# Laboratory on Privacy Total: 100 points

This lab is to be completed in pairs of two students. You will need to hand in one report per group.

Names	<b>:</b>
Date:	
•	tive: In this lab, you will gain understanding about multiple types of anonymization techniques and adeoffs. Additionally, you will learn to anonymize a dataset using k-anonymity and l-diversity.
1.5 or l	To perform this lab you will need to have a 32bit Windows and java run time environment (JRE) nigher. Additionally, you will need the software provided in the privacy-lab.zip that contains the Anoymization toolbox and the datasets that you are going to use.
Note: 7 the lab	The output of anonymization toolbox is provided. You can just use the anonymized files to complete .
Prelim	inary questions:
1.	In your own words, define what a <i>quasi-identifier</i> is
2.	Explain the difference between a <i>quasi-identifier</i> and an <i>identifier field</i>
3.	What is an equivalence class?
4.	Provide a brief overview of <i>k-anonymity</i>

Provide a brief overview of <i>L-diversity</i>
Explain the difference between $L$ -diversity and $K$ -anonymity
on questions:
: Naïve anonymization
Unzip the privacy-lab.zip file in your computer. Go to folder PartA and open the file <b>naive.txt.</b> For your convenience, each column is separated using tabs. This file contains a naively anonymized dataset. This means that the name, last name and address of people in the datase have been removed. If you know that Alice Smith is divorced and works in <i>Business_and_repair_services</i> , are you able to identify her record? If so, please write it below.
According to the previous question, do you conclude this method of anonymization guarantees any privacy?

#### Part

In this part of the lab, you will learn to anonymize datasets using different anonymization techniques. You will use the UTD anonymization toolbox provided in the zip file.

First, you will compare the effect of using different values for parameter in the anonymized data.

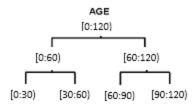
1. Open the file header.txt stored in the folder /dataset. This file contains the list of names of the fields of the dataset that you are going to analyze. List the first five fields of this file:

Open the file original-data.cvs this is the data file, each column in the data file corresponds to the fields that were described in **header.txt**. For the first record on this file, show the mapping between the header name and the data. E.g., provide the name of the column and the value of

2. Open the file config.xml. This file contains the specifications that will be used to anonymize the dataset specified in the tag <input filename>. Below, there is an example of such specification. In the example, k=1, the identifier is <id>"Zipcode" and the quasi-identifiers <qid> used are age and sex. The tag method specifies the type of anonymization technique use and the tag <ouput filename> is used to specify the output name of the anonymized data generated by the toolkit. For your convenience, these tags are underlined and shown in blue.

```
<?xml version="1.0"?>
<config method = 'Datafly' k = '6'>
    <input filename='dataset/original-data.cvs' separator=','/> <!-- If left blank, separator will</pre>
    be set as comma by default.-->
    <output filename='anonimized-k=6.txt' format ='genVals'/> <!-- Format</pre>
   options = {genVals, genValsDist, anatomy}.
        If left blank, output format will be set as genVals by default.-->
    <id><!-- List of identifier attributes, if any, these will be excluded from the output -->
        <att index='0' name='Zipcode'/>
    </id>
    <qid>
        <att index='7' name='sex'>
            <map> <!-- Mapping of a categorical domain to a discrete-valued numeric domain.-->
                <entry cat='Female' int='0' />
                <entry cat='Male' int='1' />
            <vgh value='[0:1]'>
            </vgh>
        </att>
        <att index='1' name ='age'>
            <vgh value='[0:120)'>
                <node value='[0:60)'>
                    <node value='[0:30)'/>
                    <node value='[30:60)'/>
                </node>
                <node value='[60:120)'>
                    <node value='[60:90)'/>
                    <node value='[90:120)'/>
                </node>
            </vgh>
        </att>
    </qid>
</config>
```

The age in the previous xml is generalized using the tree



- 3. From the provided xml, infer the generalization tree used for *age* and draw it in the flowing space.
- 4. Analyze the anonymization case "k=1"

The result is provided in the output folder, see the **anonymized-k=1.txt** file.

Option:

If you want to generate the anonymized file by yourself, please follow the instructions: cd to the appropriate directory and type the following in the command line java -cp sqlite.jar;anon\_toolbox.jar;. -Djava.library.path=. anonymizer.Anonymizer you will anonymize the original-data.cvs using the configuration provided in configural.

you will anonymize the **original-data.cvs** using the configuration provided in **config.xml**. After running the command, the **anonymized-k=1.txt** file is created in the directory.

Open the file. Are you able to see the zipcode of the users in the dataset?

Why?

- 5. List the first three ages of the anonymized data
- 6. Compare the original dataset with this anonymized dataset. Was the age modified substantially?
- 7. By exploring the anonymized data, given the configuration provided in **config.xml**, why was the age column modified as it was?
- 8. Plot the histogram of field age of the anonymized data with 4 bins ([0:25), [25:50), [50:75), and [75:100)).

### Tests for K=4

	Now, you need to repeat the experiment above, but changing the configuration file <b>config.xml</b> so that k=4, and then the input file is the output filename to <b>anonymized-k=4.txt.</b>		
	The result file is also provided in the output folder.		
	Open the file that was created by this action. Is the age column different from the previous anonymized data?_Why?		
10.	List the first three ages of the anonymized output		
11.	Now you need to find the equivalence classes in the anonymized data. List them and clearly separate them in the following space		
	<u> </u>		
12.	By looking uniquely the anonymized dataset (without considering the original dataset) are you able to locate the record of Anne who is 45 years old and is a widow?		
13.	If yes, list Anne's anonymized record below.		
14.	If no, why you could not identify her record?		
	<del></del>		

15.	By looking uniquely the anonymized dataset (without considering the original dataset) are you able to locate the record of Samantha who was born in Cuban?
16.	If yes, list Samantha's anonymized record below.
17.	If no, why you could not identify her record?
18.	What can you conclude about the answers of questions 12 to 17?
	ts for K=6  Now, you need to change the configuration file config.xml so that k=6 and the output filename
	to anonymized-k=6.txt.
	The result file is also provided in the output folder.
	Open the file that was created by this action. Is the age column different from the previous anonymized data?
20.	List the first three ages of the anonymized data
21.	Now you need to find the equivalence classes in the anonymized data. List them and clearly separate them in the following space

22. Are the equivalence classes different from the ones you found in point 11?	
23. Inspect uniquely the anonymized dataset (without considering the original dataset) are you a to locate the record of Anne who is 45 years old and is a widow?	able
If yes, list the record you identified as Anne's record. If not, explain why.	
24. By looking uniquely the anonymized dataset (without considering the original dataset) are you able to locate the record of Alice who is 30 years old and divorced?	
25. If yes, list Alice's anonymized record below. If not, explain why.	_
26. Looking uniquely the anonymized dataset (without considering the original dataset), are you	
able to locate the record of Pamela who is 30 years old, single?	
27. If yes, list her anonymized record below. If not, explain why.	
<ul> <li>28. You know that Jessica has never worked and she is 24. Looking uniquely the anonymized dat (without considering the original dataset), are you able to locate her record?</li> <li>29. If yes, list her anonymized record below. If not, explain why.</li> </ul>	taset
30. In this part of the lab, you have anonymized multiple datasets using k-anonymity. According your results, what value of k was less secure?and what value of k seem provide more protection?	
31. Now find the histogram of the age with two bins ([0:50) and [50:100)). Draw it below.	
32. Compare the histograms for k=1, k=4 and k=6. What k value produces the closest result to the original dataset?Why?	ne
33. Which one the worst?Why?	

# Comparing L-diversity with k-anonymity

In the following, you will compare k-anonymity with l-diversity.

l.	Does L-diversity consider the sensitivity of the parameters to be anonymized?
2.	Take the output for k=6 that you previously found. Does the output respect l-diversity if the sensitive attribute is the field place where the person was born and l=2?If it does not fulfill L-diversity. What records violate this property? List them below.
3.	If you know George mother is from Mexico and his father is form the US. Are you able to identify the record of George in the k=6 output in point 16?if so, list his record
4. 5.	Would I-diversity help to avoid this type of inference?  Describe how you could avoid this problem
6.	Anonymize the dataset manually to ensure that 1-diversity is fulfilled with k=4, 1=2 considering the <i>country_of_birth_self</i> as a sensitive attribute. Show the equivalence classes of your solution below

# Part C: Analyzing datasets

•	Go to the folder /dataset. In a text editor open the file header.txt. This file contains the list of the names of the fields of the dataset that you are going to analyze. List the fields that you think correspond to <i>identifiers</i> in the dataset and state why you think they are <i>identifiers</i>
•	Open the file /dataset/ anonymized-part3.txt List the first two rows of the dataset. Were any fields used during anonymization as identifiers? If yes, which?
	Why were you able to infer if any fields were used as identifiers during the anonymization procedure?
	The age and sex were used as a quasi-identifier in this anonymization. Their generalization tree are shown in the xml presented at the beginning of this document. By exploring the anonymized dataset, what was the $k$ specified?

5. Show the groups of equivalence classes in the following space.