ECDA 2013 - Luxembourg

# Mining preferential datasets in MCDA

Alexandru-Liviu Olteanu [1,2]     Raymond Bisdorff [1]

1. Université du Luxembourg     2. Institut Télécom, Télécom Bretagne
Université Européenne de Bretagne

10th of July 2013

# Contents

# Data mining and clustering

## Data mining and clustering

**Data**
- many forms;
  (measurements, observations, dynamics of processes, text, images, etc.)
- large quantities [GANTZ AND REINSEL 2011];
  $\approx 10^{21}$ bytes (100 TB for each person on the planet)

### Data mining

- process that **extracts information** from a data set and **transforms** it into an **understandable structure** for further use;
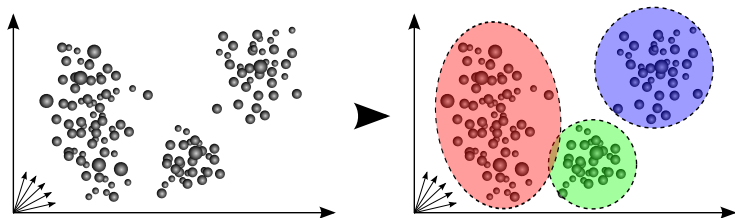
**Data** • many forms;

(measurements, observations, dynamics of processes, text, images, etc.)

• large quantities [Gantz and Reinsel 2011];

$\approx 10^{21}$ bytes (100 TB for each person on the planet)

### Data mining

• process that **extracts information** from a data set and **transforms** it into an **understandable structure** for further use;

# Multiple criteria decision aid

## Multiple Criteria Decision Aid

- aims at modelling the **preferences** of decision-makers;
- **aids** them in reaching certain **decisions**;

| Objects | Attributes | | | |
|---------|-------|--------------|--------|------|
|         | Price | Acceleration | Safety | $\cdots$ |
| Car 1   | 18,342 | 30.7s | good | $\cdots$ |
| Car 2   | 15,335 | 30.2s | medium | $\cdots$ |
| Car 3   | 16,973 | 29s | v.good | $\cdots$ |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | |

## Multiple Criteria Decision Aid

- aims at modelling the **preferences** of decision-makers;
- **aids** them in reaching certain **decisions**;

| Alternatives | Criteria | | | |
|:---:|:---:|:---:|:---:|:---:|
| | Price ↓ | Acceleration ↓ | Safety ↑ | $\cdots$ |
| Car 1 | **18,342** | **30.7s** | good | $\cdots$ |
| Car 2 | **15,335** | 30.2s | **medium** | $\cdots$ |
| Car 3 | 16,973 | **29s** | **v.good** | $\cdots$ |
| ⋮ | ⋮ | ⋮ | ⋮ | |

**Modelling preferences**

### Value functions

• aggregate all the criteria into a **score**;

• $(x_i, x_j, x_k, \ldots) \rightarrow U(x)$;

• **trade-offs** between criteria;

### Outranking relations

• x **outranks** y iff:

1) $x$ is **at least as good as** $y$ on a **weighted majority** of criteria;
2) $x$ is **not much worse** than $y$ on any criterion;
$\rightarrow$ $x \, S \, y$

• similar to **voting**;

## Modelling preferences

### Value functions

• aggregate all the criteria into a **score**;

• $(x_i, x_j, x_k, \ldots) \rightarrow U(x)$;

• **trade-offs** between criteria;

### Outranking relations

• x **outranks** y iff:

1) $x$ is **at least as good as** $y$ on a **weighted majority** of criteria;
2) $x$ is **not much worse** than $y$ on any criterion;
$\rightarrow \quad x \, S \, y$

• similar to **voting**;

### Preferential situations

| | | |
|---|---|---|
| $U(x) = U(y)$ | Indifference (I) | $x \, S \, y \wedge y \, S \, x$ |
| $U(x) > U(y)$ | Strict preference (P) | $x \, S \, y \wedge y \, \cancel{S} \, x$ |
| $U(x) \geqslant U(y)$ | Weak preference (Q) | $x \, S \, y$ |
| | Incomparability (R) | $x \, \cancel{S} \, y \wedge y \, \cancel{S} \, x$ |

**Decision problems**

Choice              Sorting              Ranking



C1 $\succ$ C2 $\succ$ C3

s

**Decision problems**

| Choice | Sorting | Ranking |
|---|---|---|



$C1 \succ C2 \succ C3$

**Clustering in MCDA**

Existing approaches:

- that use similarity measures:
- that use **preferential information**:

Formally defined using preferential relations in [MEYER, OLTEANU 2013]s

| Data mining | MCDA |
|---|---|
| • objects + attributes<br>• **similarity** | ▶ | • alternatives + criteria<br>• **indifference**, strict preference<br>incomparability |
| • clustering - formally defined<br>using similarity<br>measures | ▶ | • clustering - formally defined<br>using preferential<br>measures<br>[MEYER, OLTEANU, 2013] |
| • problem size - easily > $10^6$ | ▶ | • problem size - rarely > 100 |

# Clustering in MCDA

### Clustering in Data mining

• process that groups objects that are **similar** and separates those that are **dissimilar**;

- objects that **cannot be distinguished** are **similar**;

## Clustering in Data mining

• process that groups objects that are **similar** and separates those that are **dissimilar**;

- objects that **cannot be distinguished** are **similar**;
  **BUT**
- alternatives that **cannot be distinguished** are **indifferent**.

## Clustering in Data mining

• process that groups objects that are **similar** and separates those that are **dissimilar**;

- objects that **cannot be distinguished** are **similar**;
  **BUT**
- alternatives that **cannot be distinguished** are **indifferent**.

## Clustering in MCDA

• process that groups alternatives that are **indifferent** and separates those that are **not indifferent**;

Classical clustering

Classical clustering          Non-relational clustering

Classical clustering

Non-relational clustering

Relational clustering

Classical clustering

Non-relational clustering

Relational clustering

Ordered clustering

# Case Study: U.S. Toxic Chemicals Release Practices

**The data**

- Toxic Chemical Release Practices of facilities in the U.S.;

- $> 53,000$ facilities reporting over 25 years;

- **selected** data from 2010 ($\sim 22,000$ reports);

- chemical **toxicity** information;
- reports containing the **release amounts** of a chemical;
- reports containing the **mitigated amounts** of a chemical;

## Case Study: U.S. Toxic Chemicals Release Practices

**The data**

- Toxic Chemical Release Practices of facilities in the U.S.;
- $> 53,000$ facilities reporting over 25 years;
- **selected** data from 2010 ($\sim 22,000$ reports);

- chemical **toxicity** information;
- reports containing the **release amounts** of a chemical;
- reports containing the **mitigated amounts** of a chemical;

**The problem**

- **classifying** these practices w.r.t. their quality **without knowing** the classes a priori;
- **PREFERENCES**: handling of **less toxic** chemicals, **fewer releases** and **better mitigation** procedures;

**Structuring the problem**



■ Fictive decision-maker: bipolar-valued outranking relation [BISDORFF 2012];

**Non-relational clustering**

- used algorithms from [MEYER, OLTEANU 2013]

- selected one result to illustrate (12 clusters);



| Fitness (%) | |
|---|---|
| $f_{NR}^*$ | **80.0** |
| $f_{NR}$ | 62.7 |
| $f_{NR}^{min}$ | 0.0 |
| $f_R^*$ | 64.0 |
| $f_R$ | 55.2 |
| $f_R^{min}$ | 0.0 |

| Cluster sizes | |
|---|---|
| $K_1$ | 296 |
| $K_2$ | **1,429** |
| $K_3$ | **1,632** |
| $K_4$ | **6,237** |
| $K_5$ | **2,973** |
| $K_6$ | 167 |
| $K_7$ | **1,316** |
| $K_8$ | **2,615** |
| $K_9$ | **1,688** |
| $K_{10}$ | 540 |
| $K_{11}$ | **3,518** |
| $K_{12}$ | 356 |

**Conclusions and Perspectives**

**Conclusions:**

- highlighted clustering using preferential information;

- illustrated an application of clustering in MCDA;

**Perspectives:**

- further explore clustering in MCDA (different structures);

- methodology for using clustering when eliciting the parameters of a preference model;

- combining similarity-based and indifference-based clustering (2 layers).

# **Mining preferential datasets in MCDA**

1. **Data mining and clustering**

2. **Multiple criteria decision aid**

3. **Clustering in MCDA**

4. **Case Study: U.S. Toxic Chemicals Release Practices**

5. **Conclusions and Perspectives**

Illustrative example:

| $F$ | $i$ | $j$ | $k$ |
|---|---|---|---|
| $w$ | 1 | 1 | 1 |
| $x$ | GOOD | MEDIUM | BAD |
| $y$ | BAD | MEDIUM | GOOD |

Illustrative example:

| $F$ | $i$ | $j$ | $k$ |
|---|---|---|---|
| $w$ | 1 | 1 | 1 |
| $x$ | GOOD | MEDIUM | BAD |
| $y$ | BAD | MEDIUM | GOOD |

Similarity:

$x_i \neq y_i \; x_j = y_j \; x_k \neq y_k \rightarrow$ x,y - **dissimilar**;

Indifference:

$$x_i \succcurlyeq y_i \; x_j \succcurlyeq y_j \; x_k \not\succcurlyeq y_k \rightarrow \text{x outranks y}$$
$$y_i \not\succcurlyeq x_i \; y_j \succcurlyeq x_j \; y_k \succcurlyeq x_k \rightarrow \text{y outranks x}$$

$\left. \right\} \rightarrow$ x,y - **indifferent**.

**Comparative analysis**

Measures:

- **similarity** measures from the Manhattan distance ($S_{L_1}$), the Euclidian distance ($S_{L_2}$), from [BISDORFF,MEYER,OLTEANU 2011] ($S_{THR}$) and **indifference** measure from the outranking relation in [BISDORFF,MEYER,ROUBENS 2007] ($I_{\tilde{S}}$);

**Comparative analysis**

Measures:

- **similarity** measures from the Manhattan distance ($S_{L_1}$), the Euclidian distance ($S_{L_2}$), from [BISDORFF,MEYER,OLTEANU 2011] ($S_{THR}$) and **indifference** measure from the outranking relation in [BISDORFF,MEYER,ROUBENS 2007] ($I_{\check{S}}$);

Experiment:

- all feasible alternatives on a fixed number of criteria with fixed number of values;

- compared similarity and indifferent measures for all pairs of alternatives;

**Comparative analysis**

Measures:

- **similarity** measures from the Manhattan distance ($S_{L_1}$), the Euclidian distance ($S_{L_2}$), from [BISDORFF,MEYER,OLTEANU 2011] ($S_{THR}$) and **indifference** measure from the outranking relation in [BISDORFF,MEYER,ROUBENS 2007] ($I_{\tilde{S}}$);

Experiment:

- all feasible alternatives on a fixed number of criteria with fixed number of values;

- compared similarity and indifferent measures for all pairs of alternatives;

Results:

- in at least 25% cases **dissimilar** alternatives were in fact **indifferent**;

- **significant differences** between similarity and indifference.

$B_1^+$

| TH | TE | RA | MA |
|----|----|----|----|
| 0 | 1 | 98 | 526.4·10⁶ |
| 0 |  | 51 | 99786 |
| 0 |  | 10.0·10⁶ | 45.4·10⁶ |
|  |  | 21 | 80.0·10⁶ |
|  |  | 44 |  |

$B_2^+$

| TH | TE | RA | MA |
|----|----|----|----|
| 0 | 1 | 6 | 526.4·10⁶ |
| 0 |  | 10 | 99786 |
| 1 |  | 22 | 45.4·10⁶ |
|  |  | 9 | 80.0·10⁶ |
|  |  | 16 |  |

$B_3^+$

| TH | TE | RA | MA |
|----|----|----|----|
| 0 | 0 | 10 | 526.4·10⁶ |
| 0 |  | 3 | 30708 |
| 0 |  | 3 | 45.4·10⁶ |
|  |  | 21 | 14.0·10⁶ |
|  |  | 9 |  |

$B_4^+$

| TH | TE | RA | MA |
|----|----|----|----|
| 1 | 0 | 1 | 526.0·10⁶ |
| 1 |  | 1 | 43114 |
| 0 |  | 5 | 6.0·10⁶ |
|  |  | 21 | 80.0·10⁶ |
|  |  | 3 |  |

$B_5^+$

| TH | TE | RA | MA |
|----|----|----|----|
| 1 | 0 | 5 | 526.4·10⁶ |
| 1 |  | 6 | 99786 |
| 1 |  | 7 | 45.4·10⁶ |
|  |  | 17 | 80.0·10⁶ |
|  |  | 7 |  |

$B_6^+$

| TH | TE | RA | MA |
|----|----|----|----|
| 1 | 0 | 7 | 526.4·10⁶ |
| 1 |  | 21 | 99786 |
| 1 |  | 15 | 45.4·10⁶ |
|  |  | 0 | 80.0·10⁶ |
|  |  | 15 |  |

61% 39%   90% 10%   38% 62%   74% 26%   45% 55%   81% 19%

$K_1$   $K_2$   $K_3$   $K_4$   $K_5$   $K_6$

4% 96%   100%   100%   100%   100%   44% 56%

$B_1^-$

| TH | TE | RA | MA |
|----|----|----|----|
| 0 | 1 | 99925 | 0 |
| 0 |  | 41805 | 0 |
| 0 |  | 14.0·10⁶ | 0 |
|  |  | 0 | 2 |
|  |  | 2.0·10⁶ |  |

$B_2^-$

| TH | TE | RA | MA |
|----|----|----|----|
| 1 | 1 | 282800 | 0 |
| 1 |  | 130000 | 0 |
| 1 |  | 91000 | 0 |
|  |  | 11 | 0 |
|  |  | 8.0·10⁶ |  |

$B_3^-$

| TH | TE | RA | MA |
|----|----|----|----|
| 1 | 1 | 321565 | 0 |
| 1 |  | 505266 | 0 |
| 1 |  | 8.0·10⁶ | 0 |
|  |  | 25656 | 0 |
|  |  | 8.0·10⁶ |  |

$B_4^-$

| TH | TE | RA | MA |
|----|----|----|----|
| 1 | 1 | 2.0·10⁶ | 0 |
| 1 |  | 23.0·10⁶ | 0 |
| 1 |  | 23.0·10⁶ | 0 |
|  |  | 6.0·10⁶ | 0 |
|  |  | 2.0·10⁶ |  |

$B_5^-$

| TH | TE | RA | MA |
|----|----|----|----|
| 1 | 1 | 13205 | 0 |
| 1 |  | 60279 | 0 |
| 1 |  | 13.0·10⁶ | 0 |
|  |  | 2.0·10⁶ | 0 |
|  |  | 3.0·10⁶ |  |

$B_6^-$

| TH | TE | RA | MA |
|----|----|----|----|
| 0 | 1 | 3741 | 1 |
| 1 |  | 23372 | 1 |
| 1 |  | 717921 | 10 |
|  |  | 0 | 0 |
|  |  | 179311 |  |

$B_7^+$

| TH | TE | RA | MA |
|---|---|---|---|
| 0 | 1 | 9 | 526.4·10^6 |
| 0 |  | 2 | 15201 |
| 0 |  | 10 | 45.4·10^6 |
|  |  | 24 | 80.0·10^6 |
|  |  | 4 |  |

$B_8^+$

| TH | TE | RA | MA |
|---|---|---|---|
| 0 | 1 | 6 | 526.4·10^6 |
| 1 |  | 7 | 99786 |
| 0 |  | 7 | 9.0·10^6 |
|  |  | 43 | 80.0·10^6 |
|  |  | 11 |  |

$B_9^+$

| TH | TE | RA | MA |
|---|---|---|---|
| 0 | 1 | 209 | 526.4·10^6 |
| 1 |  | 10 | 33 |
| 0 |  | 4 | 45.4·10^6 |
|  |  | 52 | 80.0·10^6 |
|  |  | 12 |  |

$B_{10}^+$

| TH | TE | RA | MA |
|---|---|---|---|
| 1 | 0 | 46 | 526.4·10^6 |
| 1 |  | 107 | 99786 |
| 0 |  | 22670 | 45.4·10^6 |
|  |  | 112525 | 80.0·10^6 |
|  |  | 70 |  |

$B_{11}^+$

| TH | TE | RA | MA |
|---|---|---|---|
| 1 | 0 | 1 | 9.0·10^6 |
| 1 |  | 5 | 3851 |
| 0 |  | 0 | 1.0·10^6 |
|  |  | 2 | 4.0·10^6 |
|  |  | 1 |  |

$B_{12}^+$

| TH | TE | RA | MA |
|---|---|---|---|
| 1 | 1 | 624 | 526.4·10^6 |
| 1 |  | 2663 | 99786 |
| 0 |  | 796 | 45.4·10^6 |
|  |  | 4617 | 80.0·10^6 |
|  |  | 9921 |  |

46%  54%   23%  77%   1%  99%   14%  86%   43%  57%   11%  89%

$K_7$  $K_8$  $K_9$  $K_{10}$  $K_{11}$  $K_{12}$

7%  93%   100%   100%   25%  75%   100%   65%  35%

$B_7^-$

| TH | TE | RA | MA |
|---|---|---|---|
| 0 | 1 | 297 | 3 |
| 1 |  | 115510 | 0 |
| 1 |  | 25.0·10^6 | 5 |
|  |  | 1.0·10^6 | 5 |
|  |  | 13.5·10^6 |  |

$B_8^-$

| TH | TE | RA | MA |
|---|---|---|---|
| 1 | 1 | 1.0·10^6 | 0 |
| 1 |  | 12.0·10^6 | 0 |
| 1 |  | 12.0·10^6 | 0 |
|  |  | 11.0·10^6 | 0 |
|  |  | 3.0·10^6 |  |

$B_9^-$

| TH | TE | RA | MA |
|---|---|---|---|
| 1 | 1 | 7.0·10^6 | 0 |
| 1 |  | 1.0·10^6 | 0 |
| 1 |  | 3.0·10^6 | 0 |
|  |  | 960000 | 0 |
|  |  | 1.0·10^6 |  |

$B_{10}^-$

| TH | TE | RA | MA |
|---|---|---|---|
| 1 | 0 | 7.5·10^6 | 11 |
| 0 |  | 23.5·10^6 | 9 |
| 0 |  | 453.7·10^6 | 11 |
|  |  | 17.0·10^6 | 11 |
|  |  | 13.5·10^6 |  |

$B_{11}^-$

| TH | TE | RA | MA |
|---|---|---|---|
| 1 | 0 | 1.0·10^6 | 2 |
| 1 |  | 555137 | 1 |
| 1 |  | 408.0·10^6 | 1 |
|  |  | 13.0·10^6 | 0 |
|  |  | 13.0·10^6 |  |

$B_{12}^-$

| TH | TE | RA | MA |
|---|---|---|---|
| 0 | 0 | 1.4·10^6 | 11 |
| 1 |  | 3584 | 10 |
| 1 |  | 1132 | 10 |
|  |  | 6189 | 10 |
|  |  | 13563 |  |