# Security, Privacy and Explainability in Machine Learning

SS 2019

## Exercise 1

| Name | Mat.Nummer |
| --- | --- |
| Maximilian Bachl | 01100143 |
| Alexander Hartl | 01125115 |

May 11, 2019

# 1 Dataset Selection

We selected the "Airline Customers" dataset available on kaggle's online platform[1], which seems to be taken from an airline company's frequent-flyer program. Since a customer's bonus points presumably show a strong correlation with the customer's financial wealth, the dataset contains sensitive data.

## 1.1 Contained Information

The dataset contains several, partially redundant, attributes regarding a customer's flight frequency and bonus points. For demonstration purposes we reduced the attributes to

- Gender

- Age

- Work Country

- Work Province

- Work City

- First flight date

- Total flight count

- Total bonus points

## 1.2 Quasi-Identifiers

The combination of quasi-identifying attributes allows a third-party to uniquely identify an individual in the dataset. In our case the set of quasi-identifiers consists of **Gender, Age and Work Country/Province/City**. Other attributes could contribute to identifying an individual but are usually not known to a third-party.

# 2 Anonymization

We used ARX[2] for the anonymization process and deployed several anonymization techniques for the different attributes:

- Gender: Generalization, suppressing the attribute, if necessary.

---

[1] `https://www.kaggle.com/diezerg/airline-customers`
[2] `https://arx.deidentifier.org/`

Figure 1: Exemplary equivalence classes as found by ARX.

- Age: Microaggregation, reducing the exact age to intervals, if necessary.

- Work Location: Generalization, using a hierarchy to mask the data, if necessary. For this to work properly, we merged the work location attributes to one attribute in the form <Work Country> / <Work Province> / <Work City> as a preprocessing step, so that more fine-grained information will be suppressed first.
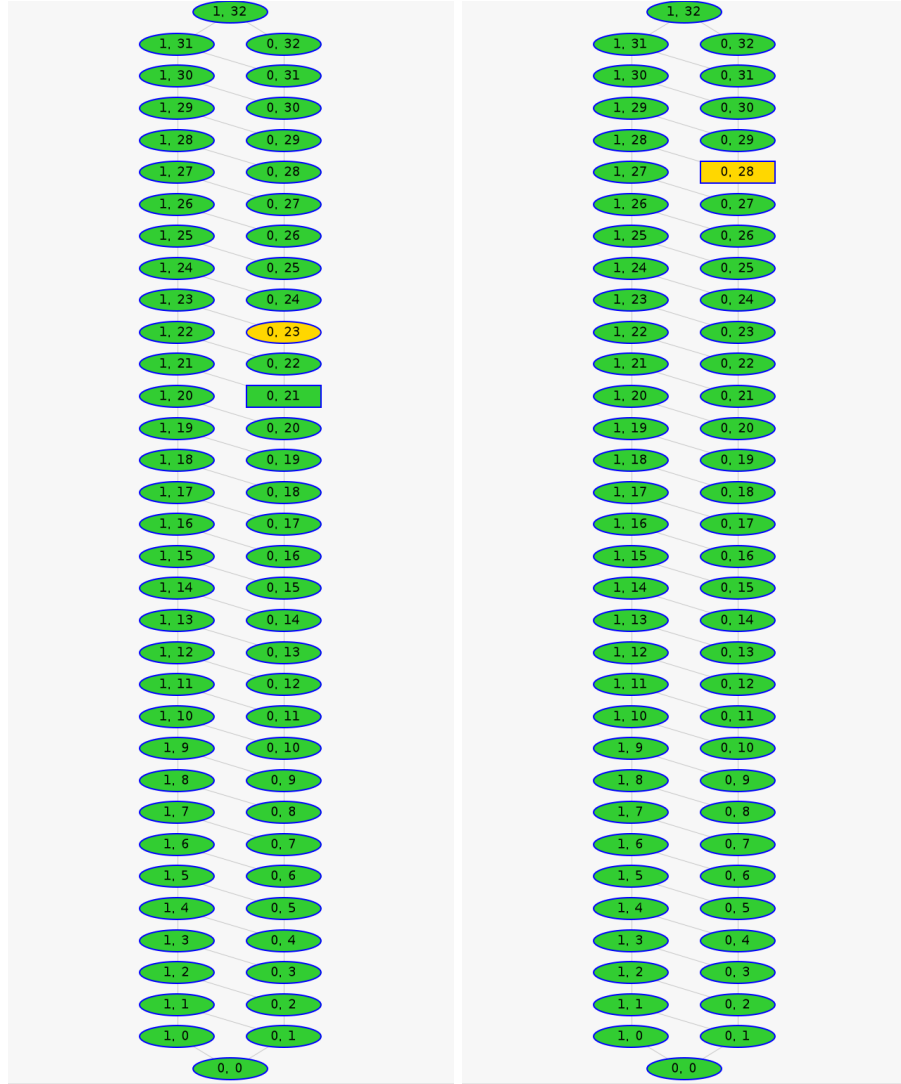
We used $k$-Anonymity with $k = 3$ and $k = 50$ as privacy model. Furthermore, as total bonus points constitutes a sensitive attribute, we used $t$-closeness with $t = 0.2$ and "EMD with ordered distance" as distance measure.

# 3 Discussion

Fig. 2 shows the anonymization paths that ARX finds for $k = 3$ and $k = 50$, respectively. To achieve the same level of anonymity, either the gender attribute can be suppressed or the contained information about the work location can be reduced, resulting in the paths shown in Fig. 2.

For different $k$, the paths don't show significant differences. However, aiming at the same level of anonymity, ARX has to suppress more information when selecting a lower $k$.

Fig. 1 shows an excerpt of an solution ARX finds for $k = 3$. Some of the equivalence classes indeed assume the desired size of 3, while others are substantially larger.

Figure 2: Anonymization paths for $k = 3$ (left) and $k = 50$ (right).