospace – Advanced Technology, [alberto.speranzon@honeywell.com](mailto:alberto.speranzon@honeywell.com). Shaunak D. Bopardikar is with United Technologies Research Center, Inc. [bopardsd@utrc.utc.com](mailto:bopardsd@utrc.utc.com).

overall system performance, enterprise level fault detection and propagation, forensic analysis, etc. Although several approaches have been proposed to address different aspects of the privatization of data within a database, this paper provides a novel framework and perspective to address the most classical version of the problem using concepts and tools from algebraic topology.

## 1.1 Literature

Among several methods that have been developed for database privacy, a classic and popular approach is *k-anonymity*, which is a mechanism for protecting privacy of individuals represented as entries in a database [3]. The idea consists of removing all attributes from a database that can be used as unique identifiers but retain a set of attributes, called *quasi-identifiers*, for which identification may be possible but these are necessary for analysis. To such quasi-identifiers – ZIP code, date of birth, GPS location, energy usage, data/time, etc. – one applies a transformation that generalizes their value so that data records become indistinguishable. More precisely, for a given value of  $k$ , the original database is modified so that at least  $k$  individuals in the database have identical quasi-identifiers. This is achieved by generalizing numeric or text attributes: for example, the ZIP code can be generalized so that a certain number of the least significant digits are suppressed as  $46532 \rightarrow 465**$ , age could be generalized to intervals,  $35 \rightarrow [30, 40]$ , and the gender could be generalized as  $\{M, F\} \rightarrow Person$ .

The problem of computing an optimal  $k$ -anonymous version of a database has been shown to be NP-hard [4]. However, efficient algorithms such as Incognito [5] and its variants or greedy clustering-based algorithms [6] have been proposed to achieve  $k$ -anonymity. A multi-dimensional extension of the greedy approaches has been addressed in [7], which results into a representation of the database that is reminiscent of classic grid-based paintings by Mondrian. In multi-dimensional settings, a data aggregation scheme based on Hilbert curves has been proposed in [8].

Algebraic topology is a branch of mathematics that leverages tools and concepts from abstract algebra to study topological spaces. For example, a simple model for sensor network is a set of points in a (multidimensional) space in which two points (or sensors) are *neighbours* if they are within a specified distance of each other. Then, concepts from algebraic topology, such as homology, have been used to detect *holes* in sensor networks [9, 10]. Distributed algorithms to localize holes in sensor networks using related concepts have been addressed in [11] and in [12]. Recently such methods have been also used for filtering and position estimation in [13, 14].

One concept of privacy which has been applied to several control problems is that of differential privacy [15]. Informally, this concept means that for a given database, if any single individual is removed from the database, then no output of a computation run on the database would become significantly more or less likely. This concept has been applied to achieve differential privacy of Kalman filtering and estimation problems [16], to ensure a level of truthfulness in electric vehicle charging applications [17], to achieve average consensus in a private manner [18], to name a few. While the concept of differential privacy is very general and is applicable to dynamic databases, the resulting mechanism relies very strongly on the type of function/query that needs to be computed on the database. In contrast,  $k$ -anonymity is a static concept but is independent of any computation to be carried out on the database and therefore suitable in the context of offline analysis. That said, there are situations when  $k$ -anonymity is not sufficient and individuals can be re-identified despite anonymization. Approaches to deal with such cases have been considered such as  $\ell$ -diversity [19] and  $t$ -closeness [20]. Even though such advanced concepts have been proposed and privacy metrics continue to be an active research area,  $k$ -anonymity is still widely used.

Extensions of the proposed approach to other metrics will be a subject of future study.

## 1.2 Contributions

This paper introduces a novel perspective to data privacy based on algebraic topology. In particular, we address the case when the data lies in a metric space. By defining two datapoints that lie within a specified radius as neighbours, we show that the representation falls within the natural setting of a Čech (or in general, a Nerve) complex [21]. By increasing the radius (generalization), we show that the sequence of Čech complexes result into a *filtration*, i.e. nested complexes. This further implies that tools such as persistent homology can be applied to efficiently obtain the entire spectrum of  $k$ -anonymity of the database for various values of the generalization. The benefit of this approach is that once the family of complexes is built, for various generalization values, we can apply fast and scalable persistent homology algorithms, such as Perseus [22] to determine the tradeoffs. Furthermore, the persistent diagram not only provides the tradeoffs between a generalization and the value of  $k$ , but also show how many equivalent classes are formed for a given generalization achieving a certain  $k$ -anonymity, a metric that has an impact on the anonymized data quality [23].

For this representation, we provide an analytic characterization of conditions under which a given representation of the dataset is  $k$ -anonymous. Finally, we discuss how this method can be extended to address the general case of a mix of categorical and metric data.

## 1.3 Organization of this paper

This paper is organized as follows. The problem formulation is presented in Section 2. Background results and concepts from algebraic topology are reviewed in Section 3. The proposed approach is presented in Section 4 for numerical data along with some simulation results. Extension of the method described in Section 4 to mixed categorical and numerical data is discussed in Section 5.

## 2 Problem formulation

Let us consider a database table  $T(A_1, A_2, \dots, A_m)$  consisting of  $N$  rows, where each  $A_i \in \mathcal{D}$  are various attributes that in general can take the form of *numeric* and/or *categorical* values, i.e., the domain  $\mathcal{D}$  can either be a set of discrete or continuous values. Without loss of generality, we can identify with  $Q_T = \{A_1, A_2, \dots, A_d\}$  a set of  $d$  attributes that we define as *quasi-identifiers*, namely attributes that can be joined with external information/databases so that private information can be obtained. Typical examples of private data that could be obtained are names of individuals, salaries, etc.

Another database table  $\bar{T}(\bar{A}_1, \bar{A}_2, \dots, \bar{A}_m)$  consisting of  $N$  rows is said to be a *generalization*  $\bar{T} = G(T)$  of the table  $T$  if, for every row  $T_j$  of  $T$ ,

$$Q_{T_j} \subset Q_{\bar{T}_j}.$$

In this paper, as said previously, we will be focusing on the concept of  $k$ -anonymity for privacy, that is formally defined as follows.

**Definition 2.1 ( $k$ -anonymity [3])** Consider a generalized database  $\bar{T}$  and a quasi-identifier set  $Q_{\bar{T}}$ . The set  $Q_{\bar{T}}$  is said to have the  $k$ -anonymity property if and only if each unique tuple in the projection of  $\bar{T}$  on  $Q_{\bar{T}}$  occurs at least  $k$  times in  $\bar{T}$ .

Given a database  $T$ , the problem of  $k$ -anonymity is thus to determine a generalization function  $G(\cdot)$  so that the resulting database  $\bar{T} = G(T)$  is  $k$ -anonymous. Clearly, one may simply generalize every entry and find the smallest set that generalizes every row of  $T$ . However, this

trivial method would completely destroy the information content in the original database. The problem is how to minimize such as *over-generalization* of the quasi-identifiers. This notion will be made precise in Section 4.

### 3 Background on Topological Methods

In this section, we provide a summary of some useful concepts from algebraic topology that will be used in our approach. Interested readers may refer to [21] and [24] for additional details on these concepts.

**Definition 3.1 (Čech complex)** *Given a collection of points  $\{x_i\} \in \mathbb{R}^n$ , the Čech complex is the abstract simplicial complex whose  $k$ -simplices are determined by unordered  $(k+1)$ -tuples of points  $\{x_i\}_0^k$  whose closed  $\epsilon$ -ball neighborhoods have a point of common intersection.*

#### 3.1 Simplicial Homology

Homology is an algebraic characterization of “holes” in a topological space. The central notion is that of a boundary homomorphism, which in the context of simplicial complexes, encodes how simplices are attached to their lower dimensional facets. To define (simplicial) homology, of a complex  $C$ , we choose an ordering of each simplex, in the same way directed graphs are ordered. Given such ordering we consider  $\mathbb{R}$ -vector spaces  $\mathcal{C}_k(C)$  with basis the oriented  $k$ -simplices. We thus have that  $\mathcal{C}_\bullet$ , forms a sequence of vector spaces, which we call chain complex. A *boundary homomorphism* is defined as the linear map  $\partial_k : \mathcal{C}_k(C) \rightarrow \mathcal{C}_{k-1}(C)$  given by associating each basis element of  $\mathcal{C}_k(C)$  to the formal sum of its (oriented) faces of dimension  $k-1$ . The boundary operator  $\partial = \{\partial_k\}$  thus encodes the assembly instructions of  $C$ . It turns out that the  $k^{\text{th}}$  homology group of the complex  $C$ ,  $H_k(C)$  is given by

$$H_k(C) = Z_k/B_k = \ker \partial_k / \text{im } \partial_{k+1}.$$

The group  $Z_k = \ker \partial_k$  is called the  $k$ -th cycle group and its elements (chains) represent  $k$ -cycles. We have that  $B_k = \text{im } \partial_{k+1}$  is the  $k$ -th boundary group whose elements are  $k$ -boundaries. The quotient space  $H_k(C)$  thus represents all the  $k$ -cycles that are not boundaries of  $k+1$  simplices, namely cycles that represent  $k$ -dimensional “holes”. The homology of the complex  $C$  is then  $H_\bullet(C) = \{H_k(C)\}$ .

In this paper, we will use the notion of dimension of the  $k$ -th homology, that is the dimension of the vector space  $H_k(C)$ ,  $\dim H_k(C)$ . In particular, the  $k$ -th homology group  $H_k(C)$  is said to be *trivial* if  $\dim H_k(C) = 0$ .

#### 3.2 Persistent Homology

Let us consider a sequence of complexes  $C^\epsilon$  with  $\epsilon = \{\epsilon_1, \epsilon_2, \dots, \epsilon_M\}$ , more specifically the sequence of Čech complexes  $\{C^{\epsilon_i}\}_{i=1}^N$ , for increasing  $\epsilon_i \in \mathbb{R}_{\geq 0}$ ,  $\epsilon_i \leq \epsilon_j$  for  $i \neq j$ . There are clearly inclusion maps between such complexes

$$C^{\epsilon_1} \hookrightarrow C^{\epsilon_2} \hookrightarrow \dots \hookrightarrow C^{\epsilon_{M-1}} \hookrightarrow C^{\epsilon_M}.$$

Rather than studying the homology of each complex for each value of the parameter  $\epsilon$ , one can then study the homology of the inclusions  $\iota : H_\bullet(C^{\epsilon_i}) \rightarrow H_\bullet(C^{\epsilon_j})$  for  $i < j$ . Such maps are important as they capture topological features that persist over the parameter space.

The dimension of the homology groups as a function of the single parameter  $\epsilon$  can be plotted in a diagram, called the *barcodes diagram*, see [24]. We show a simple example in Figure 1 where  $\epsilon$  is

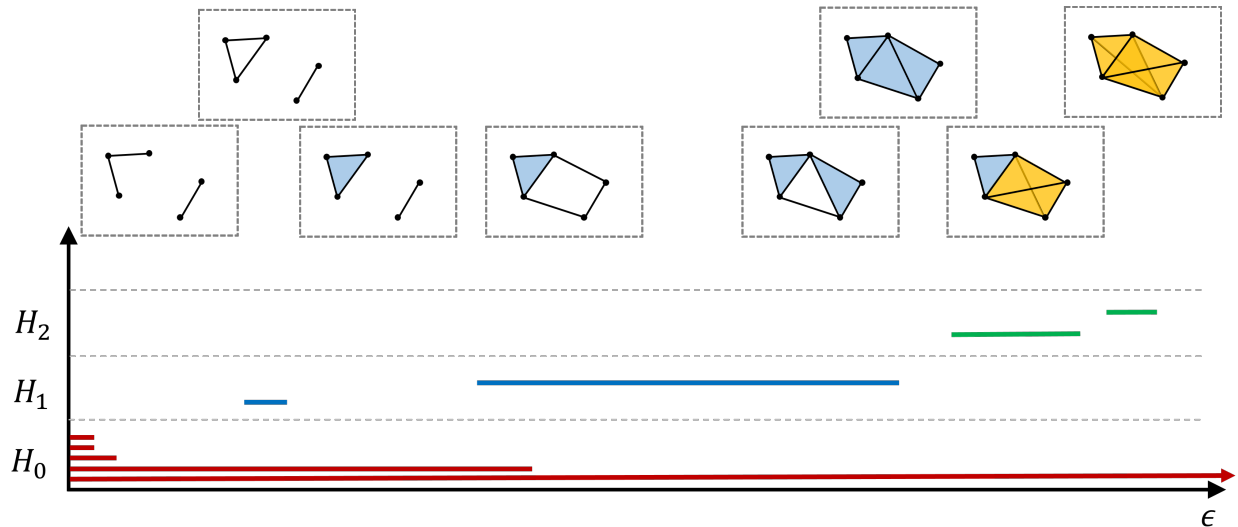


Figure 1: Example of the barcodes.

the radius of a ball around each vertex, and where there is a  $k$ -simplex whenever  $k + 1$  circles have non-empty intersections. For small values of  $\epsilon$ , we have that the number of connected components,  $\dim H_0$ , is the number of vertices (0-simplices), and as  $\epsilon$  increases, more vertices get connected, resulting in components to merge till a single connected component is obtained. For a small value of  $\epsilon$  there are no high dimensional holes, however, at some point, before the first 2-simplex (blue triangle) gets filled in, such 2-simplex is not filled and generates a hole that quickly disappears. At a later value of  $\epsilon$ , a large hole is formed at the time when the vertices form a single large connected component. As the  $\epsilon$  parameter increases further at some point the “middle” hole gets filled and the dimension of the first homology,  $\dim H_1$  becomes zero again. As  $\epsilon$  further increases, tetrahedrons appear, first with an empty volume, namely a void, that disappear. Higher dimensional holes will likely occur, but we did not depict them. As  $\epsilon$  further increases the complex will have trivial higher homology groups and only have a single connected component.

In this paper we will make use the persistent homology in the context of  $k$ -anonymity.

## 4 $k$ -Anonymity via Persistent Homology: Numerical Attributes

In this section, we describe the algebraic topological approach toward achieving  $k$ -anonymity. In this section, we will restrict our attention to the case of the attributes in the quasi-identifier set  $Q_T$  being all numeric/continuous variables, such as Age, Salaries, Taxes paid in a year, etc. In this case, we can clearly represent the set  $Q_T$  as a  $|Q_T|$ -dimensional vector. We will assume that all the vectors are elements of a real vector space.

Figure 2 shows an example of a table  $T$  with two identifiers (Name and Last Name), three quasi-identifiers  $Q_T = \{\text{Age, Total Scholarship, Money Borrowed}\}$  and the sensitive column **Current Salary**. As we show at the bottom of Figure 2, we can map the quasi-identifiers to a three-dimensional vector space (three dimensional as  $|Q_T| = 3$ ), where each entry in the table corresponds to a point in  $\mathbb{R}^3$ .

In the following, we will often refer to entries of the table  $T$  as *points* due to the previously described representation. Also, note that without any loss of generality, we can consider the data

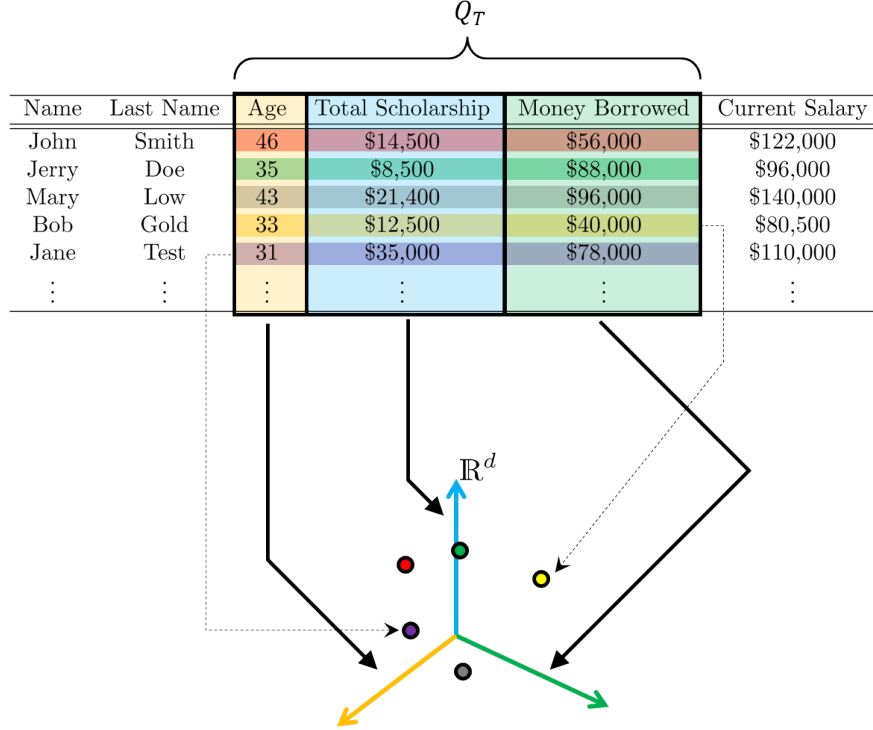


Figure 2: Example of a table where the first two columns are identifiers, the last column is the sensitive data and the middle three are the quasi-identifiers. Each entry in the table can be mapped, through the quasi-identifiers, to a point in the three-dimensional vector space,  $\mathbb{R}^3$ .

to take values within the hypercube  $\mathcal{M} = [0, 1]^{|Q_T|}$ , since one can always normalize the data accordingly.

The first definition is a direct application of the Čech complex as summarized in Section 3. In this paper, we use this structure to capture the  $k$ -anonymity property of the data.

**Definition 4.1 (Anonymity Complex)** *Given a table  $T$  with  $N$  rows and a set of quasi-identifier  $Q_T$ , let us consider the  $N$  points  $\{p_i\}_1^N \in \mathcal{M}^N$ . We define an anonymity complex  $\mathcal{C}(p)$  the simplicial complex whose  $k$ -simplices are determined by  $(k+1)$  points  $\{p_{i_0}, p_{i_1}, \dots, p_{i_k}\}_0^k$  such that closed  $\epsilon$ -ball neighborhoods centered around these points have at least one intersection point. We call the radius  $\epsilon$  the global generalization strategy.*

We now introduce an important building block.

**Definition 4.2 (Anonymity  $k$ -simplex)** *Given a global generalization  $\epsilon$ , we say that  $k$  points have the  $k$ -anonymity property if all the closed  $\epsilon$ -ball neighborhoods of the  $k$  points all intersect in at least a point. In this case, we have that the  $k$  points form a  $k$ -simplex, which we term as an anonymity  $k$ -simplex and denote with  $S_k$ .*

Figure 3(a) shows an example of an anonymity 4-simplex for a given global generalization  $\epsilon$  while Figure 3(b) shows an example where for the same value of  $\epsilon$ , 4-anonymity cannot be achieved.

The following is a useful test to determine whether  $k$ -anonymity can be achieved for a given value of  $\epsilon$  or not.

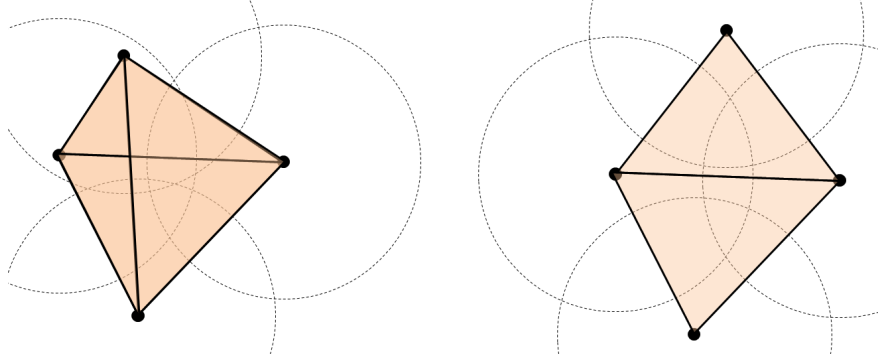


Figure 3: In (a), we show a anonymity 4-simplex representing the fact that there exists an  $\epsilon$  such that the 4 points can be anonymized. In (b), we show an example where there is no anonymity 4-simplex (indeed we have only two 2-simplices) as the selected  $\epsilon$  is not large enough, or equivalently the generalization is high enough, to ensure 4-anonymity.

**Lemma 4.1** *Given a set of points  $p = \{p_i\}_1^N$  corresponding to  $N$  rows of the table  $T$ , given a global generalization  $\epsilon$  and the corresponding anonymity complex  $C(p)$ , we have that the points  $\{p_i\}_1^N$  have the  $k$ -anonymity property if and only if*

$$C(p) = \bigcup_i S_{\ell_i}$$

where  $S_{\ell_i}$  is the  $i$ -th anonymity  $\ell_i$ -simplex with  $\ell_i \geq k$  for  $i \in \mathbb{N}_0$ . We say in this case that the anonymity complex achieves  $k$ -anonymity.

*Proof:* The points  $\{p_i\}_1^N$  have the  $k$ -anonymity property if and only if we can sub-divide the set of points into subsets such that

$$\{p_i\}_1^N = \underbrace{\{p_i\}_1^{i_1}}_{\Sigma_{\ell_1}} \cup \underbrace{\{p_i\}_{i_1+1}^{i_2}}_{\Sigma_{\ell_2}} \cup \dots \cup \underbrace{\{p_i\}_{i_{N-1}+1}^N}_{\Sigma_{\ell_N}},$$

and to each subset  $\Sigma_{\ell_i}(p)$ , we can associate an anonymity  $\ell_i$ -simplex  $S_{\ell_i}$  with  $\ell_i \geq k$  for any  $i$  (for given fixed  $\epsilon$ ) and  $S_{\ell_i} \cap S_{\ell_j} = \emptyset$  for  $i \neq j$ .

Given the previous set relations, we have that complex  $C$  associated with the set of points  $\{p_i\}$  is then given by the union of the  $S_{\ell_i}$ .  $\square$

This result then establishes a natural connection between the properties of the anonymity complex, in terms of some of its subcomplexes, and the  $k$ -anonymity property.

We further explore how topological properties of the anonymity complex are related to the  $k$ -anonymity property and how we can leverage that to find an “optimal” generalization. Let us first establish some topological properties of  $C$  and then define explicitly what we mean with “optimal” generalization.

**Proposition 4.1** *An anonymity complex  $C$ , for a given  $\epsilon$ , has the  $k$ -anonymity property if and only if its homology groups  $H_n(C)$  are trivial for any  $n > 0$ , and every connected components is an  $\ell$ -simplex with  $\ell \geq k$ . Furthermore, when this is the case the number of equivalence classes generated by using the  $\epsilon$  generalization is given by the dimension of the zero-th homology,  $\dim H_0(C)$ .*

---

**Algorithm 1**  $k$ -anonymity via Persistence Homology

---

**Inputs:**  $p = \{p_i\}_1^N \in \mathbb{R}^{d \times N}$ , Parameter:  $k$ , Radii:  $\{\epsilon_1, \dots, \epsilon_M\}$   
Re-scale the dataset into the unit cube  $[0, 1]^d \subset \mathbb{R}^{d \times N}$ .  
**for** every value of  $\epsilon \in \{\epsilon_1, \dots, \epsilon_M\}$  **do**  
    Construct the anonymity complex  $C^\epsilon(p)$   
**end for**  
Compute *weighted* persistent homology  
**return** Complete bar code diagram

---

*Proof:* (If) From Lemma 4.1, we know that we can decompose  $C$  into a finite number of disjoint anonymity  $k$ -simplices. It is known [21] that in this case  $H_n(C) = \bigoplus_i H_n(S_{\ell_i})$ , namely the  $n$ -th homology of  $C$  is given by the direct sum of the  $n$ -th homology of the anonymity  $k$ -simplices. As the  $k$ -simplices are simply connected spaces and contractible, they have trivial high order ( $n > 0$ ) homology groups and  $H_0(S_{\ell_i}) \approx \mathbb{Z}$  for every  $i$ .

(Only if) Let us assume that  $C$  has  $H_n(C) = \{0\}$  for  $n > 0$ , and  $H_0(C)$  is non-trivial, and in particular let us assume that  $H_0(C) \approx \mathbb{R}^s$ . This means that the complex  $C$  has  $s$  connected components. From the hypothesis that the connected components are  $\ell$ -simplices with  $\ell \geq k$ , we know that each component is a anonymity simplex. Thus the anonymity complex  $C$  has the  $k$ -anonymity property.  $\square$

We are now able to connect topological properties of the anonymity complex with the  $k$ -anonymity property. Of course, as it can be seen from Proposition 4.1, the result still depends on  $\epsilon$ , namely the generalization.

When anonymizing a dataset, one is typically interested in “corrupting” the data by the least amount. Indeed, if one carefully thinks about the  $k$ -anonymity problem, it is always possible to find a large enough  $k$  that makes the data anonymous, i.e., if one makes the extreme choice of  $k = n$ , then the entire data will be in one equivalence class and thus,  $k$ -anonymity will be achieved. The issue with this is that the information contained in the data will be completely lost.

In the context of this paper, as the generalization is parametrized by  $\epsilon$ , we are interested for find the smallest value of  $\epsilon$  that gives  $k$ -anonymity. We then have the following definition.

**Definition 4.3 (Minimal Anonymity Complex)** *Given an anonymity complex  $C(p)$  associated to a set of points, let us denote with  $C^\epsilon$  the anonymity complex for a given generalization  $\epsilon$ .*

*We define as minimal anonymity complex the following object*

$$C^{\epsilon^*} = \min_{\epsilon} C^\epsilon,$$

*such that  $C^{\epsilon^*}$  achieves  $k$ -anonymity.*

Even without minimization over  $\epsilon$ , the  $k$ -anonymity problem known to be an NP-hard problem, and so it is clear that we cannot easily find the minimal anonymity complex. To find an approximate solution to the problem, we instead study the *persistent homology* of  $C^\epsilon$ . In particular, in this paper, we adapt the idea of persistent homology as a tool to provide the full spectrum of  $k$ -anonymization one can obtain. Our approach is summarized in Algorithm 1.

Note that the parameter  $\epsilon$  induces a family of complexes such that  $C^{\epsilon_1} \xrightarrow{\hookrightarrow} C^{\epsilon_2} \xrightarrow{\hookrightarrow} \dots \xrightarrow{\hookrightarrow} C^{\epsilon_M}$ , where  $\epsilon_i \leq \epsilon_j$ , for any  $i < j$ , and thus we recover the same setting as in the persistent homology. Formally we have a  $\epsilon$ -based filtration [24]. The idea here is to leverage the barcodes or persistent



Age	ZIP Code	Salary	Age	ZIP Code
25	47677	\$47,000	[22-25]	[47602-47678]
22	47602	\$32,000	[22-25]	[47602-47678]
24	47678	\$52,000	[22-25]	[47602-47678]
43	47905	\$151,000	[38-52]	[47905-47909]
52	47909	\$145,000	[38-52]	[47905-47909]
38	47906	\$98,000	[38-52]	[47905-47909]
47	47605	\$110,000	[32-47]	[47605-47603]
36	47673	\$92,000	[32-47]	[47605-47603]
32	47607	\$115,000	[32-47]	[47605-47603]

(a) Example data set with  $Q_T = \{\text{Age, ZIP Code}\}$ .

(b) 3-anonymized table

Table 1: Sample data set for illustrative purpose.

diagram to extract regimes of interests, namely anonymity strategies – values of  $\epsilon$  – that lead to  $k$ -anonymity for different values of  $k$ .

Such a persistent diagram has two specific features:

- Each bar in  $H_0$  diagram has a *weight* corresponding to the number of elements in the connected components;
- Given a value  $k$ , we only consider bars that have at least  $k$  elements.

To illustrate the proposed approach, consider the sample dataset shown in part (a) of Table 1. The quasi-identifier set  $Q_T = \{\text{Age, ZIP Code}\}$ , while the salary field is the sensitive one. Figure 4 shows an example of the output of the barcode diagram. For the case of 3-anonymity in particular, the lower center figure depicts  $H_0$  while the one in top center shows the  $H_1$ . We can see that a *hole* gets created for the values of the radius approximately in  $[0.11, 0.13]$  and then it gets *filled* thereby making the dataset 3-anonymous. But again in the interval  $[0.17, 0.19]$ , a hole gets created around when the complex  $[1, 2, 3, 7, 8, 9]$  gets formed. After  $\epsilon$  increases more, this hole disappears and we are left with an anonymity 3-simplex  $[4, 5, 6]$  and an anonymity 6-simplex  $[1, 2, 3, 7, 8, 9]$ .

We can clearly see three regimes where 3-anonymity is possible. The first one (as indicated) has three classes with three elements in each. The second regime corresponds to the values of radius in the interval  $[0.19, 0.4]$  which has only two equivalence classes, and the third one (radius greater than 0.4) being the trivial solution where there is only one class with nine elements.

For minimal data quality loss, the interval  $[0.17, 0.19]$  is the best solution as 3-anonymity can be reached with largest number of classes. Mapping the first one back to the dataset yields a 3-anonymous version of the dataset shown in part (b) of Table 1.

In Figure 4 shows that 2-anonymity can be achieved for generalizations  $\epsilon > 0.08$ . For generalizations with  $\epsilon \leq 0.08$ , we have only one anonymity 2-simplex and the rest of the data will be just points, thus 2-anonymity cannot be achieved. The rightmost plot shows that 4-anonymity can be reached for  $\epsilon > 0.4$ . Note that even if there is an anonymity 6-simplex in the interval  $[0.19, 0.4]$ , there is no way for the  $[4, 5, 6]$  complex to achieve 4-anonymity and thus the generalizations in the interval  $[0.19, 0.4]$  will not yield 4-anonymity. For  $\epsilon > 0.4$ , we can clearly achieve 4-anonymity, but as everything gets into a single class, all the data will be generalized to the same record and thus data quality is compromised.

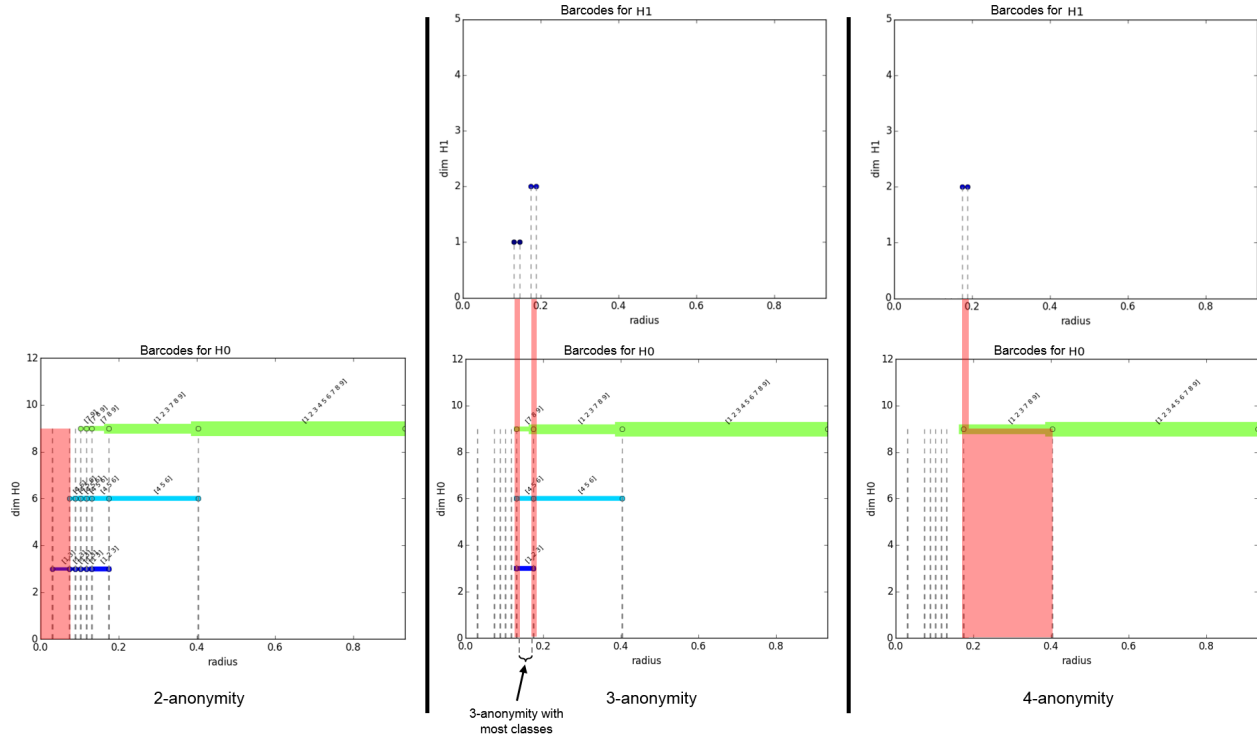
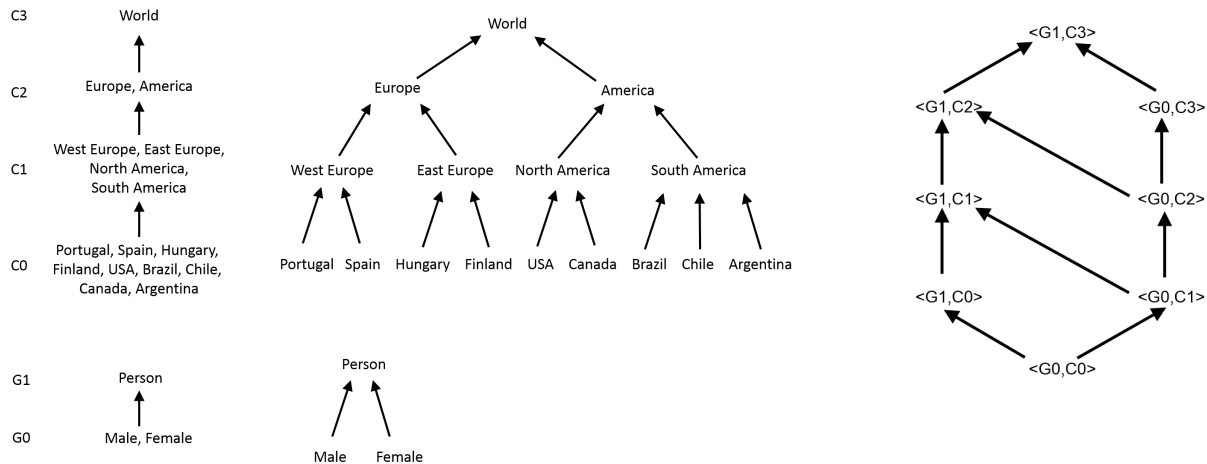


Figure 4: Weighted barcode showing the full spectrum of anonymity regimes. Although the barcode is only a single diagram, we split it here into three diagrams where we only show the simplices that meet the requirements of the  $k$ -anonymity indicated underneath each figure. With red bars we show regimes that cannot achieve 2,3,4-anonymity. These correspond to either situations where not all the  $k$ -simplices admit trivial higher order homology groups, as is the case for 3-anonymity between  $[0.11, 0.13]$  and  $[0.17, 0.19]$ , not all the simplices are  $\ell$ -simplices with  $\ell \geq k$ . This happens for example in 2-anonymity in the interval  $[0, 0.08]$  and between  $[0.19, 0.4]$  for 4-anonymity. We only show here  $\dim H_0$  and  $\dim H_1$  as higher order homology groups are trivial for this simple example.

**Remark 4.1 (Advantages of proposed approach)** *The main advantage of the method proposed is that it enables us not only to find a  $k$ -anonymization for a fixed  $k$  (if it exists), but also to provide alternative regimes that can be especially useful when  $k$ -anonymity cannot be achieved for the given  $k$ . This generally is not something that other algorithms, such as [5, 7, 8] directly provide. Indeed, one would need to run the same algorithms for various values of  $k$  to obtain the same tradeoff picture as we obtain. The persistent diagram we consider in this paper is instead computed in one shot from the filtration  $\{C^\epsilon\}$ . Furthermore, we leverage very scalable algorithms for such computation based on discrete Morse theory, see [25, 22]. Also, note that not only the persistent diagram allows us to determine the right regime that gives the desired  $k$ -anonymity, but we can also look at the other important tradeoff parameter such as the number of classes. For a given  $k$ , it is indeed possible to find various generalizations  $\epsilon$  that meet the  $k$ -anonymity requirement, however some might lead to equivalence classes with many more than  $k$  elements that is in general not desirable.*

## 5 Zig-Zag Persistent Homology for Mixed Data

The methodology we have described in the previous section has clearly nice properties but has some limitations. First, it is restricted to numerical attributes where the notion of a radius is well



(a) Example of generalization for two set of labels. On the top for *countries* and on the bottom for *gender*. (b) Generalization lattice for the attributes *countries* and *gender*.

Figure 5: In (a) we show two generalization trees and in (b) we show the corresponding generalization lattice.

defined and thus it would not work in the case where attributes are categorical in nature, such as, for example, strings, labels, social security numbers, etc. Second, growing balls (or polytopic approximations) in high dimensional spaces and determining intersections can be computationally challenging. In this section, we discuss an extension that leads to weighted persistent diagrams as the one described in the previous section and that can be used to explore anonymity tradeoffs.

For anonymization of categorical data, one needs to specify *generalization trees* that enforce a certain partial order between the various generalizations. Let us consider a simple example where there are two attributes: *Countries* and *Gender*. Examples of generalization tree are shown in Figure 5a where, as one moves from leaves to the root, the generalization becomes increasingly coarser, see for example [5]. We thus have that  $\{Portugal, Spain\} \prec \{West Europe\}$  and  $\{West Europe, East Europe\} \prec \{Europe\}$ , etc. or where  $\{Male, Female\} \prec \{Person\}$ . We may capture the generalization of numerical values within trees as well.

Formally, we have that a generalization tree is a map  $\mathcal{T} : \mathcal{N} \times [0, r] \rightarrow \mathcal{A}$  where  $\mathcal{N}$  is the set all nodes in the trees (including root and leaves) and  $[0, r]$  is the level of the generalization. For example,  $\mathcal{T}(\{USA, Canada\}, 2) = America$ . For this setting, we can extend the notion of anonymity complex as follows.

**Definition 5.1 (Generalized Anonymity Complex)** *Given a table  $T$  with  $N$  rows and a set of quasi-identifiers  $Q_T$  as well as a set of generalization trees  $\mathcal{T}_k$  for  $k = 1, \dots, |Q_T|$ , we define the generalization anonymity complex  $\Gamma$  as the simplicial complex whose  $k$ -simplices are determined by  $(k + 1)$   $|Q_T|$ -dimensional tuples  $\{A_i\}_{i=1}^k$  such that there exists a  $\bar{r}$  such that  $\mathcal{T}_i(A_i, \bar{r}) = \mathcal{T}_j(A_j, \bar{r})$ .*

The definition of a generalized anonymity  $k$ -simplex can be obtained as well. For example,  $\{Portugal, Spain, Hungary\}$  forms a generalized anonymity 3-simplex for the generalization level 2, considering the tree in Figure 5a.

**Definition 5.2 (Generalization lattice)** *Given a  $|Q_T|$ -dimensional tuple representing attributes and the associated trees  $\mathcal{T}_k$  for each attribute, we can construct a directed graph (lattice) where the*

vertices are the tuples

$$\langle \mathcal{T}_1(A_1, s), \mathcal{T}_2(A_2, s_2), \dots, \mathcal{T}_{|Q_T|}(\gamma_{|Q_T|}, s_{|Q_T|}) \rangle$$

and there is an edge between two vertices if there exists a generalization  $s_i + 1$  for an attribute  $A_i$ .

Figure 5b shows an example of a generalization lattice for the attributes  $\{Country, Gender\}$ .

Given the generalization trees and the previous definitions, we can see the similarity with the previous section, where we were interested in searching for regimes where the anonymity complexes have trivial high-order homology groups. The critical difference is, however, that these generalizations are not dependent on a single parameter ( $\epsilon$  in the previous section) anymore and thus we do not have a filtration in general.

In fact, the anonymity complexes, for various degrees of generalization, satisfy the following commutative diagram<sup>1</sup>: which is clearly not a filtration, but rather a *multi-dimensional filtration*.

$$\begin{array}{ccccccc} C^{10} & \hookrightarrow & C^{11} & \hookrightarrow & C^{12} & \hookrightarrow & C^{13} \\ \uparrow \wr & & \uparrow \wr & & \uparrow \wr & & \uparrow \wr \\ C^{00} & \hookrightarrow & C^{01} & \hookrightarrow & C^{02} & \hookrightarrow & C^{03} \end{array} \quad (1)$$

Unfortunately, computation of the barcodes associated to a multi-dimensional persistent homology is still an open question, although progress has been made in recent years, see [26] and references therein.

There are two ways we can tackle this problem, using either some further assumptions on the generalization or an approximation. The first is via *zig-zag persistent homology* [27]. In this context, we need to assume an ordering of the generalizations, something that is not uncommon, see [8]. Given the structure of the problem, a reasonable sequence of spaces to be considered would be  $C^{00}, C^{01}, C^{02}, C^{03}, C^{10}, C^{11}, C^{12}, C^{13}$  leading to the following commutative diagram:

$$\underbrace{C^{00} \hookrightarrow C^{01} \hookrightarrow C^{02} \hookrightarrow C^{03}}_{C^{00-3}} \hookrightarrow \underbrace{C^{10} \hookrightarrow C^{11} \hookrightarrow C^{12} \hookrightarrow C^{13}}_{C^{10-3}}.$$

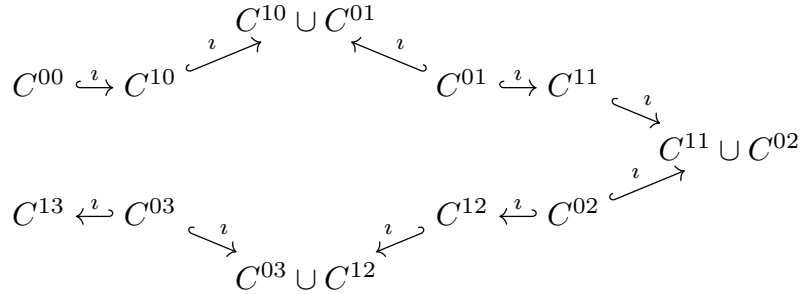
Note that the middle map is not an inclusion. It is possible to still compute persistent homology of the full chain by “joining” the two subsequences of spaces  $C^{00-3}$  and  $C^{10-3}$  through their union as follows: In this case, by using Mayer-Vietoris pyramid principle [27], we can compute the

$$\begin{array}{ccc} & C^{03} \cup C^{10} & \\ \nearrow \wr & & \nwarrow \wr \\ C^{00-3} & & C^{10-3} \end{array}$$

persistent homology of the chain and by passing through the union we obtain the barcodes for the initial chain [27]. Given the barcodes we are now in the same situation as in the previous section where we can search for different  $k$ -anonymity regimes.

Note however that the previous barcode would be different than the one built for the following sequence, say  $C^{00}, C^{10}, C^{01}, C^{11}, C^{02}, C^{12}, C^{03}, C^{13}$ . In this case we would need to apply Mayer-Vietoris pyramid principle multiple times:

<sup>1</sup>We denote with  $C^{00}$  the generalized anonymity complex corresponding to the generalization  $\langle G_0, C_0 \rangle$ , as in Figure 5b



An alternative approach to compute the barcodes for (1), is via an approximation. Specifically, one can consider the persistence equivalence theorem [28, Section 2.4.2] where, given the commutative diagram in (1), one can compute the persistent homology for the lower and upper chains. Because of the inclusion maps we have that the persistent homology modules of the two chains respect the inclusions. We can then proceed as follows: we first compute the persistent homology of the lower chain in the commutative diagram (1) and if it achieves the desired  $k$ -anonymity, then, because of the inclusion maps, such  $k$ -anonymity can be achieved also by the upper chain in (1), thus we do not need to compute the barcodes for the upper chain and stop after obtaining the barcodes for the lower one.

However, if  $k$ -anonymity is not achieved considering the lower chain, then persistent homology and corresponding barcodes need to be computed for the upper chain. Note that as we are approximating a multi-dimensional persistent homology by a sequence of (independent) persistent homology computation, if  $k$ -anonymity cannot be achieved following such process we cannot conclude that there is not a way to  $k$ -anonymize the data. This is a consequence of the natural complexity of the anonymity problem.

## 6 Conclusion

This paper introduced a new perspective to  $k$ -anonymity in data privacy based on algebraic topology. In particular, we addressed the case when the data lies in a metric space. We demonstrated how tools such as persistence homology can be applied to efficiently obtain the entire spectrum of  $k$ -anonymity of the database for various values of the radius of proximity. For this representation, we provided an analytic characterization of conditions under which a given representation of the dataset is  $k$ -anonymous. Finally, we discussed how this method can be extended to address the general case of a mix of categorical and metric data.

In future, it would be interesting to investigate other notions of privacy using these tools. In particular, applicability of such techniques to dynamic databases would be an interesting case which naturally arise in control applications.

## 7 Acknowledgments

The authors would like to thank Vidit Nanda for discussions on zig-zag persistent homology and the Perseus software used in this paper to compute persistent homology.

## References

- [1] “Terraswarm – Theme 3: Services, applications and cloud interactions,” <http://www.terraswarm.org/services/>, accessed 2015-09-27.
- [2] “Design technology for the trillion-device future,” <http://youtu.be/ViJ3SH5t4Ys>, presented at the DARPA Wait, What? conference, St. Louis, MO, 2015.
- [3] L. Sweeney, “k-anonymity: A model for protecting privacy,” *Int. J. of Uncertainty, Fuzziness and Knowledge-Based Systems*, vol. 10, no. 05, pp. 557–570, 2002.
- [4] A. Meyerson and R. Williams, “On the complexity of optimal k-anonymity,” in *ACM SIGMOD-SIGACT-SIGART symposium on Principles of database systems*, 2004, pp. 223–228.
- [5] K. LeFevre, D. J. DeWitt, and R. Ramakrishnan, “Incognito: Efficient full-domain k-anonymity,” in *ACM SIGMOD Int. Conf. on Management of Data*, 2005, pp. 49–60.
- [6] J. Byun, A. Kamra, E. Bertino, and N. Li, “Efficient k-anonymization using clustering techniques,” in *Advances in Databases: Concepts, Systems and Applications*, e. a. Kotagiri, R., Ed. Springer Berlin Heidelberg, 2007.
- [7] K. LeFevre, D. J. DeWitt, and R. Ramakrishnan, “Mondrian multidimensional k-anonymity,” in *Int. Conf. on Data Engineering*, 2006.
- [8] Y.-K. Kim, H. Lee, M. Yoon, and J.-W. Chang, “Hilbert-curve based data aggregation scheme to enforce data privacy and data integrity for wireless sensor networks,” *Int. J. of Distr. Sensor Networks*, 2013.
- [9] R. Ghrist and A. Muhammad, “Coverage and hole-detection in sensor networks via homology,” in *Int. Symp. on Information Processing in Sensor Networks*, 2005.
- [10] V. De Silva and R. Ghrist, “Homological sensor networks,” *Notices of the American Mathematical Society*, vol. 54, no. 01, pp. 10–17, 2007.
- [11] A. Muhammad and M. Egerstedt, “Control using higher order Laplacians in network topologies,” in *Int. Symp. of Mathematical Theory of Networks and Systems*, 2006, pp. 1024–1038.
- [12] A. Tahbaz-Salehi and A. Jadbabaie, “Distributed coverage verification in sensor networks without location information,” *IEEE Trans. on Aut. Control*, vol. 55, no. 8, pp. 1837–1849, 2010.
- [13] J. Derenick, A. Speranzon, and R. Ghrist, “Homological sensing for mobile robot localization,” in *IEEE Int. Conf. on Robotics and Automation*, 2013, pp. 572–579.
- [14] R. Ghrist, D. Lipsky, J. Derenick, and A. Speranzon, “Topological landmark-based navigation and mapping,” 2012. [Online]. Available: <http://www.math.upenn.edu/~ghrist/preprints/landmarkvisibility.pdf>
- [15] C. Dwork and A. Roth, “The algorithmic foundations of differential privacy,” *Found. and Trends in Theoretical Computer Science*, vol. 9, no. 3-4, pp. 211–407, 2013.
- [16] J. Le Ny and G. J. Pappas, “Differentially private filtering,” *IEEE Trans. on Aut. Control*, vol. 59, no. 2, pp. 341–354, 2014.

- [17] S. Han, U. Topcu, and G. J. Pappas, “Approximately truthful mechanism for electric vehicle charging via joint differential privacy,” in *American Control Conference*, 2015.
- [18] Y. Mo and R. Murray, “Privacy preserving average consensus,” in *IEEE International Conference on Decision and Control*, 2014.
- [19] A. Machanavajjhala, D. Kifer, J. Gehrke, and M. Venkatasubramanian, “ $\ell$ -diversity: Privacy beyond k-anonymity,” *ACM Trans. Knowl. Discov. Data*, vol. 1, no. 1, 2007.
- [20] N. Li, L. T., and S. Venkatasubramanian, “t-closeness: Privacy beyond k-anonymity and  $\ell$ -diversity,” in *IEEE Int. Conf. on Data Engineering*, 2007, pp. 106–115.
- [21] R. Ghrist, *Elementary Algebraic Topology*, 1st ed. Createspace.
- [22] “Perseus, the persistent homology software,” <http://www.sas.upenn.edu/~vnanda/perseus>, accessed 2015-09-27.
- [23] R. Dewri, I. Ray, and D. Whitley, “On the optimal selection of  $k$  in the  $k$ -anonymity problem,” in *IEEE Int. Conf. on Data Engineering*, 2008, pp. 1364–1366.
- [24] R. Ghrist, “Barcodes: the persistent topology of data,” *Bulletin of the American Mathematical Society*, vol. 45, no. 1, pp. 61–75, 2008.
- [25] K. Mischaikow and V. Nanda, “Morse theory for filtrations and efficient computation of persistent homology,” *Discrete & Computational Geometry*, vol. 50, no. 2, pp. 330–353, 2013.
- [26] G.-W. W. Kelvin Xia, “Multidimensional persistence in biomolecular data.” [Online]. Available: <http://arxiv.org/abs/1412.7679v1>
- [27] G. Carlsson, V. de Silva, and D. Morozov, “Zigzag persistent homology and real-valued functions,” in *Annual Symposium on Computational Geometry*. ACM, 2009, pp. 247–256.
- [28] V. Nanda, “Discrete Morse theory for filtrations,” Ph.D. dissertation, Department of Mathematics, Rutgers University, 2012. [Online]. Available: <http://www.sas.upenn.edu/~vnanda/source/Thesis.pdf>