
BCSE318L – DATA PRIVACY – GEETHA MARY A

DIGITAL ASSIGNMENT 1

DEMONSTRATION OF ANONYMIZATION TOOLS

AATHIRA S

21BCE3680

Contents

I. Introduction	3
II. Anonymization	3
K-anonymization	3
III. Anonymization Tools	4
ARX Data Anonymization Tool	4
IV. Demonstration	5
Steps for anonymization	5
Results of anonymization	10
V. Conclusion.....	13
VI. References.....	13

I. Introduction

There is an increasingly high amount of data being collected today which presents both opportunities and challenges. While this large warehouse of data may be used to gain valuable insights through data analysis, the privacy and security of the individuals associated with these data records are paramount. Anonymization plays a crucial role in addressing and mitigating such concerns. [1]

Anonymization works by removing and modifying personally identifiable information (PII) from the datasets. This facilitates valuable research and collaboration, as well as the sharing of sensitive information, all while protecting the individuals' privacy. [1]

The following section II discusses anonymization and its subset k-anonymization in greater detail. Section III delineates various anonymization tools available today. The capabilities of ARX, a leading anonymization tool, is demonstrated in section IV. Finally we conclude our findings in section V.

II. Anonymization

Anonymization aims to limit the ability to link (or match) released data to other external information sources while still maintaining the useability and applicability of the released data. Apart from personally identifiable information, there exists other attributes which could be used along with external information for linkage attacks. Such attributes are known as quasi-identifiers. They include explicit identifiers such as name, phone number and address, as well as attributes whose combinations can be used to uniquely identify individuals such as ZIP code, birth date, and gender [2]. A goal of anonymity is to release person-specific data such that the ability to link to other information using quasi-identifier is limited.

K-anonymization

A common approach within anonymity is k-anonymity. By applying this technique, we can ensure that any combination of identifying attributes within a dataset appears for at least k individuals, thus creating the effect of hiding the individual with a crowd and thereby making it difficult to link specific data points to one single individual. This 'crowd' is referred to as an *equivalence class*. [3]

K-anonymity achieves such protection primarily by using the following methods:

- Generalization, where specific values are replaced with broader ranges, or more generalized values. For example, the age 29 can be replaced with an age group range of [20,30).
- Suppression, where data points that could potentially lead to reidentification are removed. For example, if multiple people are born in the 20th century (1900s), but none in the same decade (example, 1990s), then their birth year data can be suppressed as 19**.

Table 1 illustrates k-anonymization of patients' medical data. Here the quasi-identifiers are {Race, Birth, Gender, ZIP} and k=2. Therefore, there exists 2 tuples for which the quasi-identifiers are the same.

Table 1: 2-anonymity in sample medical database [2]

	Race	Birth	Gender	ZIP	Problem
t1	Black	1965	M	02141	Short breath
t2	Black	1965	M	02141	Chest pain
t3	Black	1964	F	02138	Obesity
t4	Black	1964	F	02138	Chest pain
t5	White	1964	M	02138	Chest pain
t6	White	1964	M	02138	Chest pain
t7	White	1964	M	02138	Short breath

K-anonymization is not a perfect anonymization technique as it may be susceptible to attacks based on homogeneity within an equivalence class. Consider tuples *t5* and *t6* in Table 1. Both tuples belong to the same equivalence class and have the same medical problem. Hence, an attacker whose target individual is White, born in 1964, and lives in the ZIP area 02138 will have *chest pain*. [4]

III. Anonymization Tools

In order to handle and anonymize large datasets we require some mechanism for automation. For this purpose, we have a wide array of anonymization tools at our disposition. Some popular applications include ARX data anonymizer, Clover Dx's Data Anonymization Tool, Amnesia and μ -Argus.

ARX Data Anonymization Tool

The ARX Data Anonymization Tool is a comprehensive open source software for anonymizing structured individual-level health data. It has been continuously developed and extended with further functionalities over the past years and is freely available on the project website. Using anonymization or de-identification, data can be transformed in a semi-automatic manner to reduce risks to the privacy of individuals when data is processed for secondary purposes or disclosed to third parties [5].

The key features of ARX are as follows [6]:

- Transformation models: More than 10 different ways of transforming data, including various ways of generalization, top and bottom coding, suppression, sampling, and micro-aggression.
- Privacy models: ARX supports more than 10 different methods for quantifying and protecting privacy, including traditional syntactical models such as k-anonymity, statistical models for population uniqueness, and state-of-the-art semantic models including a game-theoretic approach and differential privacy.
- Utility models: The tool supports more than 10 different models for quantifying the utility of data during the anonymization process, including general-[purpose] methods reflecting data and fidelity or changes to value distributions as well as application-specific models

In the following section we illustrate how to use ARX to perform k-anonymization.

IV. Demonstration

Data Used

The dataset train.csv [7] represents the passengers on the Titanic, and their details, including which of the passengers survived its sinking.

Table 2: Data dictionary for data set [7]

Variable	Definition	Key
survival	Survival	0 = No, 1 = Yes
pclass	Ticket class	1 = 1 st , 2 = 2 nd , 3 = 3 rd
name	Passenger name	
sex	Sex	
age	Age in years	
sibsp	# of siblings / spouses aboard the Titanic	
parch	# of parents / children aboard the Titanic	
ticket	Ticket number	
fare	Passenger fare	
cabin	Cabin Number	
embarked	Port of Embarkation	C = Cherbough, Q= Queenstown, S = Southampton

Table 3: Classification of attributes in dataset [7]

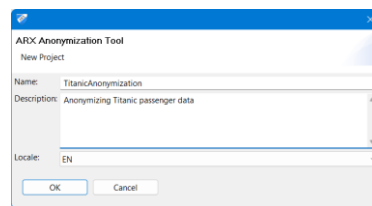
Variable	Classification
survival	Sensitive
pclass	Quasi-identifying
name	Identifying
sex	Quasi-identifying
age	Quasi-identifying
sibsp	Quasi-identifying
parch	Quasi-identifying
ticket	Identifying
fare	Quasi-identifying
cabin	Identifying
embarked	Quasi-identifying

Steps for anonymization

1. Create new project

On the taskbar ribbon at the top, go to File → New Project. Fill in the details and click OK.

Figure 1: Create new project on ARX



2. Import data

Again, go to the taskbar ribbon, File → Import Data. Select your file type. Here, we have used a CSV file. Click on *Browse* and select your file from the file explorer. Click on *Next*. Select *perform data cleansing* and then select *Next*. Confirm your data in the preview and *Finish*.

Figure 2: select input data's file type (left); Verify your data before completing import process (right).

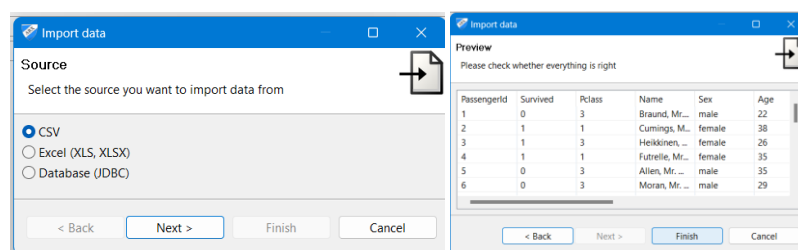
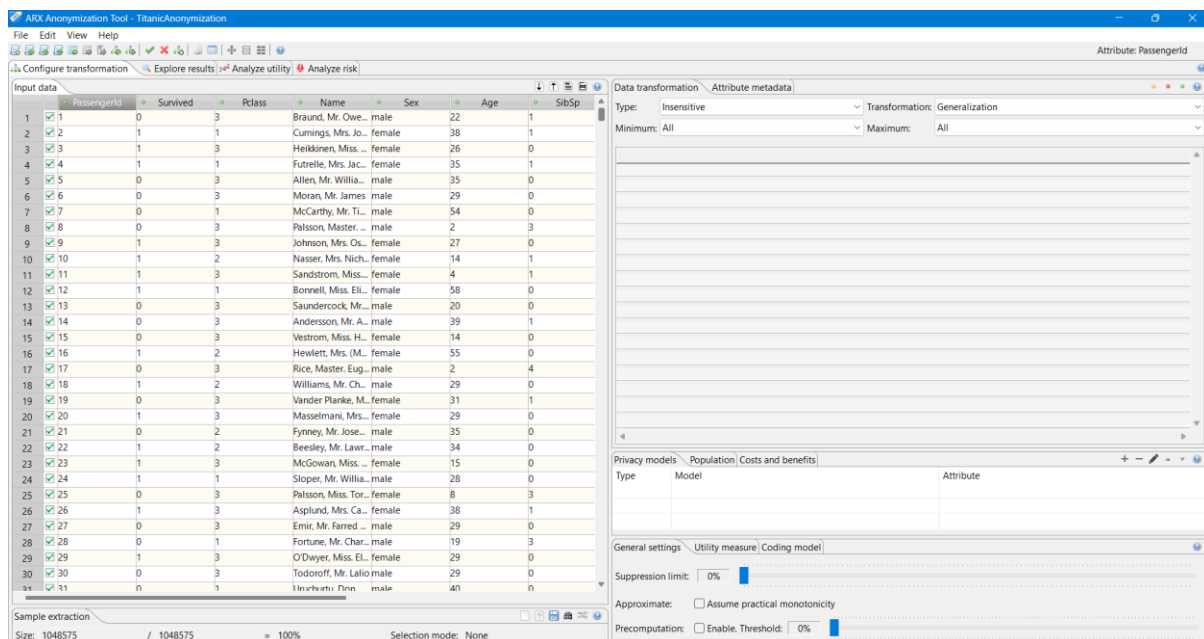


Figure 3: Data is imported into the project.

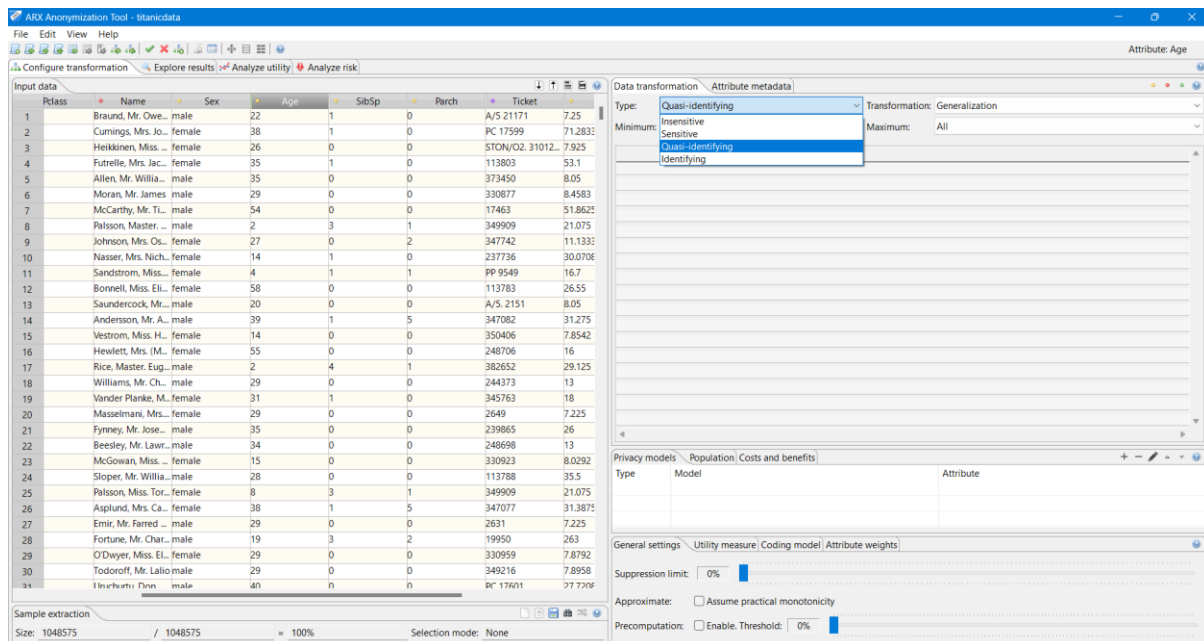


3. Set classification of attributes

For each attribute column in the dataset, set its type from one among the following:

- Insensitive
- Sensitive
- Quasi-identifying
- Identifying

Figure 4: Setting attribute type as 'Quasi-identifying' for attribute 'Age'



4. Set hierarchies for each attribute

Select the column for which you want to create a hierarchy. Now go to the taskbar ribbon select Edit→ Create Hierarchy. Now select the desired method to generate the hierarchy based on the column's data type.

- Using dates, for date type data
- Using intervals, suitable for variables with a ratio scale.
- Using ordering, suitable for variables with ordinal scale.
- Using masking, suitable for alphanumeric strings
- Using priorities, suitable for a particular priority factor such as frequency.

Creating hierarchies using intervals, ordering and masking is illustrated in figures 5, 6 and 7 respectively.

Figure 5: Creating hierarchy using intervals

Hierarchy wizard

Create a hierarchy by defining intervals

Specify the parameters. Note: Aggregate functions are only applied to interval limits.

[1, 81[[1, 81[

General Range Interval Group

Lower bound

Minimum value: 1

Bottom coding from: 1

Snap from: 1

Upper bound

Snap from: 80

Top coding from: 80

Maximum value: 81

Help... Load... Save... < Back Next > Finish Cancel

(a)

Hierarchy wizard

Create a hierarchy by defining intervals

Specify the parameters. Note: Aggregate functions are only applied to interval limits.

[1, 40[[1, 40[

General Range Interval Group

Aggregate function: Default

Function Parameter:

Min: 1

Max: 40

Help... Load... Save... < Back Next > Finish Cancel

(b)

Hierarchy wizard

Review the hierarchy

Overview of groups and values

#Groups Table

	Level-0	Level-1	Level-2
71	1	[1, 40[*
4	2	[1, 40[*
1	3	[1, 40[*
	4	[1, 40[*
	5	[1, 40[*
	6	[1, 40[*
	7	[1, 40[*
	8	[1, 40[*
	9	[1, 40[*
	10	[1, 40[*
	11	[1, 40[*

Help... Load... Save... < Back Next > Finish Cancel

(c)

Figure 6: Create hierarchy using ordering

Hierarchy wizard

Create a hierarchy by ordering and grouping items

Specify the parameters

Order

Values

C

Q

S

Move up

Move down

Order: Custom

Groups

2 Set

General Group

Aggregate function: Default

Function Parameter:

Size: 2

Help... Load... Save... < Back Next > Finish Cancel

(a)

Hierarchy wizard

Review the hierarchy

Overview of groups and values

#Groups

4

2

1

Level-0	Level-1	Level-2
	{, C}	*
C	{, C}	*
Q	{Q, S}	*
S	{Q, S}	*

Help... Load... Save... < Back Next > Finish Cancel

(b)

Figure 7: Creating hierarchy using masking

Hierarchy wizard

Create a hierarchy by masking characters

Specify the parameters

Alignment

☒ Align items to the left

☐ Align items to the right

Masking

☐ Mask characters left to right

☒ Mask characters right to left

Characters

Padding character ()

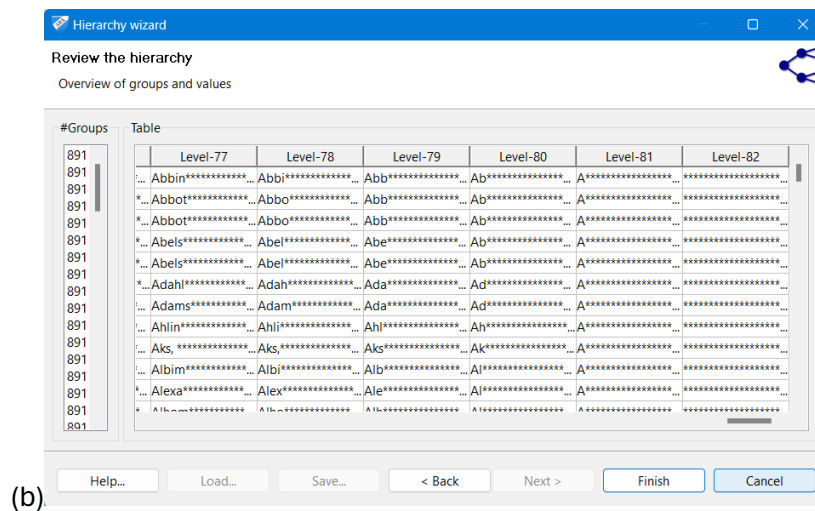
Masking character (*)

Domain properties

Domain size 891 Alphabet size 60 Max. characters 82

Help... Load... Save... < Back Next > Finish Cancel

(a)

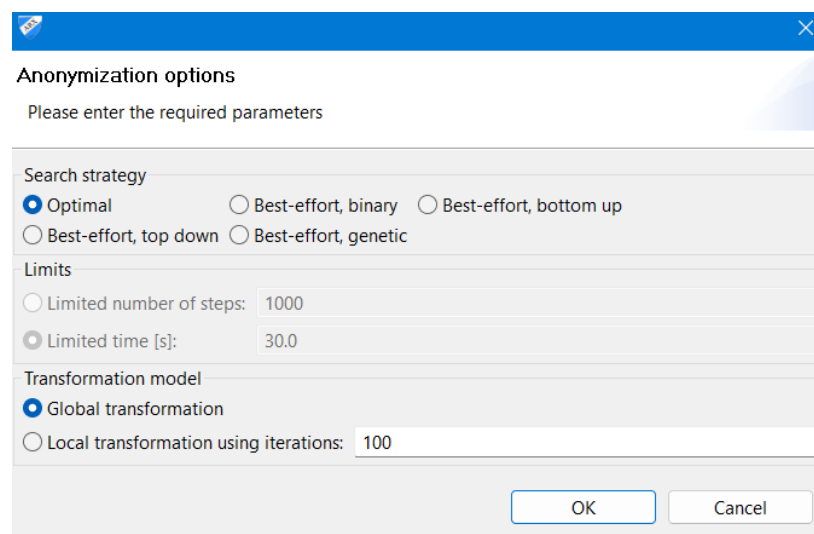


(b)

5. Anonymize the dataset

On the taskbar ribbon, select Edit→Anonymize. Choose the appropriate parameters and click OK.

Figure 8: Selecting parameters for anonymization



Results of anonymization

The results of anonymization and its effectiveness can be viewed and analysed by using the exploring the different tabs on the application interface. The results are illustrated in figures 9 to 13. We observe a decrease in overall risk of re-identification posed by the data. However there is also a drop in the utility values of the data. Moreover, the effectiveness of various attacker models against the dataset has been reduced significantly. Thus we can conclude that the results of anonymization helped protect the data from re-identification to a large extent.

Figure 9: Comparing the raw data before and after anonymization. Equivalence classes are highlighted.

ARX Anonymization Tool - titanicdata

Attribute: SibSp Transformations: 972 Selected: [0, 0, 1, 1, 1, 1, 0] Applied: [0, 0, 1, 1, 1, 1, 0]

Input data: Classification performance Quality models

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp
1	53	1	1	Harper, Mrs. Hen.	female	49	1
2	178	0	1	Isham, Miss. Ann.	female	50	0
3	195	1	1	Brown, Mrs. Jam.	female	44	0
4	196	1	1	Lurette, Miss. Elise	female	58	0
5	300	1	1	Baxter, Mrs. Jan.	female	50	0
6	320	1	1	Speiden, Mrs. Fr.	female	40	1
7	338	1	1	Burns, Miss. Eliza.	female	41	0
8	367	1	1	Warren, Mrs. Fran.	female	60	1
9	381	1	1	Bidois, Miss. Ros.	female	42	0
10	497	1	1	Eustis, Miss. Eliza.	female	54	1
11	514	1	1	Rothschild, Mrs.	female	54	1
12	524	1	1	Hippach, Mrs. Lo.	female	44	0
13	557	1	1	Duff Gordon, La.	female	48	1
14	592	1	1	Stephenson, Mrs.	female	52	1
15	880	1	1	Potter, Mrs. Tho.	female	56	0
16	35	0	1	Meyer, Mr. Edgar.	male	28	1
17	65	0	1	Stewart, Mr. Albe.	male	29	0
18	98	1	1	Greenfield, Mr.	male	23	0
19	119	0	1	Baxter, Mr. Quig.	male	24	0
20	140	0	1	Giglio, Mr. Victor	male	24	0
21	274	0	1	Natsch, Mr. Charl.	male	37	0
22	296	0	1	Lewy, Mr. Ervin G	male	29	0
23	371	1	1	Harder, Mr. Geor.	male	25	1
24	374	0	1	Ringhini, Mr. Sante	male	22	0

Summary statistics Distribution Contingency Class sizes Properties Classification models

Parameter	Value
Scale of measure	Ratio scale
Number of measures	866
Number of distinct values	7
Mode	0
Median	0
Min	0

Output data: Classification performance Quality models

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp
1	53	1	1	female	female	[40, 79]	[0, 5]
2	178	0	1	female	female	[40, 79]	[0, 5]
3	195	1	1	female	female	[40, 79]	[0, 5]
4	196	1	1	female	female	[40, 79]	[0, 5]
5	300	1	1	female	female	[40, 79]	[0, 5]
6	320	1	1	female	female	[40, 79]	[0, 5]
7	338	1	1	female	female	[40, 79]	[0, 5]
8	367	1	1	female	female	[40, 79]	[0, 5]
9	381	1	1	female	female	[40, 79]	[0, 5]
10	497	1	1	female	female	[40, 79]	[0, 5]
11	514	1	1	female	female	[40, 79]	[0, 5]
12	524	1	1	female	female	[40, 79]	[0, 5]
13	557	1	1	female	female	[40, 79]	[0, 5]
14	592	1	1	female	female	[40, 79]	[0, 5]
15	880	1	1	female	female	[40, 79]	[0, 5]
16	35	0	1	male	male	[1, 40]	[0, 5]
17	65	0	1	male	male	[1, 40]	[0, 5]
18	98	1	1	male	male	[1, 40]	[0, 5]
19	119	0	1	male	male	[1, 40]	[0, 5]
20	140	0	1	male	male	[1, 40]	[0, 5]
21	274	0	1	male	male	[1, 40]	[0, 5]
22	296	0	1	male	male	[1, 40]	[0, 5]
23	371	1	1	male	male	[1, 40]	[0, 5]
24	374	0	1	male	male	[1, 40]	[0, 5]

Summary statistics Distribution Contingency Class sizes Properties Classification models

Parameter	Value
Scale of measure	Ordinal scale
Number of measures	783
Number of distinct values	1
Mode	[0, 5]
Median	[0, 5]
Min	[0, 5]
Max	[0, 5]

Figure 10: Comparing the quality of the data before and after anonymization, based on various attributes.

ARX Anonymization Tool - titanicdata

Attribute: SibSp Transformations: 972 Selected: [0, 0, 1, 1, 1, 1, 0] Applied: [0, 0, 1, 1, 1, 1, 0]

Input data: Classification performance Quality models

Attribute-level quality

Attribute	Data type	Missings
Pclass	Integer	0%
Sex	String	0%
Age	Integer	2.80584%
SibSp	Integer	0%
Parch	Integer	0%
Fare	Decimal	0%
Embarked	String	0%

Output data: Classification performance Quality models

Attribute-level quality

Attribute	Data type	Missings	Gen. intensity	Granularity	N-U entropy	Squared error
Pclass	Integer	10.54994%	89.45006%	89.45006%	89.84448%	88.71182%
Sex	String	10.54994%	89.45006%	89.45006%	89.95443%	89.45006%
Age	String	12.12121%	45.51066%	43.64919%	13.79424%	66.90964%
SibSp	String	10.54994%	44.72503%	29.81669%	3.35172%	68.35985%
Parch	String	10.54994%	44.72503%	44.78114%	4.29457%	68.62806%
Fare	String	10.54994%	44.72503%	1.08644%	0.47457%	71.51739%
Embarked	String	10.54994%	89.45006%	89.45006%	85.58244%	91.30115%

Dataset-level quality

Model	Quality
Gen. intensity	64.00513%
Granularity	49.79718%
N-U entropy	21.63149%
Discernibility	80.40236%
Average class size	96.40094%
Record-level squared error	49.00854%
Attribute-level squared error	71.45718%
Aggregation-specific squared error	0%

Summary statistics Distribution Contingency Class sizes Properties Classification models

Feature variables

Ena.	Type	Name	Scaling
✗	✗	PassengerId	Categorical
✗	✗	Survived	Categorical
✗	✗	Pclass	Categorical
✗	✗	Name	Categorical
✗	✗	Sex	Categorical

Target variables

Ena.	Type	Name
✗	✗	PassengerId
✗	✗	Survived
✗	✗	Pclass
✗	✗	Name
✗	✗	Sex

Classifier: Logistic regression

Parameter	Value
Alpha	1
Decay exponent	0.2
Lambda	0.00001
Learning rate	1
Prior function	L1

Figure 11: Comparing the distribution of risks across the data points before and after anonymization

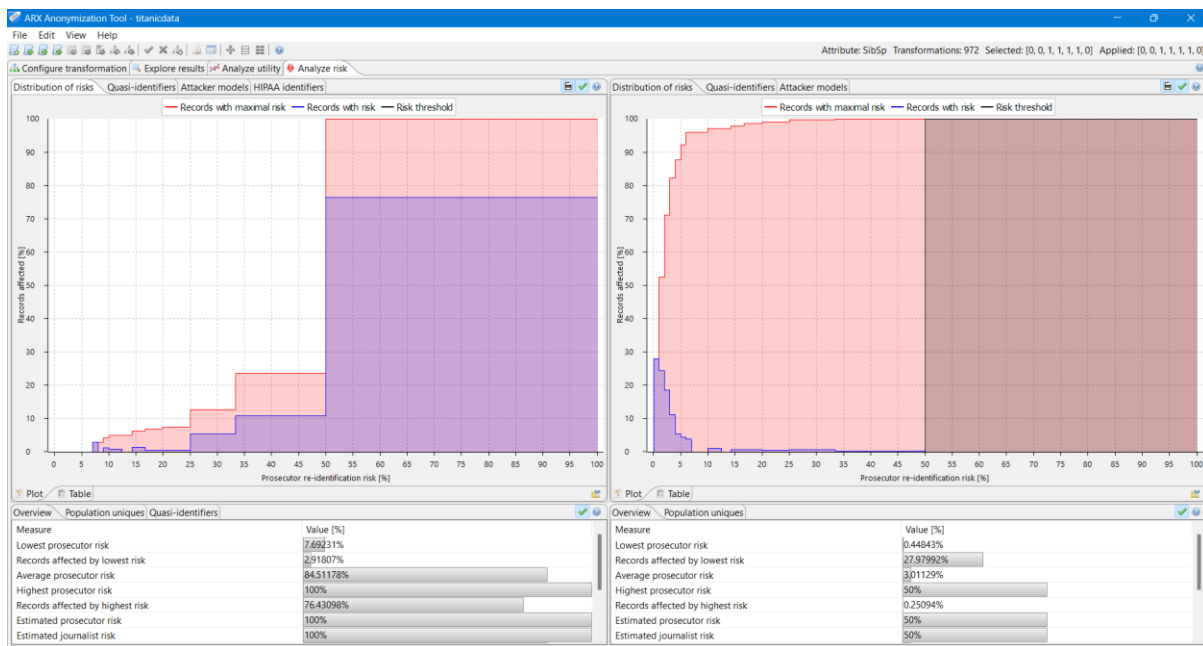


Figure 12: Comparing the effectiveness of attacker models before and after anonymization

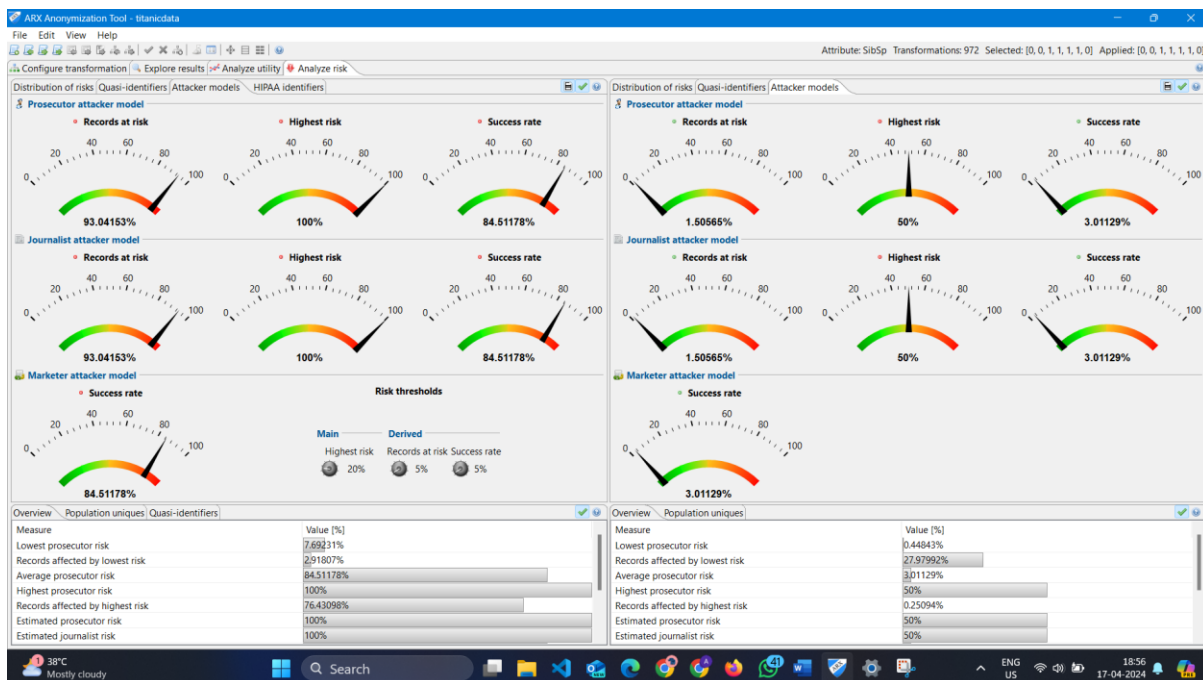
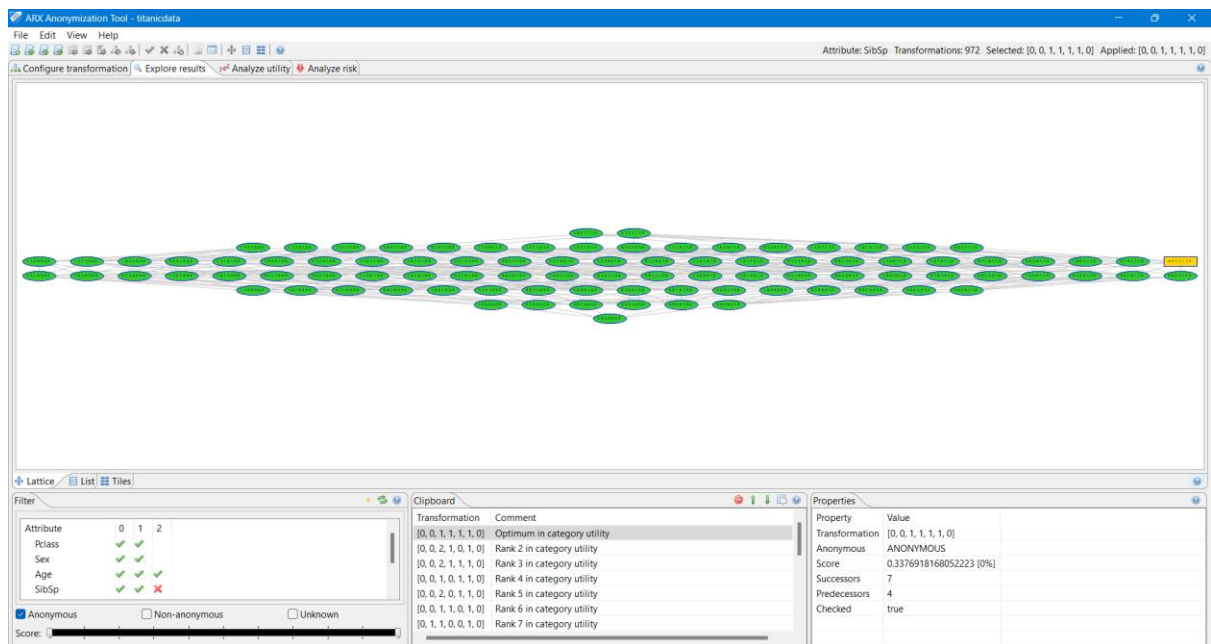


Figure 13: Exploring the results as nodes in a lattice of data in the anonymized dataset



V. Conclusion

This report has demonstrated the capabilities of ARX as a valuable tool for anonymizing data. Through its diverse anonymization techniques, ARX empowers users to achieve a strong balance between data privacy and analytical utility.

The demonstration showcased how ARX can be leveraged to achieve k-anonymity, a widely used technique for protecting individual privacy within datasets. By offering various methods like generalization and suppression, ARX allows users to tailor the anonymization process to their specific needs.

VI. References

- [1] A. Chia, "Data Anonymization Explained: What You Need To Know," 25 May 2023. [Online]. Available: https://www.splunk.com/en_us/blog/learn/data-anonymization.html.
- [2] L. Sweeney, "Achieving k-anonymity privacy protection using generalization and suppression," *International Journal on Uncertainty, Fuzziness and Knowledge-based Systems*, pp. 571-588, 2002.
- [3] L. Sweeney, "k-Anonymity: a model for protecting privacy," *International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems*, 2002.
- [4] A. S. Nataraj Venkataramanan, *Data Privacy: Principles and Practices*, Taylor and Francis.
- [5] ARX Deidentifier, "Overview of supported anonymization methods," 2024. [Online].

Available: <https://arx.deidentifier.org/overview/>.

[6] ARX - Data Anonymisation Tool, "Overview of Features of the ARX Data Anonymization Tool," 28 March 2021. [Online]. Available: https://www.youtube.com/watch?v=mkfJ-td-B94&ab_channel=ARX-DataAnonymizationTool.

[7] W. Cukierski, "Titanic - Machine Learning from Disaster," Kaggle, 2012.