

A dissimilarity approach to the characterization of datasets

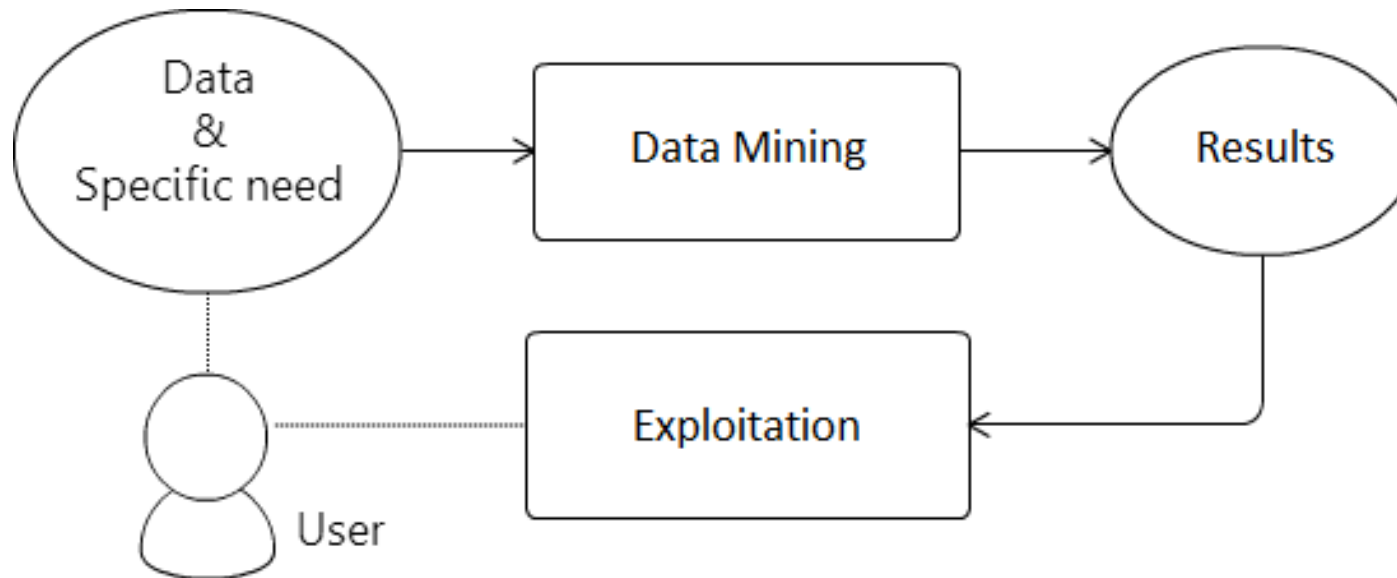
William Raynaut

Chantal Soule-Dupuy, Nathalie Valles-Parlangeau

Cedric Dray, Philippe Valet



Introduction



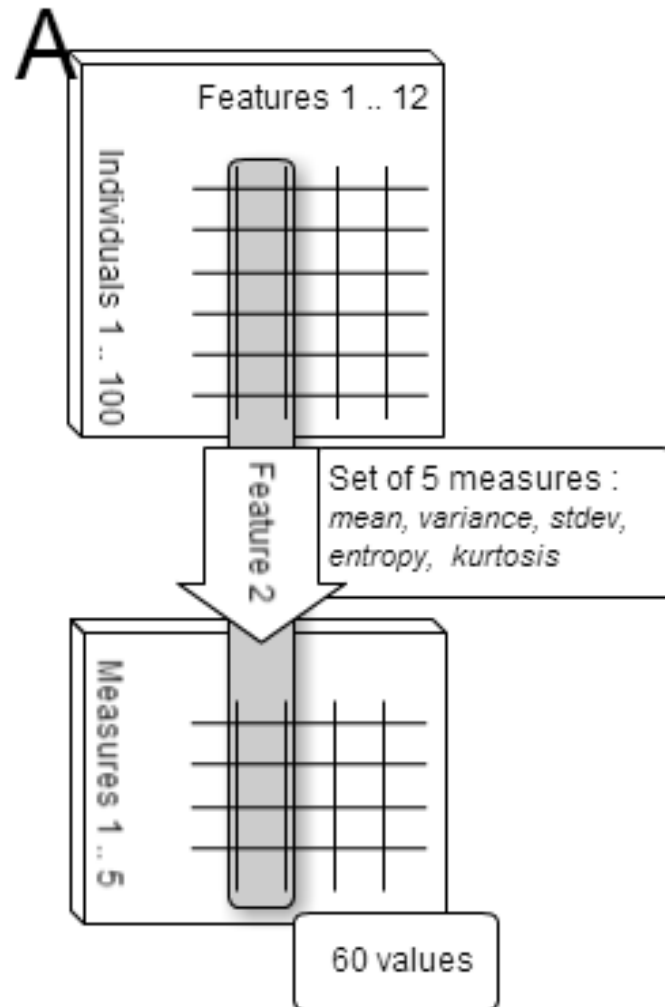
Data analysis is a **complex** process.
The typical user needs **assistance**.

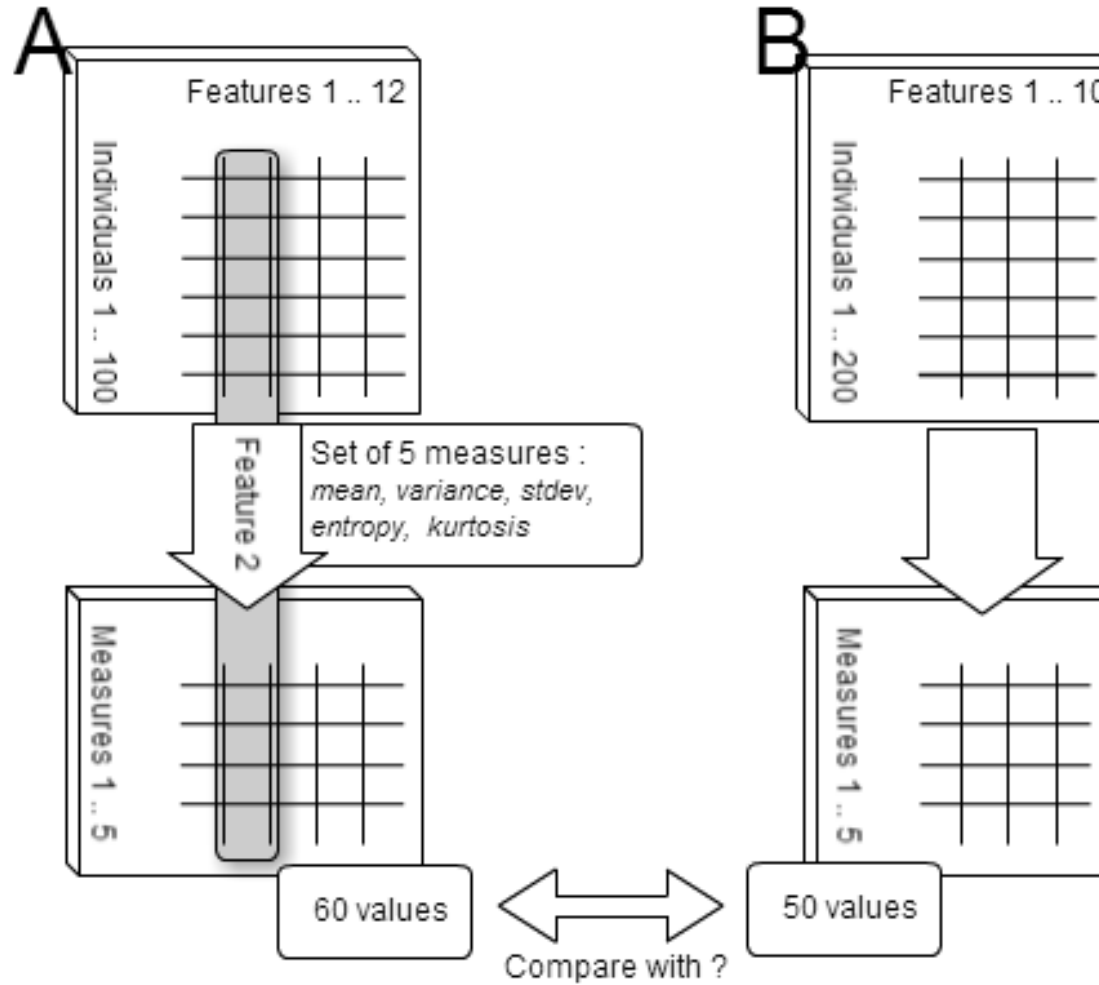
⇒ many **Intelligent Discovery Assistants (IDA)**

- ▶ Expert systems
[Sleeman et al., 1995]
- ▶ Meta-learning methods
[Kalousis, 2002, Giraud-Carrier et al., 2004, Sun and Pfahringer, 2013]
- ▶ Case bases
[Goble et al., 2010]
- ▶ Heuristic search methods
[Escalante et al., 2009, Sun, 2014]
- ▶ DM workflows planners
[Zakova et al., 2011, Serban, 2013, Nguyen et al., 2014]

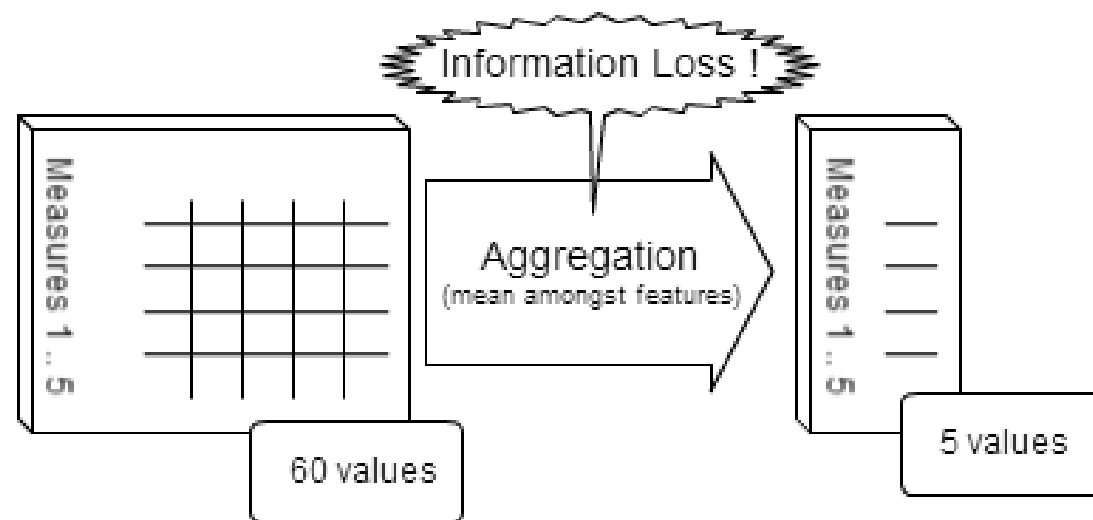
- ▶ IDAs help the user design processes or workflows that are **appropriate** for his **dataset** and objectives.
 - ▶ This implies a **characterization** of the user's dataset that fits the particular mechanics of the IDA.
- We investigate the possibility that limitations in the classical characterization of datasets are among the main obstacles to well performing IDAs

Dataset characterization

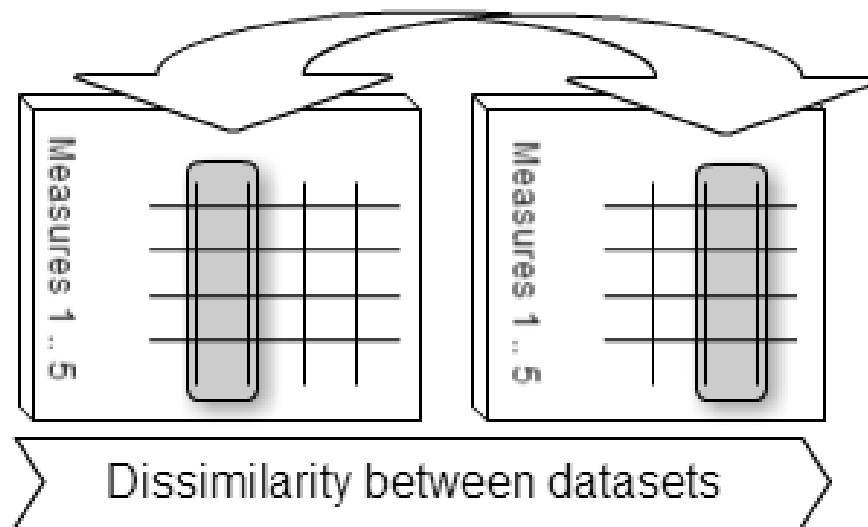




The usual solution is to average the measures along the features.



We propose to compare the features by most similar pairs.



Dissimilarity function

$$d_{\omega}^{ubr}(x, y) = \frac{d_{G(\omega)}^{ubr}(x, y)}{\max_{(x', y') \in \omega^2} (d_{G(\omega)}^{ubr}(x', y'))} + \frac{d_{F(\omega)}(x, y)}{\max_{(x', y') \in \omega^2} (d_{F(\omega)}(x', y'))}$$

The proposed dissimilarity take into account :

- ▶ Global properties of the datasets such as the number of instances and attributes, or the performance of landmarks (simple learners) on them...
- ▶ Properties of the individual features of the datasets, such as entropy or kurtosis...

And compares the features of the datasets by most similar pairs.

Theoretical validation

- ▶ We use the *strong* (ϵ, γ) -goodness from [Wang et al., 2009]
- ▶ Characterizes the usefulness of a dissimilarity function for learning a concept from a dataset (here we try to learn the *appropriateness* of classifiers to different datasets)
- ▶ Gives upper bound of error rate for a simple classifier using the dissimilarity

| | 200 examples | | 1000 examples | | 5000 examples | |
|--------------------|--------------|-------------|---------------|-------------|---------------|-------------|
| | δ | error bound | δ | error bound | δ | error bound |
| d_{ω}^{ubr} | 0,973 | 1,023 | 0,871 | 0,921 | 0,501 | 0,551 |
| Euclidean based | 0,989 | 1,039 | 0,945 | 0,995 | 0,755 | 0,805 |
| Manhattan based | 0,990 | 1,040 | 0,952 | 1,002 | 0,783 | 0,833 |

| | 10000 examples | | 20000 examples | | 50000 examples | |
|--------------------|----------------|-------------|----------------|-------------|----------------|-------------|
| | δ | error bound | δ | error bound | δ | error bound |
| d_{ω}^{ubr} | 0,251 | 0,301 | 0,063 | 0,113 | 0,001 | 0,051 |
| Euclidean based | 0,570 | 0,620 | 0,325 | 0,375 | 0,060 | 0,110 |
| Manhattan based | 0,613 | 0,663 | 0,375 | 0,425 | 0,086 | 0,136 |

Table: Error bound achievable with probability $1 - \delta$ by dissimilarity based classifiers for different numbers of examples

Experimental validation

Extracting data from OpenML :

- ▶ Accuracy of a set of classifiers over a range of datasets

| | <i>classifier₁</i> | <i>classifier₂</i> | ... | <i>classifier₉₃</i> |
|------------------------------|-------------------------------|-------------------------------|-----|--------------------------------|
| <i>dataset₁</i> | 0.8 | 0.9 | ... | ... |
| <i>dataset₂</i> | 0.9 | 0.7 | ... | ... |
| ... | ... | ... | ... | ... |
| <i>dataset₄₃₄</i> | ... | ... | ... | ... |

Experimental validation

Extracting data from OpenML :

- ▶ Accuracy of a set of classifiers over a range of datasets

| | <i>classifier₁</i> | <i>classifier₂</i> | ... | <i>classifier₉₃</i> |
|------------------------------|-------------------------------|-------------------------------|-----|--------------------------------|
| <i>dataset₁</i> | 0.8 | 0.9 | ... | ... |
| <i>dataset₂</i> | 0.9 | 0.7 | ... | ... |
| ... | ... | ... | ... | ... |
| <i>dataset₄₃₄</i> | ... | ... | ... | ... |

- ▶ Characterization of the datasets

| | <i>NumberOfInstances</i> | <i>NumberOfFeatures</i> | ... | <i>MetaAttribute₁₀₅</i> |
|------------------------------|--------------------------|-------------------------|-----|------------------------------------|
| <i>dataset₁</i> | 100 | 62 | ... | ... |
| <i>dataset₂</i> | 5000 | 13 | ... | ... |
| ... | ... | ... | ... | ... |
| <i>dataset₄₃₄</i> | 360 | 20 | ... | ... |

Building a Meta-dataset for a classification meta-problem

| | <i>NumberOfInstances</i> | ... | <i>MetaAttribute</i> ₁₀₅ | <i>Class</i> |
|-------------------------------|--------------------------|-----|-------------------------------------|---------------------------------|
| <i>dataset</i> ₁ | 100 | ... | 4 | <i>classifier</i> ₁₈ |
| <i>dataset</i> ₂ | 5000 | ... | 92 | <i>classifier</i> ₇ |
| ... | ... | ... | ... | ... |
| <i>dataset</i> ₄₃₄ | 360 | ... | 13 | <i>classifier</i> ₆₃ |

foreach *Dataset instance dataset*; **do**

Exclude *dataset*; from the metadataset

Apply an attribute selection algorithm on the metadataset

Learn a classification model from the reduced metadataset
using a classification algorithm

Use this model to predict a class label *classifier*; for *dataset*;

x = actual criterion value achieved on *dataset*; by *classifier*;

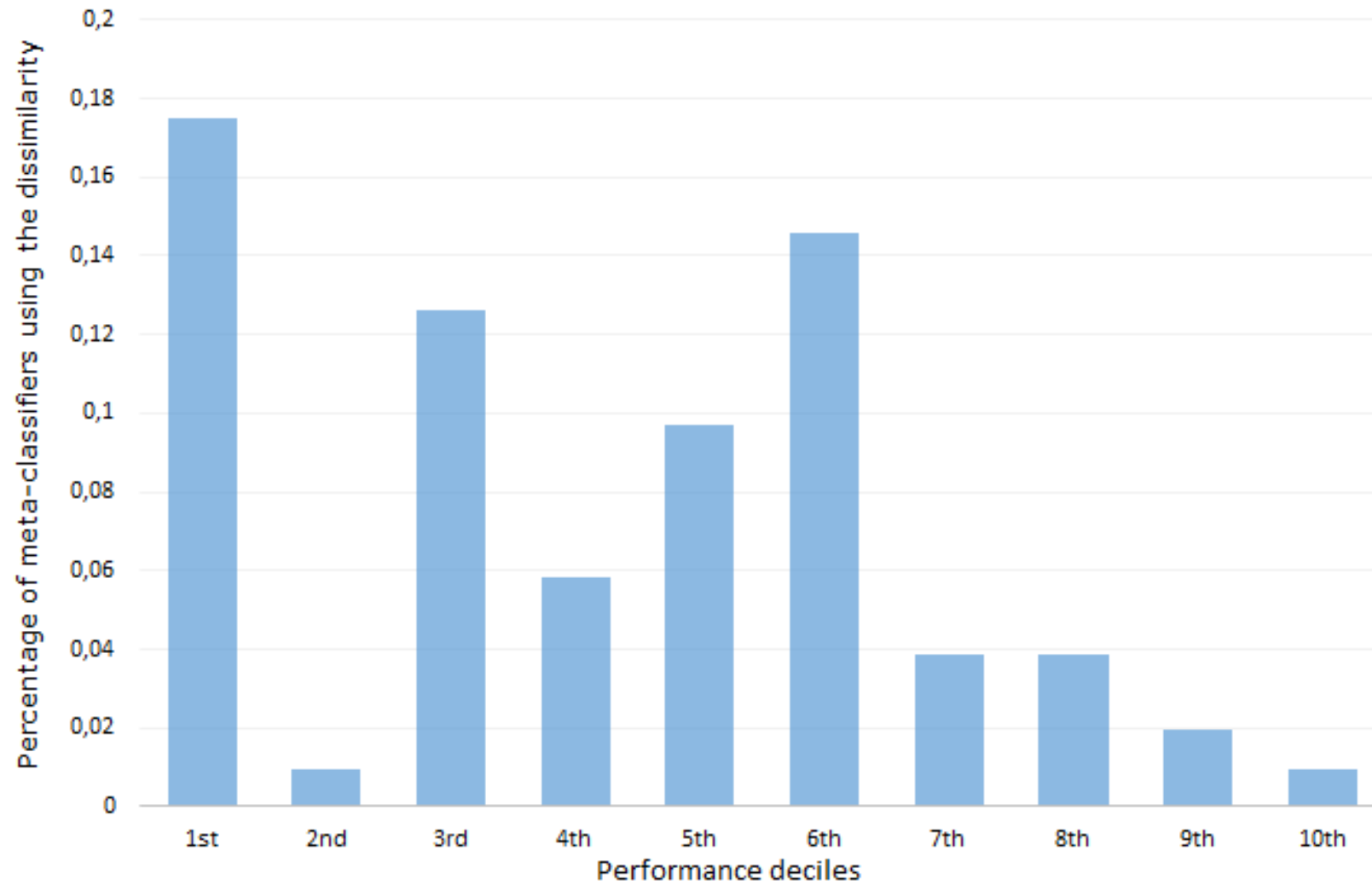
best = best criterion value achieved on *dataset*; among the classifiers
*classifier*_{1..m}

def = actual criterion value achieved on *dataset*; by the default
classifier (majority class classifier)

$$perf(run) = 1 - \frac{|best - x|}{|best - def|}$$

- ▶ 434 datasets
 - ▶ 11 evaluation criteria
 - ▶ > 2600 (meta)-classifiers
 - ▶ > 60 (meta)-feature selection methods
- ⇒ over 200k experiments, 700M ML algorithms runs

How does the dissimilarity perform relatively to all other approaches ?



Contribution

- ▶ A dissimilarity accounting for the inner topology of the datasets attributes
- ▶ With theoretical guarantees when characterizing the appropriateness of classifiers
- ▶ Able to perform well in a meta-learning framework

Thank you for your attention

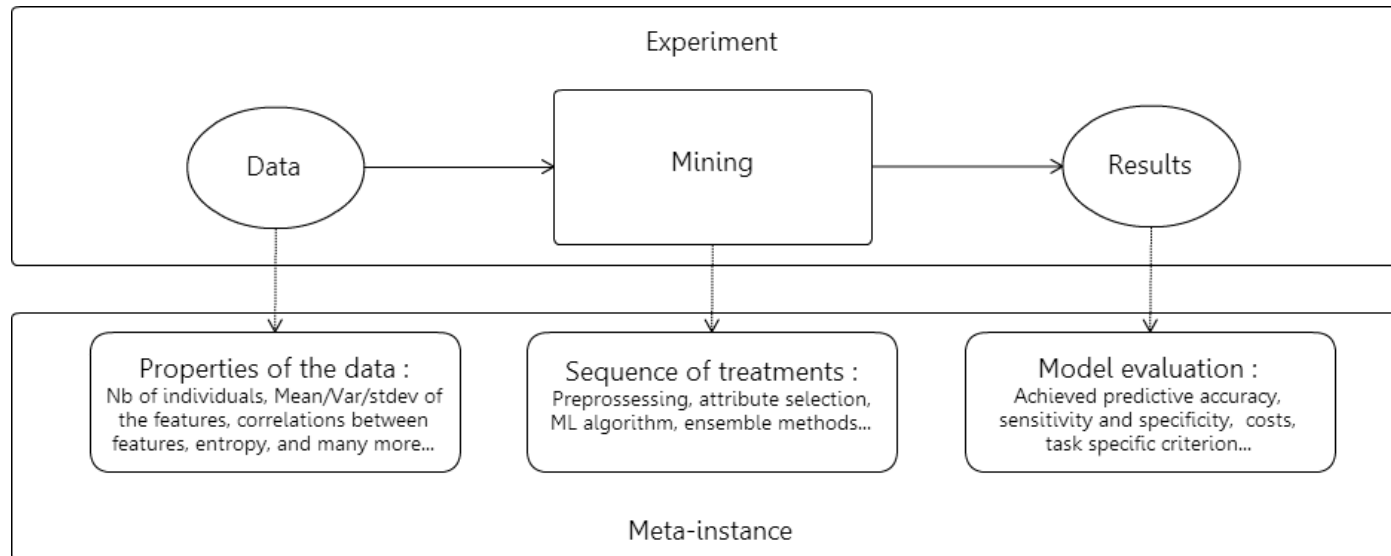
$$D = \{ \text{classification datasets} \}, A = \{ \text{classifiers} \}$$

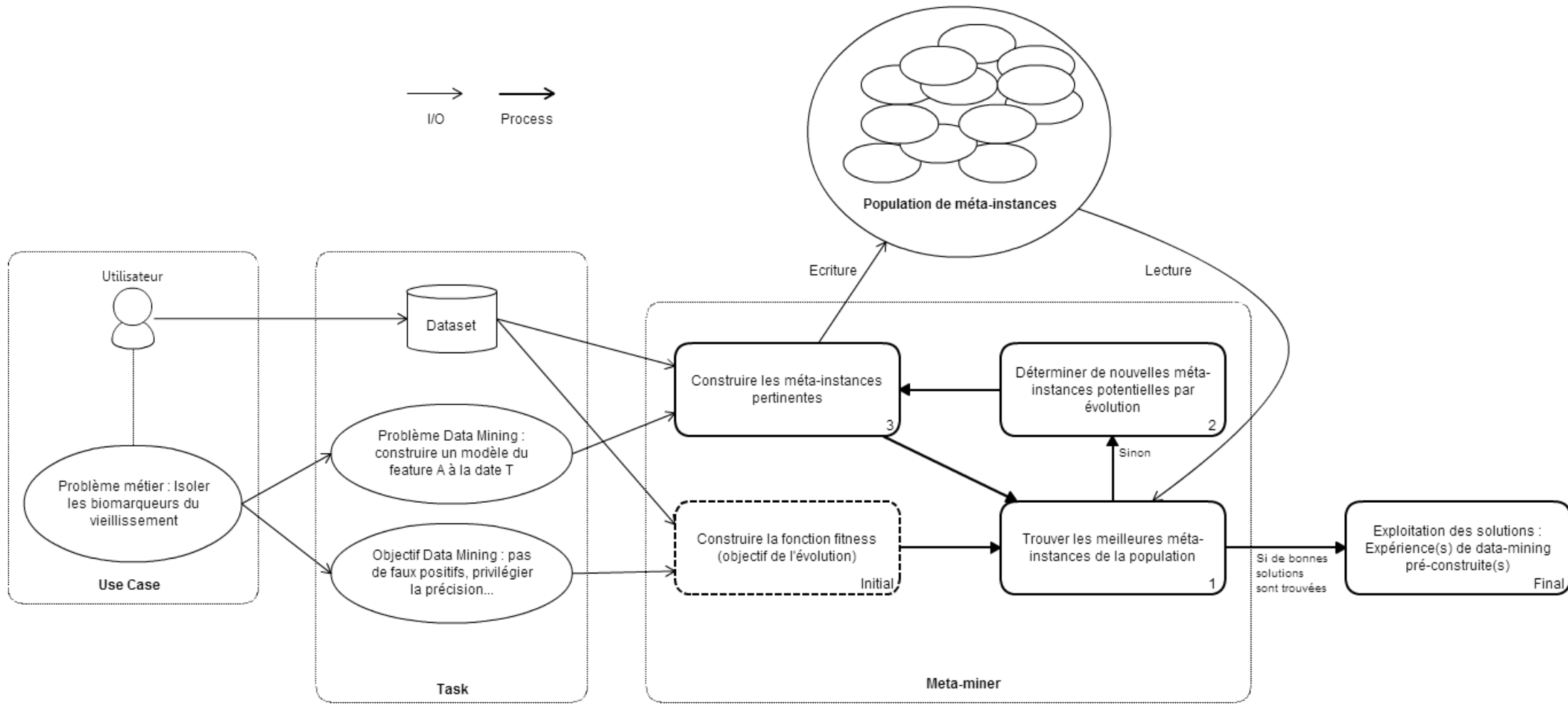
- ▶ We execute every classifier of A on every dataset of D and measure a performance criterion c of the resulting model.
- ▶ For each dataset x , we define the set A_x of the algorithms that are appropriate on this dataset along our performance criterion as those at most one standard deviation away from the best :

$$A_x = \{ a \in A \text{ such that } | \max_{a' \in A} (c(a', x)) - c(a, x) | \leq \sigma_x \}$$

We can then consider, for each algorithm $a \in A$, the binary classification problem where instances are the datasets $x \in D$, and their class label stating whether a is appropriate on them.

1 Experiment = 1 Meta-instance





References I



Escalante, H. J., Montes, M., and Sucar, L. E. (2009).
Particle swarm model selection.
The Journal of Machine Learning Research, 10:405–440.



Giraud-Carrier, C., Vilalta, R., and Brazdil, P. (2004).
Introduction to the special issue on meta-learning.
Machine learning, 54(3):187–193.



Goble, C. A., Bhagat, J., Aleksejevs, S., Cruickshank, D., Michaelides, D., Newman, D., Borkum, M.,
Bechhofer, S., Roos, M., Li, P., et al. (2010).
myexperiment: a repository and social network for the sharing of bioinformatics workflows.
Nucleic acids research, 38(suppl 2):W677–W682.



Kalousis, A. (2002).
Algorithm selection via meta-learning.
PhD thesis, Universite de Geneve.



Nguyen, P., Hilario, M., and Kalousis, A. (2014).
Using meta-mining to support data mining workflow planning and optimization.
Journal of Artificial Intelligence Research, pages 605–644.



Serban, F. (2013).
Toward effective support for data mining using intelligent discovery assistance.
PhD thesis.

References II



Sleeman, D., Rissakis, M., Craw, S., Graner, N., and Sharma, S. (1995).
Consultant-2: Pre-and post-processing of machine learning applications.



Sun, Q. (2014).
Meta-learning and the full model selection problem.
Unpublished PhD thesis, University of Waikato.



Sun, Q. and Pfahringer, B. (2013).
Pairwise meta-rules for better meta-learning-based algorithm ranking.
Machine learning, 93(1):141–161.



Wang, L., Sugiyama, M., Yang, C., Hatano, K., and Feng, J. (2009).
Theory and algorithm for learning with dissimilarity functions.
Neural computation, 21(5):1459–1484.



Zakova, M., Kremen, P., Zelezny, F., and Lavrac, N. (2011).
Automating knowledge discovery workflow composition through ontology-based planning.
Automation Science and Engineering, IEEE Transactions on, 8(2):253–264.