



Sharing Knowledge.Preserving Privacy.

## **UnKnown: A Data Anonymization Tool**

# Carnegie Mellon University

## Engineering Privacy in Software

**Team Members:** Hana Habib, Preethi Josephina Mudialba, Siddharth Nair, Bill Quan, Chelse Swoopes, Lidong Wei  
**Project Advisor:** Rebecca Balebako

# Index

- Goal of Tool.....3
- Data Privacy.....4-5
- Suppression..... 6-7
- K-Anonymization.....8-12
- References & Additional Resources.....13-14

# Goal of Tool

Our goal is to protect the privacy of data subjects included in a particular dataset of interest to the WPRDC, as well as other datasets containing similar data. Attributes that our tool can anonymize are those which could lead to re identification because of personally identifiable information (PII).

# Data Privacy

PII is any information which can either be used as is or combined with additional information to identify an individual. PII can be used maliciously to negatively affect the life of an identified individual [1]. Besides an outside party knowing private information about the data subject, the outside party could use the information to negatively affect the identified data subject.

# Data Privacy

- E-mail addresses, dates of birth, addresses, and phone numbers, or combinations of attributes may lead to reidentification of individuals. Depending on the type of data included in the dataset, anonymization will reduce the risk of both objective harms, including identity theft or credit card fraud, and subjective harms, including unwanted attention and discomfort associated with the fear of being reidentified from the release of a dataset.
- To further emphasize the importance of dataset privacy, see the following video: <http://www.tripwiremovie.net>



# Suppression

- When using **unKnown**, users have the option to select whether an attribute is 'Sensitive,' 'Insensitive,' or 'Identifying.'
- When a user opts to select an attribute as 'Identifying,' **unKnown** suppresses the information associated with that attribute. Suppressing an attribute removes private information to protect an individual or group from being inferred or identified.

# Suppression

Example: After suppression, rows 1 and 3 are identical and rows 2 and 4 are identical [2].

**Original Dataset**

first	last	age	race
Harry	Stone	34	Afr-Am
John	Reyser	36	Cauc
Beatrice	Stone	34	Afr-Am
John	Delgado	22	Hisp



**Suppressed Dataset**

first	last	age	race
*	Stone	34	Afr-Am
John	*	*	*
*	Stone	34	Afr-Am
John	*	*	*



# The Theory of K-Anonymity

## 1. Who created k-anonymization? Why?

The theory of k-anonymity was developed by Dr. Latanya Sweeney. It is based on the the observation that 87% of the U.S. population is uniquely identified by date of birth, gender, postal code [3].

## 2. Why not use just suppression?

As shown in the example above, suppressing an attribute/column completely redacts information about that particular attribute. This can reduce the quality and usability of the dataset [3].

# The Theory of K-Anonymity

## 3. What is the purpose of K-Anonymity?

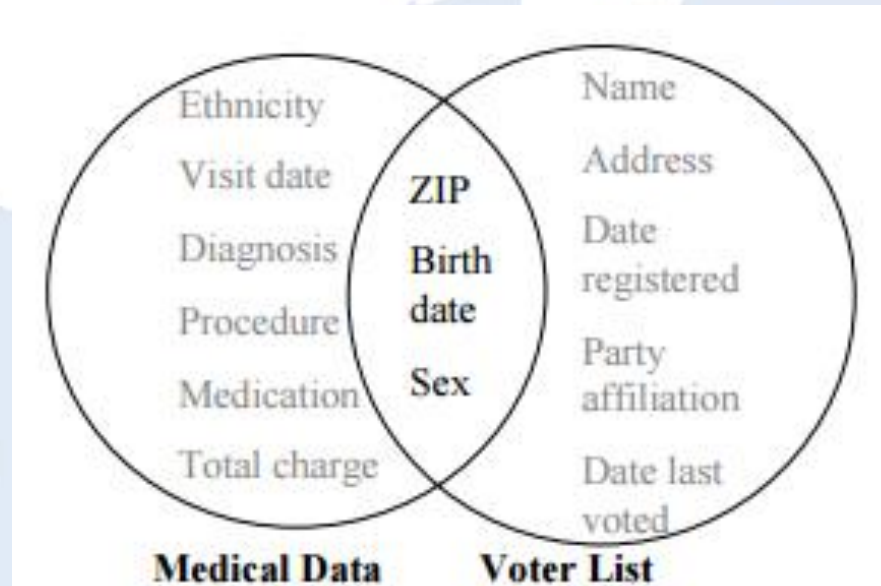
The purpose is to reduce linking of data subjects across different datasets while also keeping a useful dataset [3].

## 4. What causes linking?

Quasi-identifiers - “set of data elements in entity-specific data that in combination associates uniquely or almost uniquely to an entity and therefore can serve as a means of directly or indirectly recognizing the specific entity that is the subject of the data” [3]

# The Theory of K-Anonymity

Example: From medical data and voter list, the quasi-identifiers that can be used to link an individual to the data include zip code, birthdate, and sex [3]



# The Theory of K-Anonymity

## 5. How does k-anonymity work?

Let's say you have a table, represented as  $RT$ , with  $A_1, \dots, A_n$  attributes. The quasi-identifiers associated with your table are represented as  $QI_{RT}$  [4]. k-anonymity is satisfied if and only if "each sequence of values in  $RT[QI_{RT}]$  appears with at least  $k$  occurrences in  $RT[QI_{RT}]$ " [4].

# The Theory of K-Anonymity

Example : A table that satisfies k-anonymity:  $k=2$  and  $QI = \{\text{Race, Birth, Gender, Zip}\}$  [4]

	Race	Birth	Gender	ZIP	Problem
t1	Black	1965	m	0214*	short breath
t2	Black	1965	m	0214*	chest pain
t3	Black	1965	f	0213*	hypertension
t4	Black	1965	f	0213*	hypertension
t5	Black	1964	f	0213*	obesity
t6	Black	1964	f	0213*	chest pain
t7	White	1964	m	0213*	chest pain
t8	White	1964	m	0213*	obesity
t9	White	1964	m	0213*	short breath
t10	White	1967	m	0213*	chest pain
t11	White	1967	m	0213*	chest pain

Notice: For each tuples in the table  $[T]$ , the values of the tuple that comprise the quasi-identifier appear at least twice in the table [4].

→  $t1[QI_T] = t2[QI_T]$ ,  $t3[QI_T] = t4[QI_T]$ ,  $t5[QI_T] = t6[QI_T]$ ,  $t7[QI_T] = t8[QI_T] = t9[QI_T]$ , and  $t10[QI_T] = t11[QI_T]$



# References

1. B. Raghunathan. The Complete Book of Data Anonymization From Planning to Implementation. Chapter 1. Published 2013. [http://www.odbms.org/wp-content/uploads/2014/03/The-Complete-Book-of-Data-Anonymization\\_Chap\\_1.pdf](http://www.odbms.org/wp-content/uploads/2014/03/The-Complete-Book-of-Data-Anonymization_Chap_1.pdf)
2. Author Unknown. "K-Anonymity". Published 2007. <https://www.cs.cmu.edu/~jblocki/Slides/K-Anonymity.pdf>
3. L. Sweeney. "Simple Demographics Often Identify People Uniquely. Published 2000. <https://dataprivacylab.org/projects/identifiability/paper1.pdf>
4. L. Sweeney. "k-anonymity: a model for protecting privacy". Received May 2002. [https://epic.org/privacy/reidentification/Sweeney\\_Article.pdf](https://epic.org/privacy/reidentification/Sweeney_Article.pdf)



# Additional Resources

**For additional information on data privacy, suppression, and k-anonymity, please see the following references.**

- P. Samarati and L. Sweeney. “Generalizing Data to Provide Anonymity when Disclosing Information”. Publish Date NA.  
<https://pdfs.semanticscholar.org/5c15/b11610d7c3ee8d6d99846c276795c072eec3.pdf>
- L. Ohno-Machado, S. Venterbo, and S. Dreiseitl. “Effects of Data Anonymization by Cell Suppression on Descriptive Statistics and Predictive Modeling Performance”. Published Nov 2002.  
<https://www.ncbi.nlm.nih.gov/pmc/articles/PMC419433/>



Sharing Knowledge.Preserving Privacy.

## UnKnown: A Data Anonymization Tool