

Privacy-Preserving Data Publishing

Assignment #4

Diogo Cordeiro (up201705417) Hugo Sales (up201704178)

2022/06/02

Attribute classification

We classified the attributes as follows:

Table 1: Attribute classifications

Attribute	Classification
age	QID
workclass	Insensitive
fnlwgt	Insensitive
education	QID
education-num	QID
marital-status	QID
occupation	QID
relationship	QID
race	QID
sex	QID
capital-gain	Sensitive
capital-loss	Sensitive
hours-per-week	QID
native-country	Insensitive
prediction	Insensitive

Justifications

The vast majority of attributes present extremely low values of distinction. We speculate this may be an TODO

age

According to HIPPA recommendations, and together with it's very high separation value (99.87%), this attribute is classified as a QID.

workclass

This attribute presents a relatively low separation value (49.71%), and given how generic it is, it's deemed Insensitive.

fnlwgt

Despite high values of distinction (66.48%) and separation (99.99%) the **fnlwgt** column is not a QID because it represents a weight, not a count of individuals in the same equivalence class in the original dataset. This can be seen with the results below. Additionally, it's not easily connected to another auxiliary info dataset.

```
$ tail -n '+2' adult_data.csv | awk -F',' '{count[$10] += $3;} \
    END {for(sex in count){print sex, count[sex]}}'
```

Resulting in:

Table 2: Sum of **fnlwgt** for each **sex**

Sex	Sum
Female	2000673518
Male	4178699874

The sum of these values is 6,179,373,392. This value is much larger than the population of the U.S.A., the origin of the dataset, which implies this attribute is not a count, as stated.

We also note there are substantially more Male than Female records (more than double the **fnlwgt**).

education

This attribute presents a separation of 80.96%, which is quite high, so this attribute is classified as a QID.

education-num

We exported the anonymized dataset and used the following command to verify there weren't any discrepancies between the **education** and **education-num** columns:

```
$ cat anonymized.csv | sed -r 's/,([\^ ])/\t\1/g' | awk -F';' '{print $5, $4}' | sort -un
```

Since there was a one-to-one mapping, we concluded this was just a representation of the **education** attribute. As such, this attribute receives the same classification, which is backed by the equally high separation value of 80.96%, so it's qualified as a QID.

marital-status

With a relatively high separation value of 66.01%, together with the fact that it could be cross referenced with other available datasets, we classify this attribute as a QID.

occupation

With a separation of 90.02%, this attribute is classified as a QID.

relationship

Given it's separation value of 73.21%, this attribute is classified as a QID.

race

This column presents some weirdly specified values (Amer-Indian-Eskimo), but has a separation of 25.98%; given the fact that this attribute could be cross referenced with other datasets, it is classified as a QID, so it may be transformed into more generic values.

sex

Despite the low separation value of 44.27%, this attribute is canonically classified as a QID, since it can be easily cross referenced with other datasets.

We noted this dataset seems to more males than females. See tbl. 2 and the following table

Table 3: Number of records with each **education** for each **sex**

education	Female	Male
Preschool	16	35
1st-4th	46	122
5th-6th	84	249
7th-8th	160	486
9th	144	370
10th	295	638
11th	432	743
12th	144	289
HS-grad	3390	7111
Some-college	2806	4485
Assoc-voc	500	882
Assoc-acdm	421	646
Bachelors	1619	3736
Masters	536	1187
Prof-school	92	484

education	Female	Male
Doctorate	86	327

capital-gain & capital-loss

With a separation of 15.93% and 9.15% respectively, these attributes are not QIDs. They're qualified as Sensitive, as the individuals may not want their capital gains and losses publicly known.

A t-closeness privacy model was chosen for these attributes, with a value of t of 0.2. This reasoning is discussed in [Applying anonymization models > k-Anonymity > Effect of parameters](#)

hours-per-week

This attribute has a relatively high separation (76.24%) and since it had really unique values, it could be cross referenced with another dataset to help identify individuals, so it's classified as QID.

native-country

While this attribute might be regarded as a QID, it presents really low separation values (19.65%) in this dataset, so it's qualified as Insensitive.

prediction

This is the target attribute, the attribute the other attributes predict, and is therefore Insensitive.

Privacy risks in the original dataset

In the original dataset, nearly 40% of records have a more than 50% risk of re-identification by a prosecutor. In general, we see a stepped distribution of the record risk, which indicates some privacy model was already applied to the dataset, however to a different standard than what we intend.

All records had really high uniqueness percentage even for small sampling factors, according to the Zayatz, Pitman and Dankar methods. Only SNB indicated a low uniqueness percentage for sampling factors under 90%. What this means, is that with a fraction of the original dataset, a very significant number of records was sufficiently unique that it could be distinguished among the rest, which means it's potentially easier to re-identify the individuals in question.

All attacker models show a success rate of more than 50%, which is not acceptable.

Applying anonymization models

k-Anonymity

We opted for 8-anonymity, for its tradeoff between maximal risk and suppression. t-closeness was chosen for `capital-gain` and `capital-loss` (sensitive attributes).

Re-identification risk

The average re-identification risk dropped to nearly 0%, whereas the maximal risk dropped to 12.5%. The success rate for all attacker models was reduced drastically, to 1.3%.

Utility

The original Classification Performance, a measure of how well the attributes predict the target variable (`prediction`) was 83.24% and it remained at 82.45%. 10.07% of attributes are missing from the anonymized dataset. This value being equal across all attributes suggests entire rows were removed, rather than select values from separate rows. The only exception is the `occupation` attribute, which was entirely removed.

Effect of parameters

At a suppression limit of 0%, the same accuracy is maintained, but the vast majority of QIDs are entirely removed.

At a suppression limit of 5%, roughly the same prediction accuracy is maintained, with around 4.5% of values missing, however with really high Generalization Intensity values for some attributes (e.g. 95.42% for `sex`, 93.87% for `race` and 91.47% for `education` and `education-num`). `occupation` was entirely removed.

At a suppression limit of 10%, the prediction accuracy is maintained, with around 9.8% of values missing. However, the Gen. Intensity drops to around 90%.

At a suppression limit of 20%, accuracy is maintained, once again, with around 10% of values missing, indicating this would be the optimal settings, as the same results are achieved with a limit of 100%.

At a t-closeness for `capital-gain` and `capital-loss` t value of 0.001 (the default), anonymization fails, not producing any output.

At a t value of 0.01, accuracy drops to 75% and most attributes have missing values of 100%.

At a t value of 0.1, classification accuracy is nearly 81%, but missings values are around 20%.

At a t value of 0.2, the chosen value, the accuracy is 82.5% with lower Gen. Intensity values.

At a t value of 0.5, the classification accuracy goes to 82.2% with increased Generalization Intensity values.

Adjusting the coding model had no significant effects.

(ϵ, δ) -Differential Privacy

With the default ϵ value of 2 and a δ value of 10^{-6} , the performance was really good.

Re-identification risk

All indicators for risk by each attacker model was between 0.1% and 0.9%.

Utility

The original Classification Performance was 83.24% and it remained at 80.97%.

Nearly 16% of attributes are missing, with the exception of `age` and `education-num`, which are 100% missing.

Effect of parameters

An ϵ value of 3 maintained the accuracy at 80.5% with missings values rounding 32%.

An increase of δ to 10^{-5} resulted in a classification performance of 82.05% and a missings value of 21.02% for all attributes.

A further increase of δ to 10^{-4} resulted in an increased accuracy of 82.32%, but a maximal risk of 1.25%.

Results

The 8-anonymity model was chosen as it resulted in a broader distribution of attribute values like `age`, whereas with Differential Privacy, they were split into only 2 categories.

Observations

We noted that the contingency between `sex` and `relationship` maintained the same distribution after anonymization, meaning that these changes don't mean `relationship` can identify an individual's `sex` any more than in the original dataset.

With the following commands, we noted some possible errors in the original dataset, where the `sex` and `relationship` attributes didn't map entirely one to one: there was one occurrence of (Husband, Female) and two of (Wife, Male). It's possible this is an error in the original dataset.

```
$ cat adult_data.csv | tail -n +2 | sed -r 's/,([\ ])/\t\1/g' | cut -d',' -f8,10 | sort | un
```

```
1 Husband, Female
2 Wife, Male
430 Other-relative, Female
551 Other-relative, Male
792 Unmarried, Male
1566 Wife, Female
2245 Own-child, Female
2654 Unmarried, Female
2823 Own-child, Male
3875 Not-in-family, Female
4430 Not-in-family, Male
13192 Husband, Male
```

```
$ cat anonymized.csv | tail -n +2 | sed -r 's/,([\ ])/\t\1/g' | cut -d',' -f8,10 | sort | un
```

```
1295 {Husband, Wife} Female
2264 {Other-relative, Own-child} Female
2981 {Other-relative, Own-child} Male
3280 * *
4391 {Unmarried, Not-in-family} Male
5713 {Unmarried, Not-in-family} Female
12637 {Husband, Wife} Male
```

Since there were occurrences of (Wife, Male), “({Husband, Wife}, Male)” does not undo the transformation of the `relationship` attribute.