

# The Experiment of K-anonymity

M11215032 葉品和

M11215052 陳奕帆

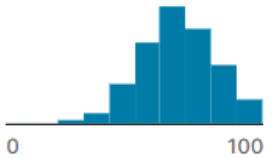
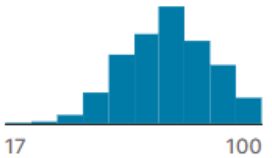
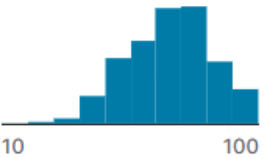
M11215066 鄭宜珊

## 1. Dataset

We use the dataset “**Students Performance in Exams**” from Kaggle. There are a total of 8 attributes. The first 5 attributes are quasi-identifier, and the last 3 attributes are target identifiers.

|          |     |                  |     |                          |     |              |     |
|----------|-----|------------------|-----|--------------------------|-----|--------------|-----|
| ▲ gender | ≡   | ▲ race/ethnicity | ≡   | ▲ parental level of e... | ≡   | ▲ lunch      | ≡   |
| female   | 52% | group C          | 32% | some college             | 23% | standard     | 65% |
| male     | 48% | group D          | 26% | associate's degree       | 22% | free/reduced | 36% |
|          |     | Other (419)      | 42% | Other (552)              | 55% |              |     |

|                         |     |   |     |  |     |   |     |
|-------------------------|-----|---|-----|--|-----|---|-----|
| ▲ test preparation c... | ≡   | # math score  | ≡   | # reading score  | ≡   | # writing score   | ≡   |
| none                    | 64% |  |     |  |     |  |     |
| completed               | 36% | 0   | 100 | 17   | 100 | 10  | 100 |

## 2. K-anonymity Algorithm

In the "race/ethnicity" attribute, there are a total of 5 categories, namely "group A" to "group E." Among them, "group C" and "group D" are the most frequent, accounting for 32% and 26%, respectively. As a result, we have modified the remaining 3 categories to "other" to achieve anonymity.

There are 6 different categories in attribute “parental level of education”. “some college” appears the most frequently(23%), “associate's degree” is the second, it has about 22%. We combine the rest of categories into two combinations. “some high school” and “high school” to “before college”. “bachelor's degree” and “master's degree” to “bachelor's and master's degree”.

### 3. Result Data

The value of K is 9.

|   | gender | race_ethnicity | parental_level_of_education | lunch        | test_preparation_course |
|---|--------|----------------|-----------------------------|--------------|-------------------------|
| 0 | female | group B        | bachelor's degree           | standard     | none                    |
| 1 | female | group C        | some college                | standard     | completed               |
| 2 | female | group B        | master's degree             | standard     | none                    |
| 3 | male   | group A        | associate's degree          | free/reduced | none                    |
| 4 | male   | group C        | some college                | standard     | none                    |
| 5 | female | group B        | associate's degree          | standard     | none                    |

(a) original data

|   | gender | race_ethnicity | parental_level_of_education    | lunch        | test_preparation_course |
|---|--------|----------------|--------------------------------|--------------|-------------------------|
| 0 | female | other          | bachelor's and master's degree | standard     | *                       |
| 1 | female | group C        | some college                   | standard     | *                       |
| 2 | female | other          | bachelor's and master's degree | standard     | *                       |
| 3 | male   | other          | associate's degree             | free/reduced | *                       |
| 4 | male   | group C        | some college                   | standard     | *                       |
| 5 | female | other          | associate's degree             | standard     | *                       |

(b) anonymized data

### 4. Machine Learning Models

In the preprocessing phase, we first drop the attributes “math score” and “test preparation course”. After dropping, we map the attributes “race/ethnicity” and “gender” by one-hot encoding method. We then map the attribute “parental level of education” from 1 to 5 according to the corresponding levels of education. Next, we map the attribute “lunch” with boolean values. Last, we map two “score” attributes to 3 categories each, the bottom 33%, the middle 33%, and the top 33%, with values 1, 2, 3, respectively.

We tried 3 ML algorithms to predict the attribute “writing\_score”, the figure below shows the results of original data and data after anonymization.

## 5. Results

predict target : writing score

|               |          | Accuracy | F1-Score | Recall |
|---------------|----------|----------|----------|--------|
| SVM           | original | 81.50%   | 81.52%   | 81.50% |
|               | after    | 81.50%   | 81.52%   | 81.50% |
| Random Forest | original | 77.60%   | 77.66%   | 77.60% |
|               | after    | 78.70%   | 78.69%   | 78.70% |
| MLP           | original | 81.30%   | 81.39%   | 81.30% |
|               | after    | 81.30%   | 81.35%   | 81.30% |