

Privacy shield –The Data Anonymizer

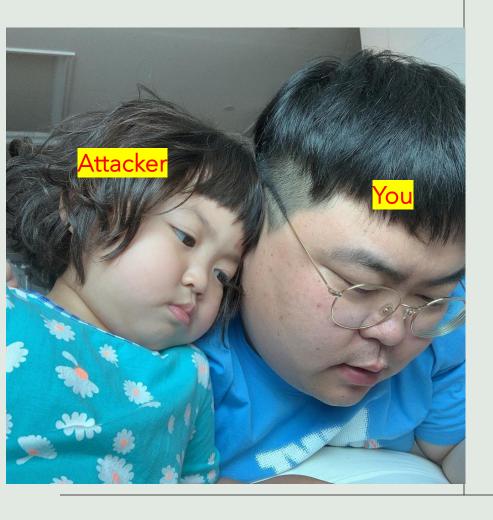
Data Anonymizer

 Overview: This document contains 5 sections, section-1 is about introduction to the project, section-2 and 3 contains the design and development of project, section-4 contains the user manual and section-5 contains future work





1.1. Context





The digital world of now is composed of personalized data services and there is a huge demand for data, especially for personal data



Large volumes of data can now be employed quite efficiently to gain new insights on a certain phenomenon, and data mining algorithms can give a lot of information about persons, events or entities



The purposes for which one requires meaningful personal data are very diversified, for example, research or public health policy purposes, to develop new services and products, to enhance the efficiency and effectiveness of new drugs, or even to simply condition people's behaviour, when they make a purchase on an online store

1.2. Motivation



In recent years, escalation of technology led to an increase in the capability to record and store personal data about consumers and individuals



To mitigate these issues, some de-identification methodologies have recently been proposed that, in some well controlled circumstances, allow for the re-use of personal data in privacypreserving ways



We will present a web tool that provides a way to easily anonymize data



Anonymization

Anonymity, means simply that a person is not identifiable

The anonymization process is intended to irreversibly remove the association between an individual and some information that can identify this person

Personal data is anonymized to protect the privacy of subjects when storing or disclosing data

1.3. Scope of Project



The main objective of this project is to create web application which make it easy for a user to anonymize the data



It means the user can upload the dataset and can choose configuration or create a new configuration with in the web interface for de-identification, after that the de-identified or anonymised data can be downloaded



The important task of this project was to create an interface for the user such that a new configuration can be developed on the dataset



2. Architecture



In order to give the web application some structure, a layered architecture was created, grouping the different components in sections



This layered separation is often used on web platforms, on this particular architecture is separated in three layers



Each of these layers is responsible for a particular and unique processing, they communicate with each other through different frame works and languages, so, aggregating all those tiers in a single framework simplifies hugely this interconnection process

Presentation Layer





This layer contains the user oriented functionalities, they are responsible for managing user interaction with the system, and generally consists of components that provide a common bridge to the business layer core

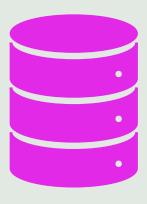
Presentation tier components allow users to interact with the application

Business Tier

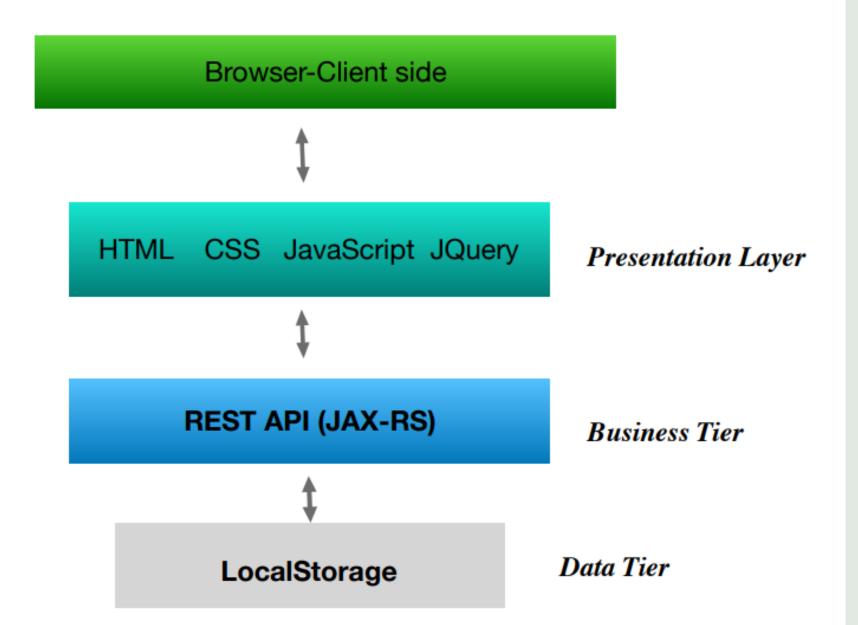
This layer implements the core functionality of the system and encapsulates the relevant business logic

It generally consists of elements that execute rules and separates them from the user interface and data access

Data Tier



The data tier provides access to the databases and storage devices used by a three-tier applications. It is responsible for retrieving data and transforming into a suitable format for the rest of the application. Essentially, this layer stores and retrieves information to the business tier for processing and eventually to the presentation tier. On adapting the ARX API there was no need to use a database, because saving the files locally was sufficient enough to perform the anonymization.



Technologies

Client-side

- In the application front-end, we used five technologies: HyperText Markup Language (HTML), Cascading Style Sheets (CSS), JavaScript, Jquery, and Asynchronous JavaScript and XML (AJAX). The web page presented on the user's browser was built on HTML, created from scratch, exhibiting a simple but intuitive interface that guides the user through a configuration process that sometimes could be a complex task. In order to make the HTML more attractive, we applied CSS, giving a better visual presentation to the interface.
- Another technology used in order to send and receive data from the server, is called JavaScript Object Notation (JSON), and is useful for parsing JavaScript objects into text. This is needed when the client wants to exchange data with the server because this data has to be in text format. JSON converts JavaScript objects into text, and vice-versa, the data is saved as pairs of keys/values and they are kept as an array of strings. This technology was used in all the AJAX requests created, parsing the data sent to the server, or making the reverse process for the received array of data.

JERSEY: Jersey RESTful Web Services framework is open source, production quality framework for developing RESTful Web Services in Java that provides support for JAX-RS APIs

REST: A web application without resources don't serve for much, they are usually something that can be stored on a computer, e.g., an electronic document, a dataset or the result of executing an algorithm

Tomcat: To deploy the services and resources, apache tomcat is being used as server

Server-side

ARX API



The ARX open source tool allows the user to change structured sensitive personal data, using Statical Disclosure Control, into data that can be shared securely



ARX API has been chosen for anonymization process because of its simplicity and organised code, with a good documentation, plus a deductive online javadoc with all packages and classes included in the API



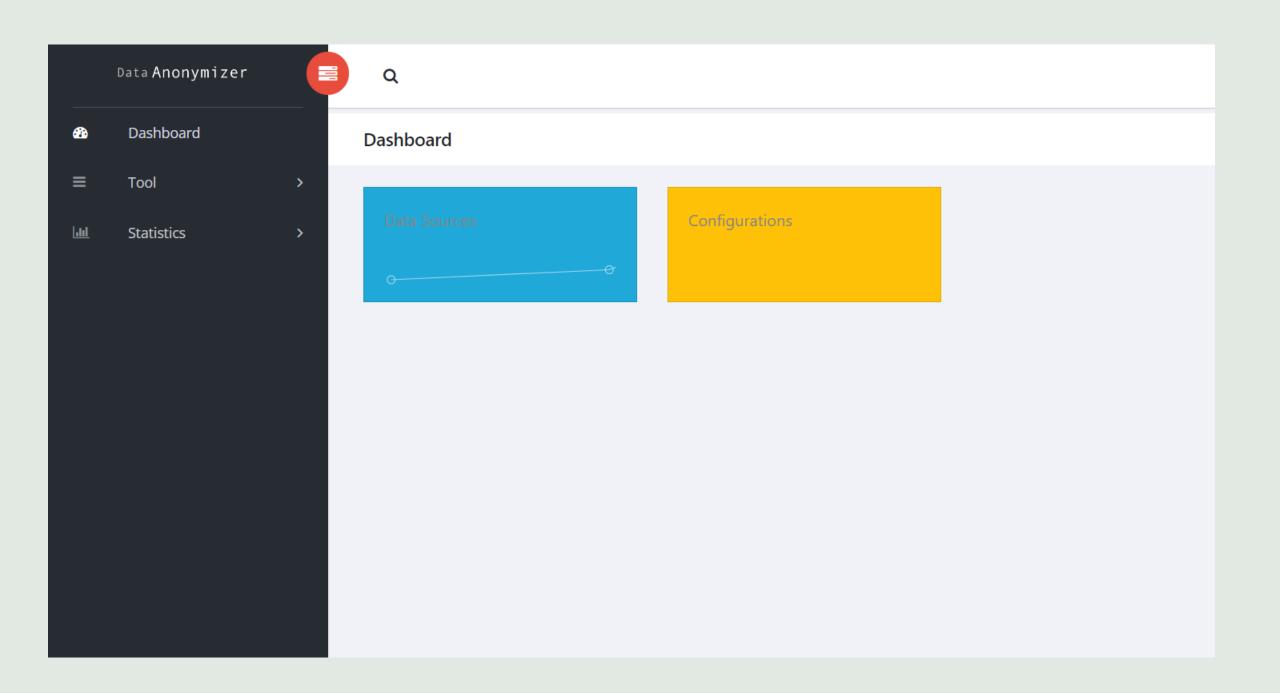
The ARXConfiguration defines a group of settings that are sent to the ARXAnonymizer, this configuration allows multiple parameters, offering the user an opportunity to create suitable deidentified datasets

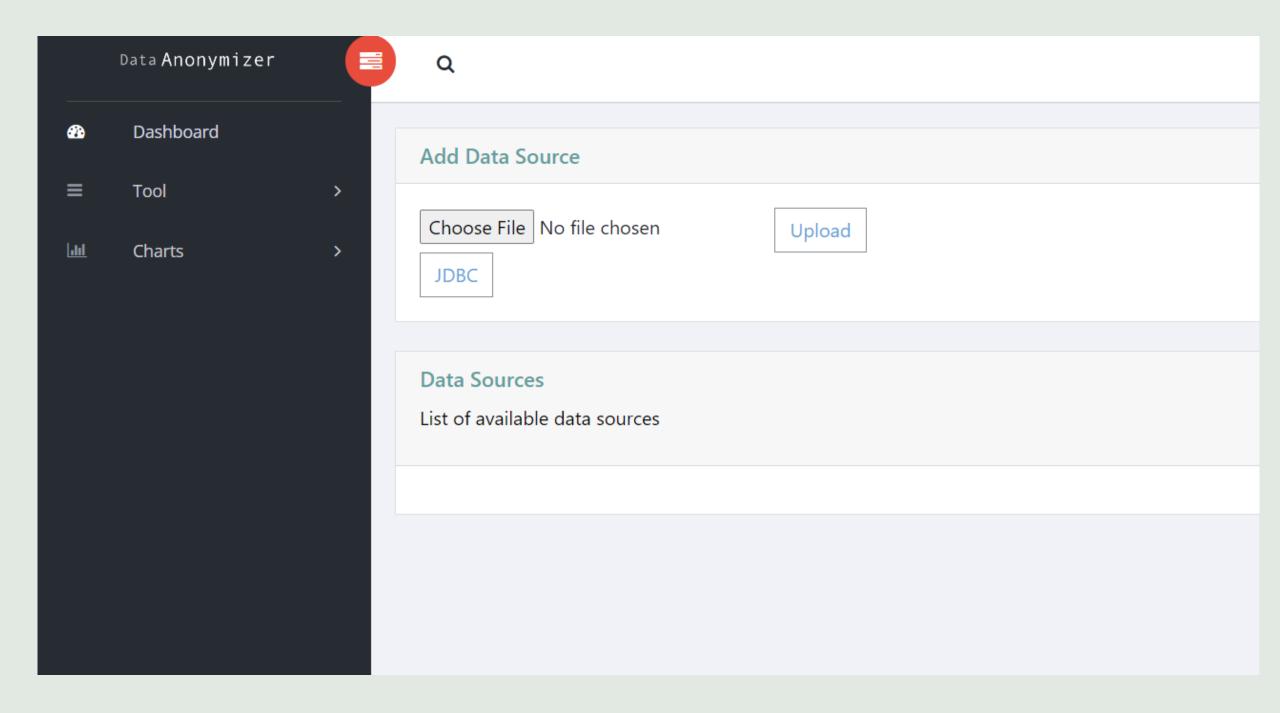
3.4. Project Folder Structure and Dependencies The data that is uploaded to the server or data anonymizer will reside in a directory named Anonymization_resources of root folder Inside the directory it again contains folders named data sources, configurations, hierarchies for their respective files

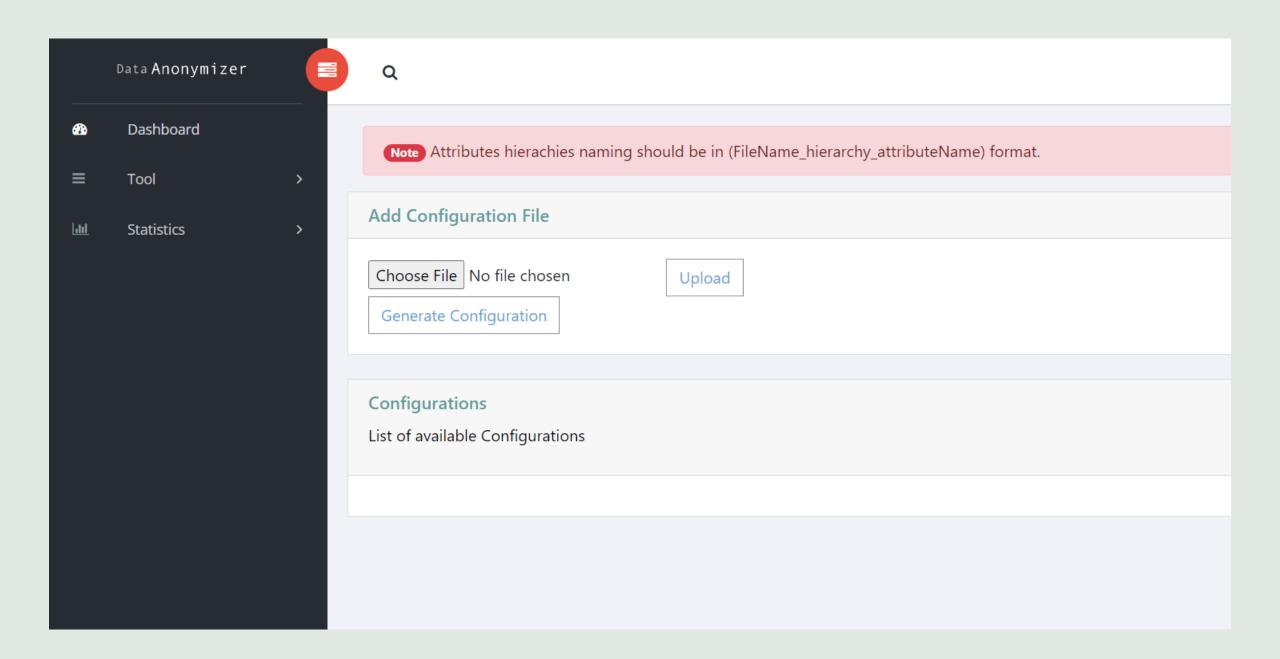
Also the output folder contains the recent anonymized file

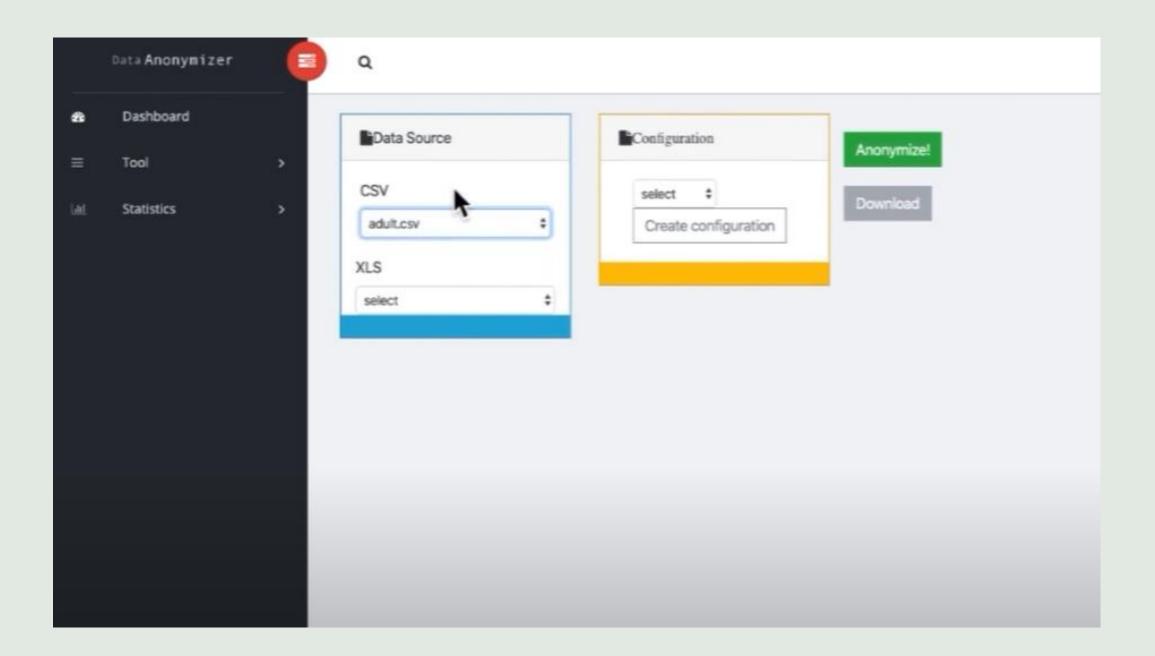
4. Interface/User Manual

Data Sources: By clicking the data sources card the user will be redirected to the data sources page Configurations: By clicking the configurations card from the index page, the user will be redirected towards configurations page









4.3. Privacy Model

Typically most of the anonymization algorithms are based on Generalization and Suppression of QID attributes

They are

k-Anonymity



This algorithm requires that each QID tuple appear in at least k-1 records, this will ensure at minimal, that released data processed with k-anonymity will be difficult to re-identify



The QIDs contains information that is more likely to find over the dataset, so this type of attribute is more vulnerable to re-identification



k-anonymity uses generalization and suppression methods

1-diversity



This algorithm requires a high entropy on the distribution of SAs for each QID



The main idea behind I-diversity, is the well balancing dispersion of SAs between all the groups included on the datasets



Vulnerability- Similarity attacks, that happens because it considers the diversity of SAs in the group, but it is not concerned with the semantic proximity of the values

t-closeness



Due to the fact that previous privacy models had some vulnerabilities, a new one emerged, the tcloseness algorithm



It requires that the distribution of a SA in any equivalence class must be similar to the attributes distribution in the overall dataset, this way, the chances of learning individual's information are lower



In order to introduce and manage gaps between values of SAs, tcloseness uses the Earth Mover Distance metric, receiving a precise distance between the two distributions

Hierarchies



Defining the intervals is important in the configuration process, all the solutions below use this mechanism, but sometimes is not clear what their function is, neither the right way to configure them



Hierarchies are normally used for categorical attributes, to increase the utility of anonymized datasets, "categorisation" is often combined with tuple suppression, i.e., data records inconsistent with privacy criteria are automatically removed from the dataset



Data and generalization hierarchies can be imported from many different types, providing compatibility with a wide range of data processing tools

Hierarchies

 In this tool the hierarchies can be supplied by the user in the form of 'CSV' files, which each file contains the column headers specifying the level fo the hierarchy and rows contain intervals or generalised attributes



	Α	В	С)	Е		F	G	Н	1		J	K
1	sex;age;ra	ce;marital	l-status;e	education	;native-	country	y;work	class;c	occupatio	n;salary-cl	ass			
2	Male;39;W	Vhite;Neve	er-marrie	ed;Bachel	ors;Unit	ed-Stat	tes;Sta	te-gov	;Adm-cle	rical;<=50k				
3	Male;50;W	/hite;Marr	ried-civ-s	pouse;Ba	chelors	;United	l-State	s;Self-	emp-not-	inc;Exec-m	anageria	l;<=50	K	
4	Male;38;W	/hite;Divo	rced;HS-	grad;Unit	ed-State	es;Priva	ite;Hai	ndlers-	cleaners;	<=50K				
5	Male;53;Bl	lack;Marri	ied-civ-sp	pouse;11t	h;Unite	d-State	s;Priva	ate;Hai	ndlers-cle	aners;<=50	K			
6	Female;28	;Black;Ma	rried-civ	-spouse;E	Bacheloi	rs;Cuba	;Privat	te;Prof	-specialty	;<=50K				
7	Female;37	;White;Ma	arried-civ	v-spouse;	Masters	;United	d-State	es;Priva	ate;Exec-ı	nanageria	;<=50K			
8	Female;49	;Black;Ma	rried-spo	ouse-abse	nt;9th;J	lamaica	;Priva	te;Oth	er-service	;<=50K				
9	Male;52;W	/hite;Marr	ried-civ-s	pouse;HS	-grad;U	nited-S	tates;	Self-en	np-not-in	;Exec-ma	nagerial;>	50K		
10	Female;31	;White;Ne	ever-mar	ried;Mast	ers;Uni	ted-Sta	tes;Pri	ivate;P	rof-specia	alty;>50K				
11	Male;42;W	/hite;Marr	ried-civ-s	pouse;Ba	chelors	;United	l-State	s;Priva	te;Exec-r	nanagerial	>50K			
12	Male;37;Bl	lack;Marri	ied-civ-sp	pouse;Sor	ne-colle	ege;Uni	ted-St	ates;Pr	rivate;Exe	c-manage	ial;>50K			
13	Male;30;A	sian-Pac-I	slander;I	Married-c	iv-spous	se;Bach	elors;	India;S	tate-gov;	Prof-speci	alty;>50K			
14	Female;23	;White;Ne	ever-mar	ried;Bach	elors;U	nited-S1	tates;P	Private	;Adm-cler	ical;<=50K				
15	Male;32;Bl	lack;Neve	r-marrie	d;Assoc-a	cdm;Un	ited-St	ates;P	rivate;	Sales;<=5	OK				
16	Male;34;A	mer-India	n-Eskimo	;Married	-civ-spo	use;7th	n-8th;N	Mexico	;Private;T	ransport-r	noving;<=	50K		
17	Male;25;W	Vhite;Neve	er-marrie	ed;HS-grad	d;United	d-States	s;Self-e	emp-no	ot-inc;Far	ming-fishir	g;<=50K			
18	Male;32;W	Vhite;Neve	er-marrie	ed;HS-grad	d;United	d-States	s;Priva	te;Mad	chine-op-	nspct;<=5	OK			
19	Male;38;W	/hite;Marr	ried-civ-s	pouse;11	th;Unite	ed-State	es;Priv	ate;Sa	les;<=50K					
20	Female;43	;White;Div	vorced;N	/lasters;U	nited-St	ates;Se	lf-emp	p-not-i	nc;Exec-n	nanagerial	>50K			
21	Male;40;W	/hite;Marr	ried-civ-s	pouse;Do	ctorate	;United	d-State	es;Priva	ate;Prof-s	pecialty;>5	OK			
22	Female;54	;Black;Sep	oarated;F	HS-grad;U	nited-St	tates;Pr	ivate;	Other-	service;<=	50K				
23	Male;35;Bl	lack;Marri	ied-civ-sp	pouse;9th	;United	-States	;Feder	al-gov	;Farming-	fishing;<=5	OK			
24	Male;43;W	/hite;Marr	ried-civ-s	pouse;11	th;Unite	ed-State	es;Priv	ate;Tra	ansport-n	noving;<=5	OK			
25	Female;59	;White;Div	vorced;H	IS-grad;Ur	nited-St	ates;Pri	ivate;T	Tech-su	ipport;<=	50K				
26	Male;56;W	/hite;Marr	ried-civ-s	pouse;Ba	chelors	;United	l-State	s;Loca	l-gov;Tec	h-support;	>50K			
27	Male;19;W	Vhite;Neve	er-marrie	ed;HS-grad	d;United	l-States	s;Priva	te;Craf	ft-repair;<	=50K				
28	Male;39;W	/hite;Divo	rced;HS-	grad;Unit	ed-State	es;Priva	te;Exe	ec-man	agerial;<	=50K				
_29	Male:49:W	/hite:Marr adult	ried-civ-s +	spouse:HS	-grad:U	nited-S	tates:	Private	:Craft-rei	air:<=50K				

Adult hierarchy age

Adult hierarchy native country

		Addit meraleny hative country
A B C	_	
1 1;0-4;0-9;0-19;*	1	United-States;North America;*
2 2;0-4;0-9;0-19;*	2	Cambodia;Asia;*
3 3;0-4;0-9;0-19;*	3	England;Europe;*
4 4;0-4;0-9;0-19;*	4	Puerto-Rico;North America;*
5 5;0-4;0-9;0-19;*	5	Canada;North America;*
	6	Germany;Europe;*
7 7;5-9;0-9;0-19;*	7	Outlying-US(Guam-USVI-etc);North America;*
8 8;5-9;0-9;0-19;*	8	India;Asia;*
9 9;5-9;0-9;0-19;*	9	Japan;Asia;*
10 10;5-9;0-9;0-19;*		Greece;Europe;*
11 11;10-14;10-19;0-19;*		South;Africa;*
12 12;10-14;10-19;0-19;*	12	China;Asia;*
13 13;10-14;10-19;0-19;*		Cuba;North America;*
14 14;10-14;10-19;0-19;*		Iran;Asia;*
15 15;10-14;10-19;0-19;*		Honduras;North America;*
16 16;15-19;10-19;0-19;*		Philippines;Asia;*
17 17;15-19;10-19;0-19;*		Italy;Europe;*
18 18;15-19;10-19;0-19;*		Poland;Europe;*
19 19;15-19;10-19;0-19;*		Jamaica;North America;*
20 20;15-19;10-19;0-19;*		Vietnam;Asia;*
21 21;20-24;20-29;20-39;*		Mexico;North America;*
22 22;20-24;20-29;20-39;*		Portugal;Europe;*
23 23;20-24;20-29;20-39;*		Ireland;Europe;*
24 24;20-24;20-29;20-39;*		France;Europe;*
25 25;20-24;20-29;20-39;*		Dominican-Republic;North America;*
26 26;25-29;20-29;20-39;*		Laos;Asia;*
27 27;25-29;20-29;20-39;*		Ecuador;South America;*
		Taiwan;Asia;*
28 28;25-29;20-29;20-39;*	29	Haiti:North America:*
29 29:25-29:20-29:20-39:*		adult_hierarchy_native-country

Attack scenarios

Identity Re-identification:

- Attack Scenario: Malicious actors attempt to re-identify individuals in the anonymized dataset to discover sensitive information.
- Prevention: Techniques like k-anonymity, l-diversity, and t-closeness introduce generalization and suppression methods, making it challenging to identify specific individuals.

Homogeneity Attacks:

- Attack Scenario: Attackers exploit patterns in the data to reveal that a particular attribute is the same for a group of individuals.
- Prevention: k-Anonymity helps protect against homogeneity attacks by ensuring that each quasi-identifier tuple appears in at least k-1 records, reducing the risk of revealing identical sensitive information.

4.4. Anonymize



When the user has both Data source and configuration selected or in place, then by clicking the anonymize button the anonymization process can be initiated

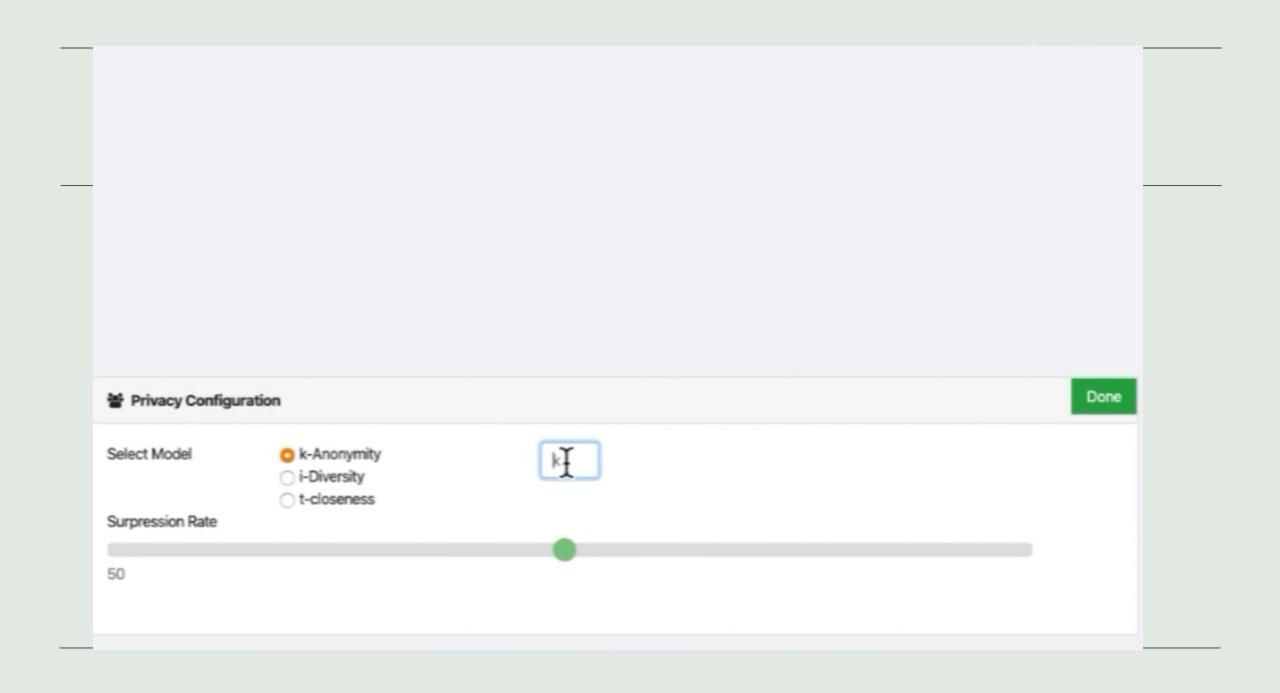


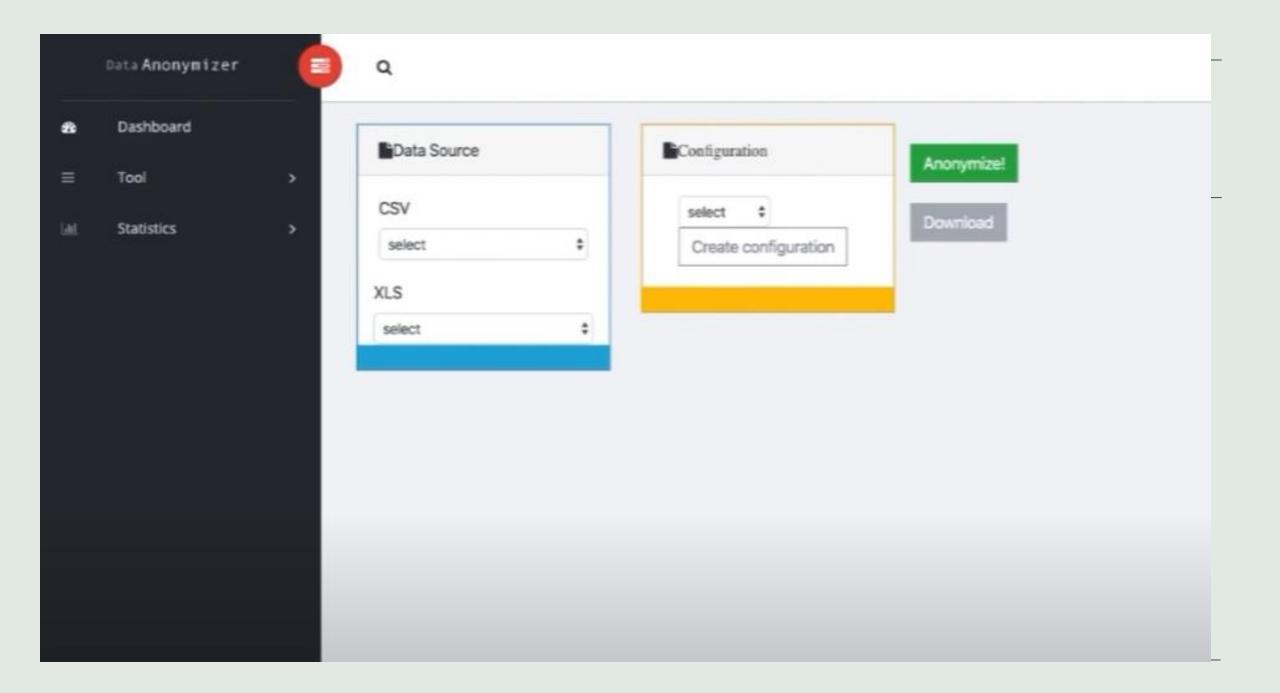
After the anonymization the user can download the anonymized data by clicking the 'Download' button which is below on the 'anonymize' button



After the anonymization process the Data Source file will be automatically deleted from the server in order for the system to maintain privacy

Configurtion Table Show 10 * entries Search: Attribute Name AttributeType DataType Heirarchy Choose file No file chosen / Please select age Please select Insensitive education Identifying Choose file No file chosen Please select Quasi identifying marital-status Choose file No file chosen Please select Please select native-country Choose file No file chosen Please select Please select Choose file No file chosen occupation Please select Please select Choose file No file chosen race Please select Please select salary-class Choose file No file chosen Please select Please select Choose file No file chosen sex Please select Please select workclass Choose file No file chosen Please select Please select





	Α	В	С		D	Е	F	G	Н		
1	age	sex	native-	cou	marital-sta	race	education	salary-clas	workclass	occupatio	n
2	20-39	Male	North A	Αme	Never-mai	White	Higher edu	<=50K	State-gov	Other	
3	40-59	Male	North A	Αme	Married-ci	White	Higher edu	<=50K	Self-emp-r	Nontechni	ical
4	20-39	Male	North A	Αme	Divorced	White	Secondary	<=50K	Private	Nontechni	ical
5	40-59	Male	North A	Αme	Married-ci	Black	Secondary	<=50K	Private	Nontechni	ical
6	*	*	*		*	*	*	*	*	*	
7	20-39	Female	North A	Αme	Married-ci	White	Higher edu	<=50K	Private	Nontechni	ical
8	40-59	Female	North A	Αme	Married-sp	Black	Secondary	<=50K	Private	Other	
9	40-59	Male	North A	Αme	Married-ci	White	Secondary	>50K	Self-emp-r	Nontechni	ical
10	20-39	Female	North A	Αme	Never-mai	White	Higher edu	>50K	Private	Technical	
11	40-59	Male	North A	Αme	Married-ci	White	Higher edu	>50K	Private	Nontechni	ical
12	20-39	Male	North A	Αme	Married-ci	Black	Higher edu	>50K	Private	Nontechni	ical
13	20-39	Male	Asia		Married-ci	Asian-Pac-	Higher edu	>50K	State-gov	Technical	
14	20-39	Female	North A	Αme	Never-mai	White	Higher edu	<=50K	Private	Other	
15	20-39	Male	North A	Αme	Never-mai	Black	Higher edu	<=50K	Private	Nontechni	ical
16	20-39	Male	North A	Αme	Married-ci	Amer-India	Secondary	<=50K	Private	Other	
17	20-39	Male	North A	Αme	Never-mai	White	Secondary	<=50K	Self-emp-r	Other	
18	20-39	Male	North A	Αme	Never-mai	White	Secondary	<=50K	Private	Technical	
19	20-39	Male	North A	Αme	Married-ci	White	Secondary	<=50K	Private	Nontechni	ical
20	*	*	*		*	*	*	*	*	*	
21	20-39	Male	North A	Αme	Married-ci	White	Higher edu	>50K	Private	Technical	
22	40-59	Female	North A	Αme	Separated	Black	Secondary	<=50K	Private	Other	
23	20-39	Male	North A	Αme	Married-ci	Black	Secondary	<=50K	Federal-go	Other	
24	40-59	Male	North A	Αme	Married-ci	White	Secondary	<=50K	Private	Other	
25	40-59	Female	North A	Αme	Divorced	White	Secondary	<=50K	Private	Technical	
26	40-59	Male	North A	Αme	Married-ci	White	Higher edu	>50K	Local-gov	Technical	
27	0-19	Male	North A	Αme	Never-mai	White	Secondary	<=50K	Private	Technical	
28	20-39	Male	North A	Αme	Divorced	White	Secondary	<=50K	Private	Nontechni	ical
29	40-59	Male	North A	Αme	Married-ci	White	Secondary	<=50K	Private	Technical	

5. Future work



Improving Interface design- The design of this version is composed of single wizard which contains all the features



Database storage- In this version the files uploaded by the user are being stored in a local storage due to ARX API in the server which is a kind of limitation



Input and output views- When the file has been uploaded by the user, there is no option in the present system to view the input



Hierarchy Creation- Similar to creating the configuration there can be an interface for each attribute to specify the intervals or levels or generalisations



ARX - Open Source Data Anonymization Software



A Web Anonymizer Platform for Datasets with Personal Information