# Three Preprocessing Approaches to Fairness Correction

## Visualization-Empowered Human-in-the-Loop AI Focus Period

Vladimiro González-Zelaya

19 May 2025

Linköping University — Campus Norrköping

Fairness Notions

Parametrised Data Sampling

Fair and Private Data Correction

Genetic Pipeline Optimisation

# Fairness Notions

## Definition

Feature of a dataset that is prone to an unjustified discriminatory decision

## Examples

- ▶ Race or Skin Colour
- ▶ Sex or Gender
- ▶ Age
- ▶ Income Level
- ▶ Education
- ▶ Nationality

## Definition

Feature of a dataset that is prone to an unjustified discriminatory decision

## Examples

► Race or Skin Colour

► Sex or Gender

► Age

► Income Level

► Education

► Nationality

**Protected Attribute**
- ► Favoured
- ► Unfavoured

**Class**
- ► Positive
- ► Negative

$$U+ \quad F+$$

$$U- \quad F-$$

## Protected Attribute
- ► Favoured
- ► Unfavoured

## Class
- ► Positive
- ► Negative

$$U+ \quad F+$$

$$U- \quad F-$$

**Protected Attribute**
- Favoured
- Unfavoured

**Class**
- Positive
- Negative

$U+$  $F+$

$U-$  $F-$

## Individual Fairness

**Similar individuals** should be treated in a **similar way**

## Demographic Parity

Same **positive rate** across *PA* groups

## Equalised Odds

$\hat{Y}$ and *PA* are **independent**, conditional on *Y*
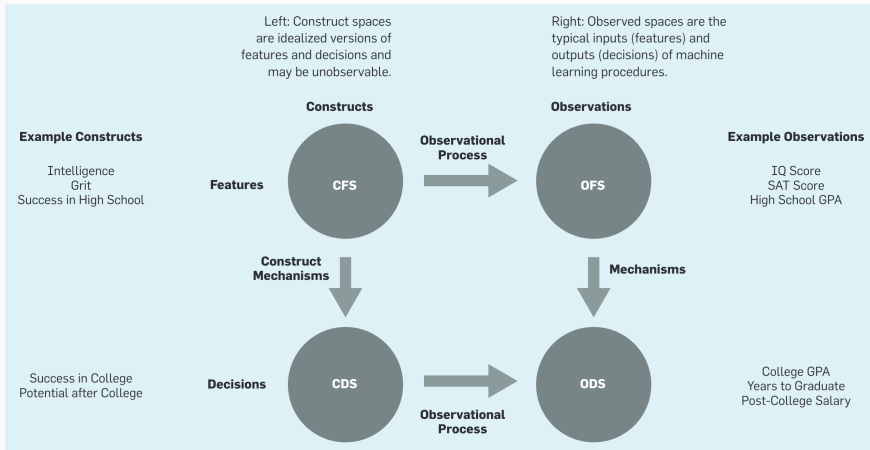
## Individual Fairness

$$d(x_1, x_2) \leqslant \delta \Rightarrow \hat{Y}(x_1) \approx \hat{Y}(x_2)$$

## Demographic Parity

$$P(\hat{Y} = 1 \mid PA = 0) = P(\hat{Y} = 1 \mid PA = 1)$$

## Equalised Odds

$$P(\hat{Y} = 1 \mid PA = 0, Y = y) = P(\hat{Y} = 1 \mid PA = 1, Y = y), \quad y \in \{0, 1\}$$
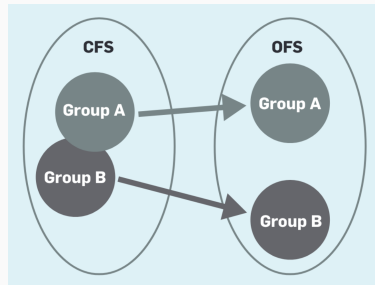
# The (Im)possibility of Fairness[1]



Left: Construct spaces are idealized versions of features and decisions and may be unobservable.

Right: Observed spaces are the typical inputs (features) and outputs (decisions) of machine learning procedures.

**Constructs**

**Observations**

**Example Constructs**

Intelligence
Grit
Success in High School

**Features**

CFS

**Observational Process**

OFS

**Example Observations**

IQ Score
SAT Score
High School GPA

**Construct Mechanisms**

**Mechanisms**

Success in College
Potential after College

**Decisions**

CDS

ODS

**Observational Process**

College GPA
Years to Graduate
Post-College Salary

[1]Sorelle A Friedler, Carlos Scheidegger, and Suresh Venkatasubramanian (2021). "The (Im)possibility of Fairness". In: *Communications of the ACM*.

## WYSIWYG

Construct space and observed space maintain the relative position of individuals w.r.t. the task. Aligns with individual fairness.

## We're All Equal

Within a given construct space all groups are essentially the same. Aligns with group fairness.



---

[1]Sorelle A Friedler, Carlos Scheidegger, and Suresh Venkatasubramanian (2021). "The (Im)possibility of Fairness". In: *Communications of the ACM*.

**Pre-Processing**  Modify the training set to "sample from a better world"

**In-Processing**  Add *constraints* or *regularisation terms* to improve fairness

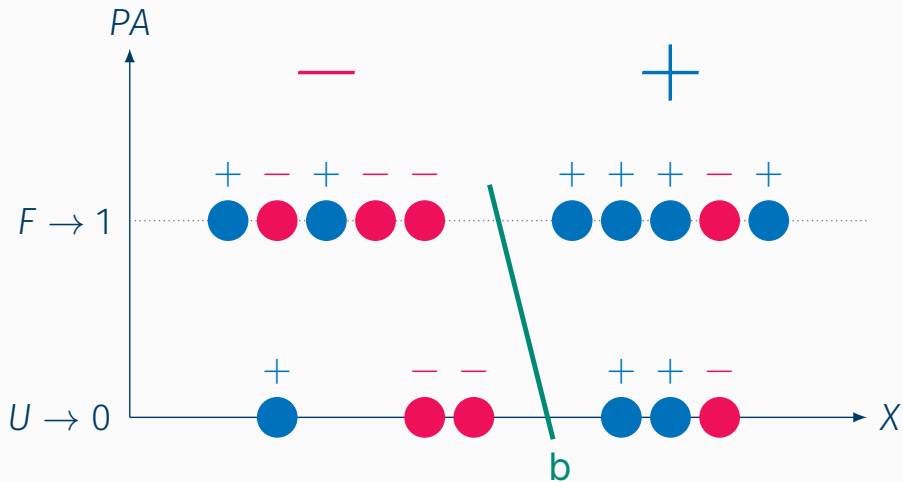**Post-Processing**  Adjust the predictions after fitting the model

**Pre-Processing**  Modify the training set to "sample from a better world"

**In-Processing**  Add *constraints* or *regularisation terms* to improve fairness

**Post-Processing**  Adjust the predictions after fitting the model

# Parametrised Data Sampling

Original Data     Undersampling     Oversampling     Both

[2]Vladimiro González-Zelaya, Julián Salas, Dennis Prangle, and Paolo Missier (2021). "Optimising Fairness through Parametrised Data Sampling". In: *International Conference on Extending Database Technology*.
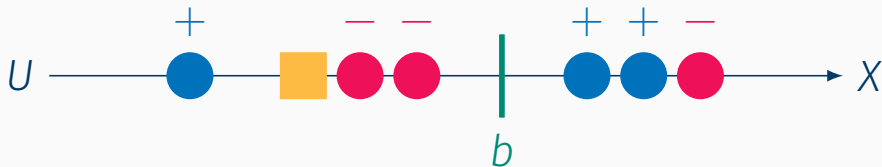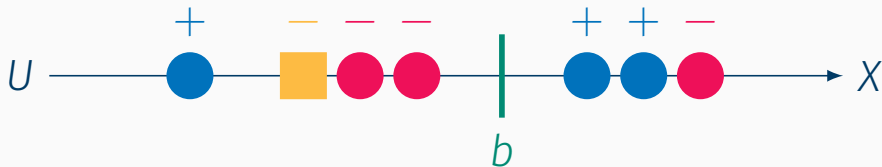
The positive predictions for *U* increase by:
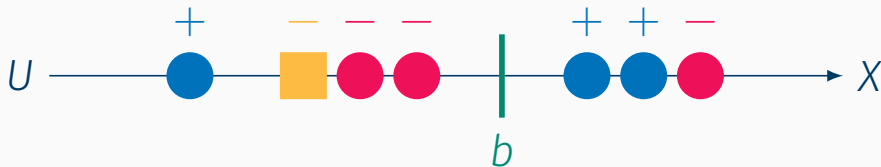
► Undersampling negative instances
► Oversampling positive instances

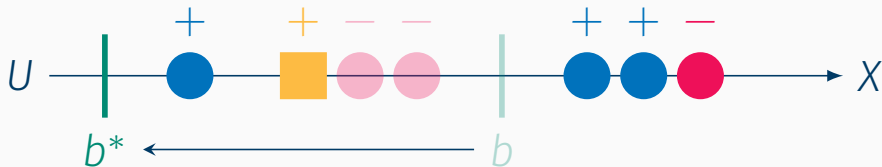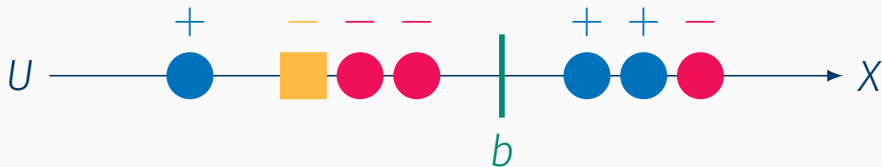|  | Age | Country | Gender | Ethnicity | Combined PA |
|---|---|---|---|---|---|
|  | (20-30] | Mexico | Male | Latin | Unfavoured |
|  | (30-40] | Canada | Female | White | Favoured |

|  | Age | Country | Gender | Ethnicity |
|---|---|---|---|---|
| Subgroup PR |  |  |  |  |
| Dataset PR | 0.3 | 0.3 | 0.3 | 0.3 |
| Difference |  |  |  |  |

| | Age | Country | Gender | Ethnicity | Combined PA |
|---|---|---|---|---|---|
| | (20-30] | Mexico | Male | Latin | Unfavoured |
| | (30-40] | Canada | Female | White | Favoured |

| | | | | |
|---|---|---|---|---|
| Subgroup PR | | | | |
| Dataset PR | 0.3 | 0.3 | 0.3 | 0.3 |
| Difference | | | | |

| | Age | Country | Gender | Ethnicity | Combined PA |
|---|---|---|---|---|---|
| | (20-30] | Mexico | Male | Latin | Unfavoured |
| | (30-40] | Canada | Female | White | Favoured |
| Subgroup PR | 0.2 | 0.3 | 0.4 | 0.1 | |
| Dataset PR | 0.3 | 0.3 | 0.3 | 0.3 | |
| Difference | | | | | |

# Dealing with Multiple PAs

| | Age | Country | Gender | Ethnicity | Combined PA |
|---|---|---|---|---|---|
| | (20-30] | Mexico | Male | Latin | Unfavoured |
| | (30-40] | Canada | Female | White | Favoured |
| Subgroup PR | 0.2 | 0.3 | 0.4 | 0.1 | |
| Dataset PR | 0.3 | 0.3 | 0.3 | 0.3 | |
| Difference | | | | | |

| | Age | Country | Gender | Ethnicity | Combined PA |
|---|---|---|---|---|---|
| | (20-30]<br>(30-40] | Mexico<br>Canada | Male<br>Female | Latin<br>White | Unfavoured<br>Favoured |
| Subgroup PR | 0.2 | 0.3 | 0.4 | 0.1 | |
| Dataset PR | 0.3 | 0.3 | 0.3 | 0.3 | |
| Difference | −0.1 | +0.0 | +0.1 | −0.2 | |

| | Age | Country | Gender | Ethnicity | Combined PA |
|---|---|---|---|---|---|
| | (20-30] | Mexico | Male | Latin | Unfavoured |
| | (30-40] | Canada | Female | White | Favoured |
| Subgroup PR | 0.2 | 0.3 | 0.4 | 0.1 | |
| Dataset PR | 0.3 | 0.3 | 0.3 | 0.3 | |
| Difference | $-0.1$ | $+0.0$ | $+0.1$ | $-0.2$ | Sum $= -0.2$ |

# Dealing with Multiple PAs

| | Age | Country | Gender | Ethnicity | Combined PA |
|---|---|---|---|---|---|
| | (20-30] | Mexico | Male | Latin | Unfavoured |
| | (30-40] | Canada | Female | White | Favoured |
| Subgroup PR | 0.2 | 0.3 | 0.4 | 0.1 | |
| Dataset PR | 0.3 | 0.3 | 0.3 | 0.3 | |
| Difference | $-0.1$ | $+0.0$ | $+0.1$ | $-0.2$ | Sum $= -0.2$ |

| | Age | Country | Gender | Ethnicity | Combined PA |
|---|---|---|---|---|---|
| | (20-30] | Mexico | Male | Latin | Unfavoured |
| | (30-40] | Canada | Female | White | Favoured |

| | Age | Country | Gender | Ethnicity | |
|---|---|---|---|---|---|
| Subgroup PR | | | | | |
| Dataset PR | 0.3 | 0.3 | 0.3 | 0.3 | |
| Difference | | | | | |

|  | Age | Country | Gender | Ethnicity | Combined PA |
|---|---|---|---|---|---|
| | (20-30] | Mexico | Male | Latin | Unfavoured |
| | (30-40] | Canada | Female | White | Favoured |
| Subgroup PR | 0.4 | 0.4 | 0.1 | 0.4 | |
| Dataset PR | 0.3 | 0.3 | 0.3 | 0.3 | |
| Difference | $+0.1$ | $+0.1$ | $-0.2$ | $+0.1$ | Sum $= +0.1$ |

# Dealing with Multiple PAs

| | Age | Country | Gender | Ethnicity | Combined PA |
|---|---|---|---|---|---|
| | (20-30]<br>(30-40] | Mexico<br>Canada | Male<br>Female | Latin<br>White | Unfavoured<br>Favoured |
| Subgroup PR | 0.4 | 0.4 | 0.1 | 0.4 | |
| Dataset PR | 0.3 | 0.3 | 0.3 | 0.3 | |
| Difference | $+0.1$ | $+0.1$ | $-0.2$ | $+0.1$ | $\text{Sum} = +0.1$ |

# Fair and Private Data Correction

Both aim at concealing sensitive information while preserving data utility:

**Fairness**  To prevent *classifier* behaviours related to sensitive data

**Privacy**  To protect sensitive data from disclosure to *adversaries*

**Quasi-Identifiers** Collection of features such that their values may be used to **re-identify** an individual

$k$-**Anonymity** Dataset records' are **indistinguishable** from at least $k-1$ other records w.r.t. QIs

$k$-**Group** Set of **indistinguishable** records in a $k$-anonymous dset

$t$-**Closeness** The **PA distributions** of the $k$-groups are **similar** to the whole dataset's

**Quasi-Identifiers** Collection of features such that their values may be used to **re-identify** an individual

$k$-**Anonymity** Dataset records' are **indistinguishable** from at least $k - 1$ other records w.r.t. QIs

$k$-**Group** Set of **indistinguishable** records in a $k$-anonymous dset

$t$-**Closeness** The **PA distributions** of the $k$-groups are **similar** to the whole dataset's

**Quasi-Identifiers** Collection of features such that their values may be used to **re-identify** an individual

$k$**-Anonymity** Dataset records' are **indistinguishable** from at least $k-1$ other records w.r.t. QIs

$k$**-Group** Set of **indistinguishable** records in a $k$-anonymous dset

$t$**-Closeness** The **PA distributions** of the $k$-groups are **similar** to the whole dataset's
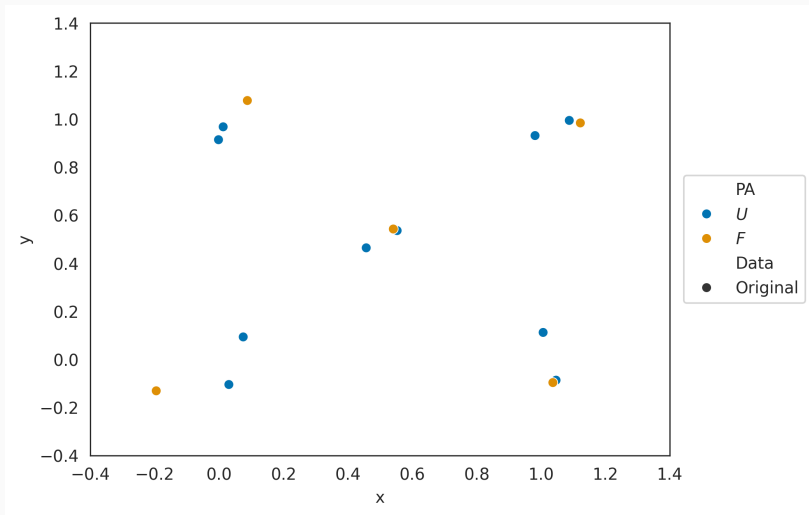
**Quasi-Identifiers**  Collection of features such that their values may be used to **re-identify** an individual

**$k$-Anonymity**  Dataset records' are **indistinguishable** from at least $k - 1$ other records w.r.t. QIs
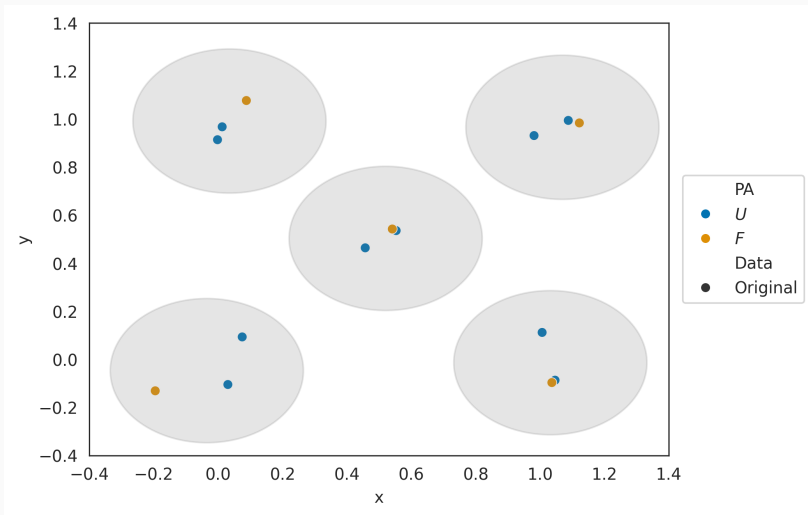
**$k$-Group**  Set of **indistinguishable** records in a $k$-anonymous dset

**$t$-Closeness**  The **PA distributions** of the $k$-groups are **similar** to the whole dataset's
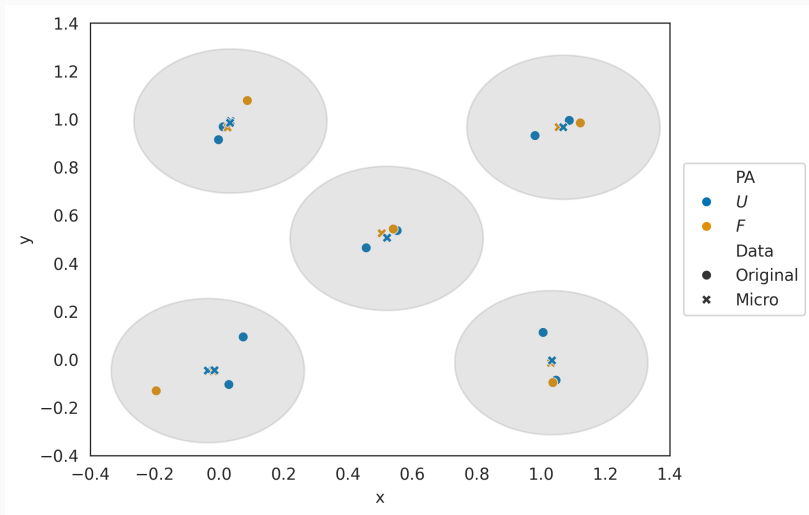
# MDAV — Maximum Distance to AVerage[3]

[3] Josep Domingo-Ferrer and Vicenç Torra (2005). "Ordinal, Continuous and Heterogeneous k-Anonymity through Microaggregation". In: *Data Mining and Knowledge Discovery* 11, pp. 195–212.

# MDAV — Maximum Distance to AVerage[3]
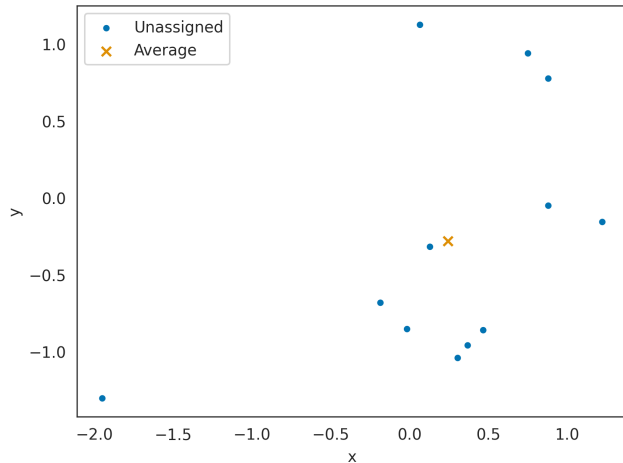


[3] Josep Domingo-Ferrer and Vicenç Torra (2005). "Ordinal, Continuous and Heterogeneous k-Anonymity through Microaggregation". In: *Data Mining and Knowledge Discovery* 11, pp. 195–212.
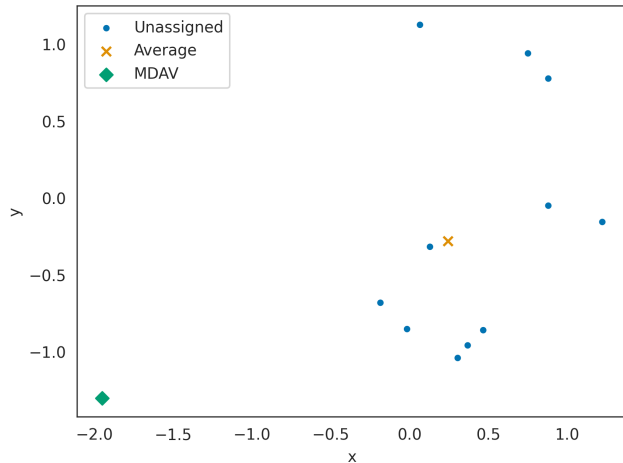
# MDAV — Maximum Distance to AVerage[3]



[3] Josep Domingo-Ferrer and Vicenç Torra (2005). "Ordinal, Continuous and Heterogeneous k-Anonymity through Microaggregation". In: *Data Mining and Knowledge Discovery* 11, pp. 195–212.
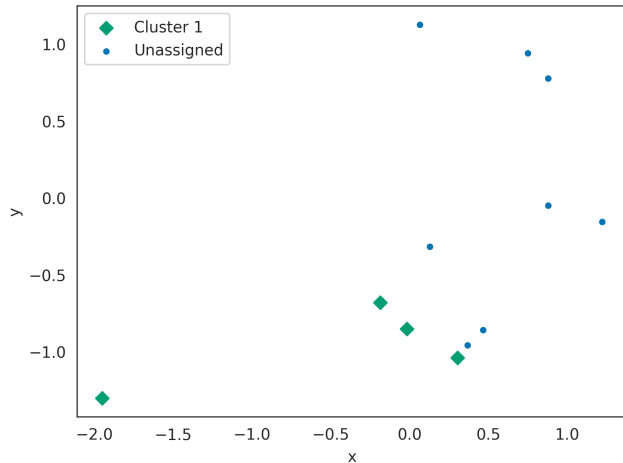
# MDAV — Maximum Distance to AVerage[3]

[3] Josep Domingo-Ferrer and Vicenç Torra (2005). "Ordinal, Continuous and Heterogeneous k-Anonymity through Microaggregation". In: *Data Mining and Knowledge Discovery* 11, pp. 195–212.

# MDAV — Maximum Distance to AVerage[3]

[3] Josep Domingo-Ferrer and Vicenç Torra (2005). "Ordinal, Continuous and Heterogeneous k-Anonymity through Microaggregation". In: *Data Mining and Knowledge Discovery* 11, pp. 195–212.
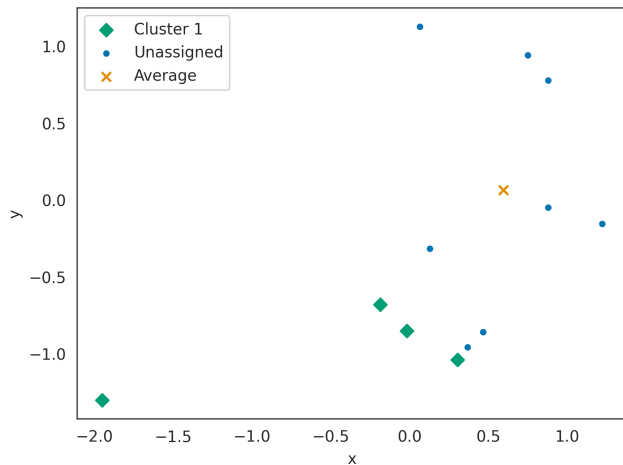
[3] Josep Domingo-Ferrer and Vicenç Torra (2005). "Ordinal, Continuous and Heterogeneous k-Anonymity through Microaggregation". In: *Data Mining and Knowledge Discovery* 11, pp. 195–212.
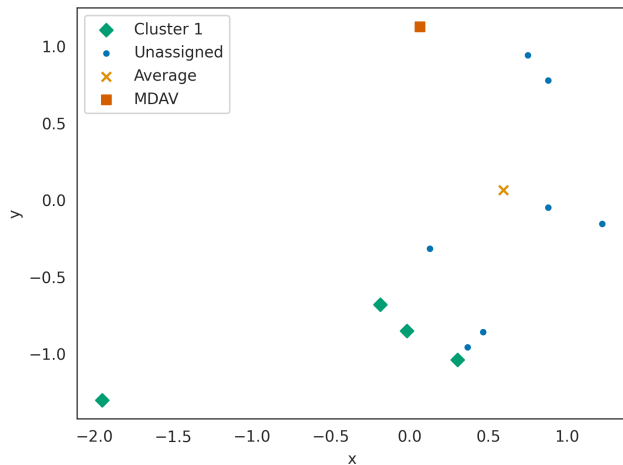
# MDAV — Maximum Distance to AVerage[3]



[3]Josep Domingo-Ferrer and Vicenç Torra (2005). "Ordinal, Continuous and Heterogeneous k-Anonymity through Microaggregation". In: *Data Mining and Knowledge Discovery* 11, pp. 195–212.

# MDAV — Maximum Distance to AVerage[3]



[3]Josep Domingo-Ferrer and Vicenç Torra (2005). "Ordinal, Continuous and Heterogeneous k-Anonymity through Microaggregation". In: *Data Mining and Knowledge Discovery* 11, pp. 195–212.
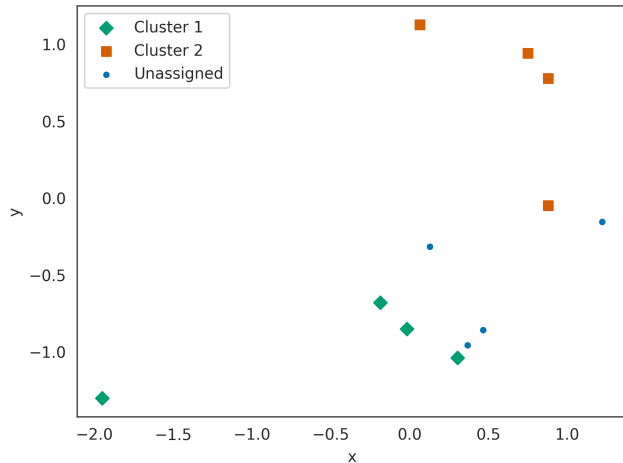
[3] Josep Domingo-Ferrer and Vicenç Torra (2005). "Ordinal, Continuous and Heterogeneous k-Anonymity through Microaggregation". In: *Data Mining and Knowledge Discovery* 11, pp. 195–212.
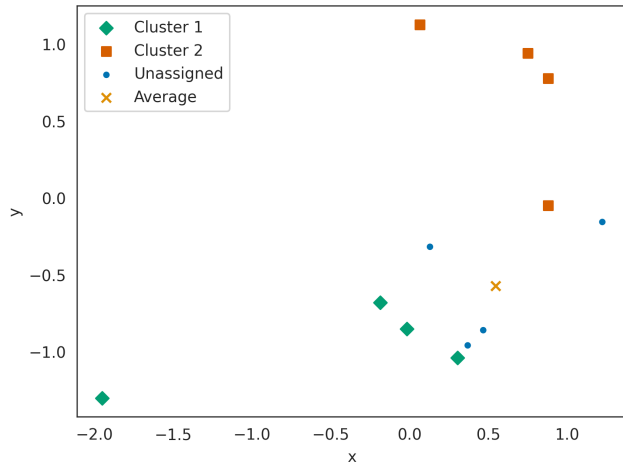
# MDAV — Maximum Distance to AVerage[3]



[3] Josep Domingo-Ferrer and Vicenç Torra (2005). "Ordinal, Continuous and Heterogeneous k-Anonymity through Microaggregation". In: *Data Mining and Knowledge Discovery* 11, pp. 195–212.
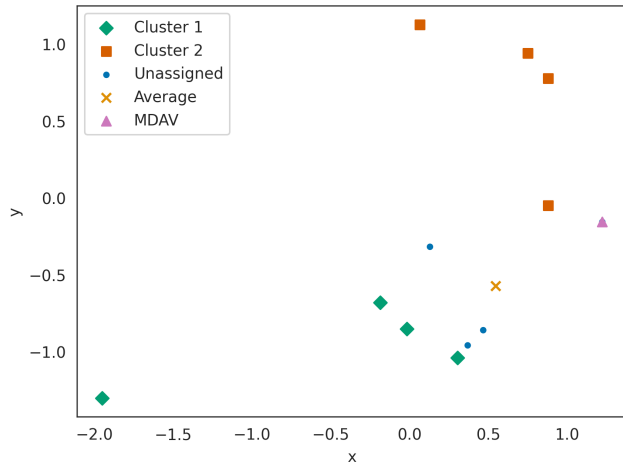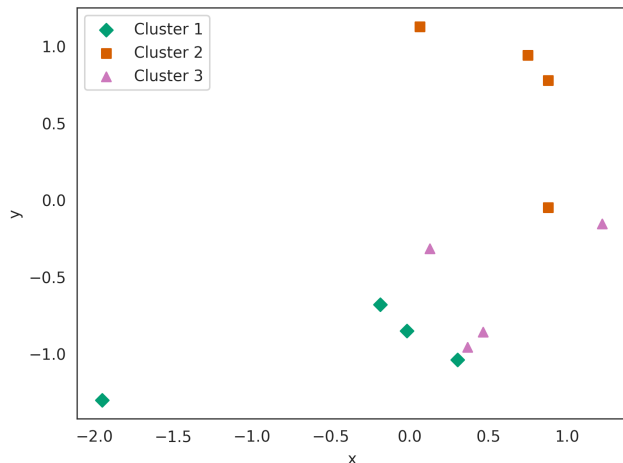
1. **Cluster** $D$ into $(m, n)$-*fairlets* (sets of $m$ unfavoured and $n$ favoured records) using the MDAV algorithm, where

$$\frac{m}{n} \approx \frac{|U|}{|F|}, \text{ subject to } m + n = k$$

2. **Microaggregate** the feature values with their corresponding fairlet's mean/mode, except for PA and Class, whose original values are kept

3. **Locally correct the fairness** of each fairlet by relabelling its records depending on their PA values, so that

$$\frac{|U^+|}{|U|} \geqslant \tau \cdot \frac{|F^+|}{|F|}, \text{ where } \tau \text{ modulates the correction}$$

---

[4]Vladimiro González-Zelaya, Julián Salas, David Megías, and Paolo Missier (2023). "Fair and Private Data Preprocessing through Microaggregation". In: *ACM Transactions on Knowledge Discovery from Data*.

1. **Cluster** $D$ into $(m, n)$-*fairlets* (sets of $m$ unfavoured and $n$ favoured records) using the MDAV algorithm, where

$$\frac{m}{n} \approx \frac{|U|}{|F|}, \text{ subject to } m + n = k$$

2. **Microaggregate** the feature values with their corresponding fairlet's mean/mode, except for PA and Class, whose original values are kept

3. **Locally correct the fairness** of each fairlet by relabelling its records depending on their PA values, so that

$$\frac{|U^+|}{|U|} \geqslant \tau \cdot \frac{|F^+|}{|F|}, \text{ where } \tau \text{ modulates the correction}$$

---

[4]Vladimiro González-Zelaya, Julián Salas, David Megías, and Paolo Missier (2023). "Fair and Private Data Preprocessing through Microaggregation". In: *ACM Transactions on Knowledge Discovery from Data*.

1. **Cluster** $D$ into $(m, n)$-*fairlets* (sets of $m$ unfavoured and $n$ favoured records) using the MDAV algorithm, where

$$\frac{m}{n} \approx \frac{|U|}{|F|}, \text{ subject to } m + n = k$$

2. **Microaggregate** the feature values with their corresponding fairlet's mean/mode, except for PA and Class, whose original values are kept

3. **Locally correct the fairness** of each fairlet by relabelling its records depending on their PA values, so that

$$\frac{|U^+|}{|U|} \geqslant \tau \cdot \frac{|F^+|}{|F|}, \text{ where } \tau \text{ modulates the correction}$$

---

[4]Vladimiro González-Zelaya, Julián Salas, David Megías, and Paolo Missier (2023). "Fair and Private Data Preprocessing through Microaggregation". In: *ACM Transactions on Knowledge Discovery from Data*.

# Fair-MDAV Example ($k = 3$, $\tau = 1$)

## Example Data

| id | $X$ | PA | Class |
|----|-----|-----|-------|
| a | 1 | $F$ | 1 |
| b | 2 | $U$ | 0 |
| c | 3 | $U$ | 1 |
| d | 11 | $F$ | 0 |
| e | 12 | $F$ | 0 |
| f | 13 | $F$ | 1 |
| g | 14 | $F$ | 1 |

## (1, 2)-Fairlets

| id | $X$ | $X_{ma}$ | PA | Class | Fair Class |
|----|-----|----------|-----|-------|------------|
| a | 1 | 4.67 | $F$ | 1 | 1 |
| b | 2 | 4.67 | $U$ | 0 | 1 |
| c | 3 | 9.33 | $U$ | 1 | 1 |
| d | 11 | 4.67 | $F$ | 0 | 0 |
| e | 12 | 9.33 | $F$ | 0 | 0 |
| f | 13 | 9.33 | $F$ | 1 | 1 |
| g | | *Dropped* | | | |

The following parameter values were tested over three benchmark datasets [5]:

| Description | Parameter | Values |
|---|---|---|
| Fairlet Size | $m$ | $\left.\begin{array}{l}\left\lceil k \cdot \frac{\|U\|}{\|D\|}\right\rceil \\ k - m\end{array}\right\}$ for $k \in \{10, 20, \cdots, 100\}$ |
| | $n$ | |
| Fairness Correction | $\tau$ | $0, 0.1, \cdots, 1$ |
| Microaggregation | $ma$ | `True`, `False` |

[5] *Census Income*, *COMPAS*, and *German Credit*, available online

# Fairness/Accuracy Trade-Off (Census Income Dataset)

# Genetic Pipeline Optimisation

▶ Steps that transform the raw input data into its final form as a training set

▶ Some are *required* by the classification framework:

  ▶ Encoding categorical variables

  ▶ Imputing missing data

▶ Others may *optionally* be deployed:

  ▶ Class balancing

  ▶ Feature selection

  ▶ Feature scaling

▶ Steps usually combined into *pipelines* based on best-practice considerations, with model *performance* as the main objective
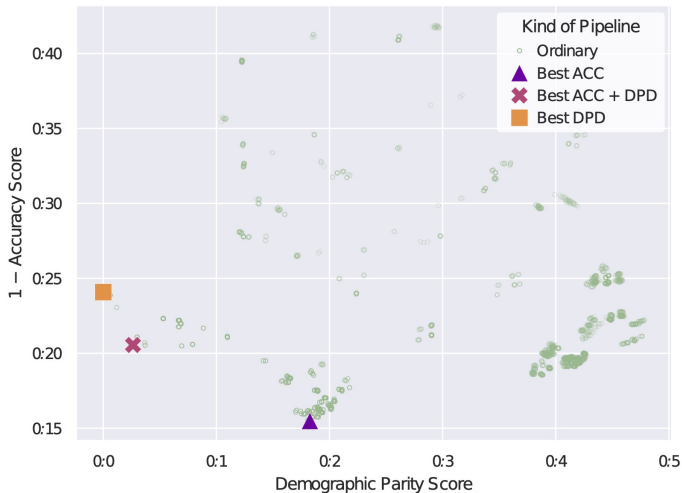
► Steps that transform the raw input data into its final form as a training set
► Some are *required* by the classification framework:
  ► Encoding categorical variables
  ► Imputing missing data
► Others may *optionally* be deployed:
  ► Class balancing
  ► Feature selection
  ► Feature scaling
► Steps usually combined into *pipelines* based on best-practice considerations, with model *performance* as the main objective
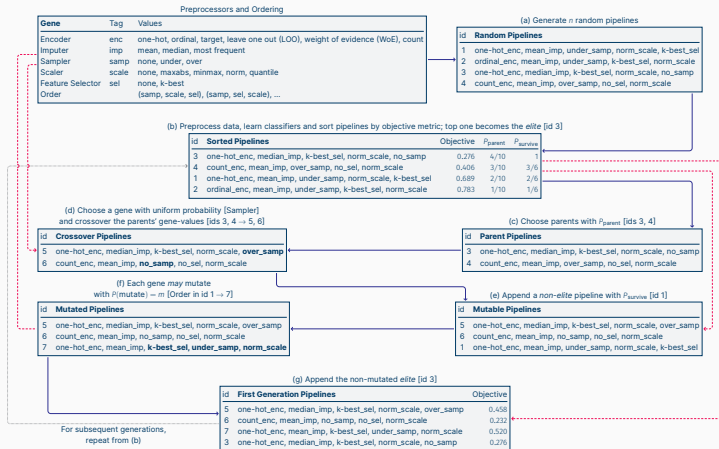
► Steps that transform the raw input data into its final form as a training set
► Some are *required* by the classification framework:
  ► Encoding categorical variables
  ► Imputing missing data
► Others may *optionally* be deployed:
  ► Class balancing
  ► Feature selection
  ► Feature scaling
► Steps usually combined into *pipelines* based on best-practice considerations, with model *performance* as the main objective

- Steps that transform the raw input data into its final form as a training set
- Some are *required* by the classification framework:
  - Encoding categorical variables
  - Imputing missing data
- Others may *optionally* be deployed:
  - Class balancing
  - Feature selection
  - Feature scaling
- Steps usually combined into *pipelines* based on best-practice considerations, with model *performance* as the main objective
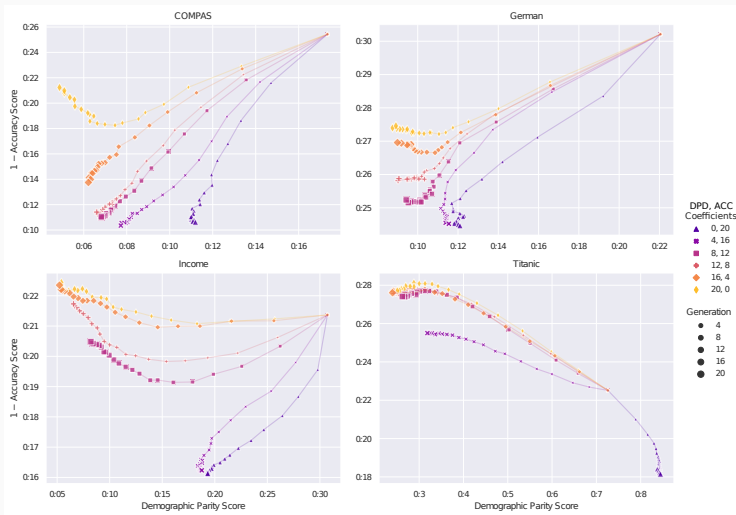
# Pipeline-Space Fairness/Accuracy and Pareto Front

[6]Vladimiro González-Zelaya, Julián Salas, Dennis Prangle, and Paolo Missier (2023). "Preprocessing Matters: Automated Pipeline Selection for Fair Classification". In: *MDAI 2023*.

# Evolution Toward Optimal Solutions

▶ There are multiple incompatible fairness definitions

▶ Definitions should align with world-views

▶ Possible to correct fairness through data pre-processing

▶ Can use fairness-specific methods or optimise pipeline choices

▶ Fairness and privacy can be achieved simultaneously

► There are multiple incompatible fairness definitions

► Definitions should align with world-views

► Possible to correct fairness through data pre-processing

► Can use fairness-specific methods or optimise pipeline choices

► Fairness and privacy can be achieved simultaneously

► There are multiple incompatible fairness definitions

► Definitions should align with world-views

► Possible to correct fairness through data pre-processing

► Can use fairness-specific methods or optimise pipeline choices

► Fairness and privacy can be achieved simultaneously

▶ There are multiple incompatible fairness definitions

▶ Definitions should align with world-views

▶ Possible to correct fairness through data pre-processing

▶ Can use fairness-specific methods or optimise pipeline choices

▶ Fairness and privacy can be achieved simultaneously

► There are multiple incompatible fairness definitions

► Definitions should align with world-views

► Possible to correct fairness through data pre-processing

► Can use fairness-specific methods or optimise pipeline choices

► Fairness and privacy can be achieved simultaneously