



REFINE&MINE: ANYTIME SUBGROUP DISCOVERY IN NUMERICAL DOMAINS WITH GUARANTEES

AUTHORS.

- Aimene BELFODIL ● Adnene BELFODIL
- Mehdi KAYTOUE

\°/ All icons used are courtesy of flaticon.com

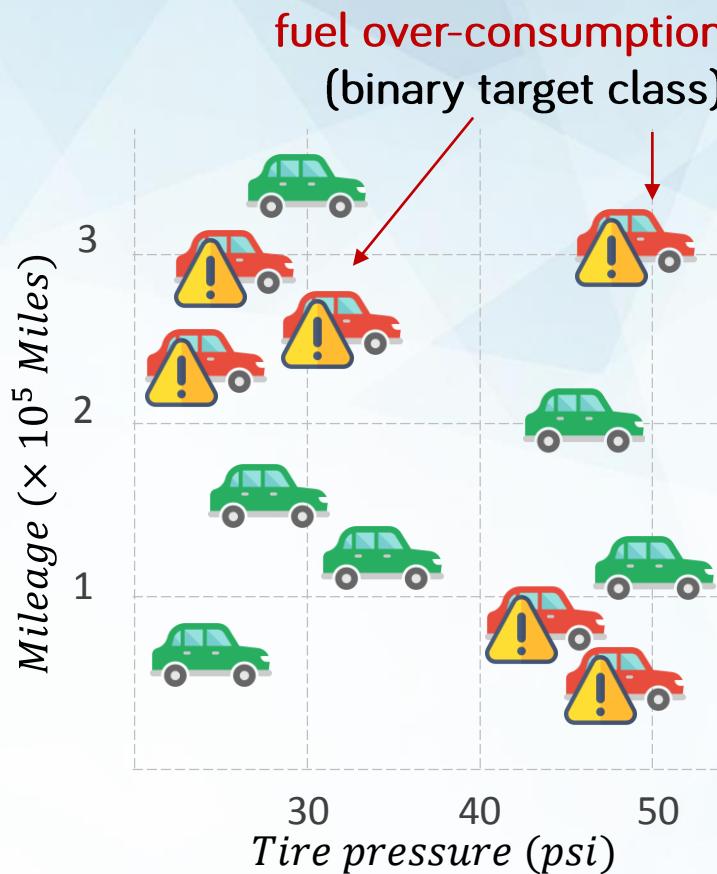


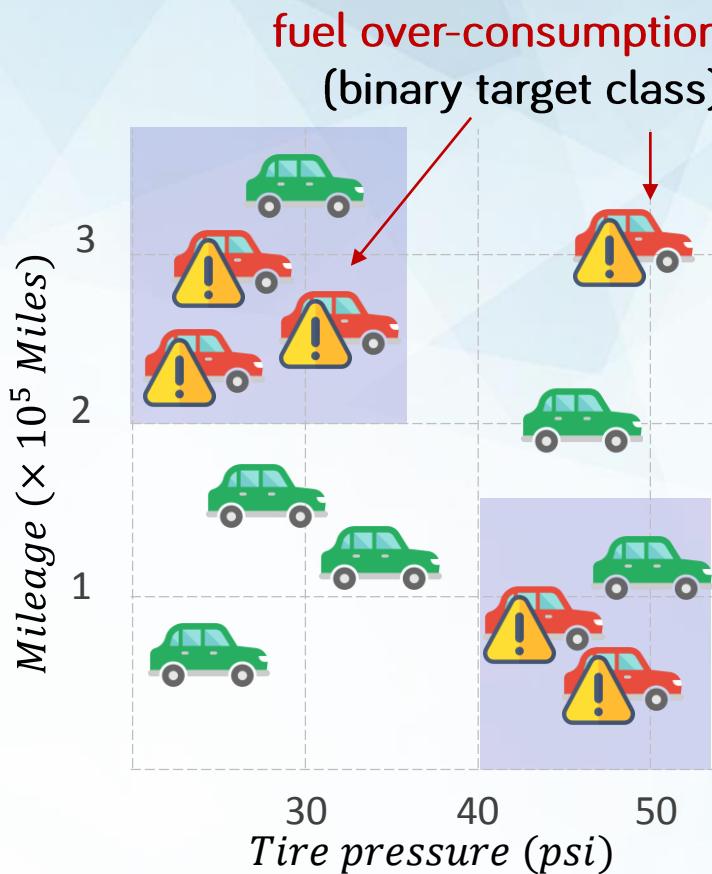


INTRODUCTION

PROBLEM DEFINITION

Subgroup Discovery



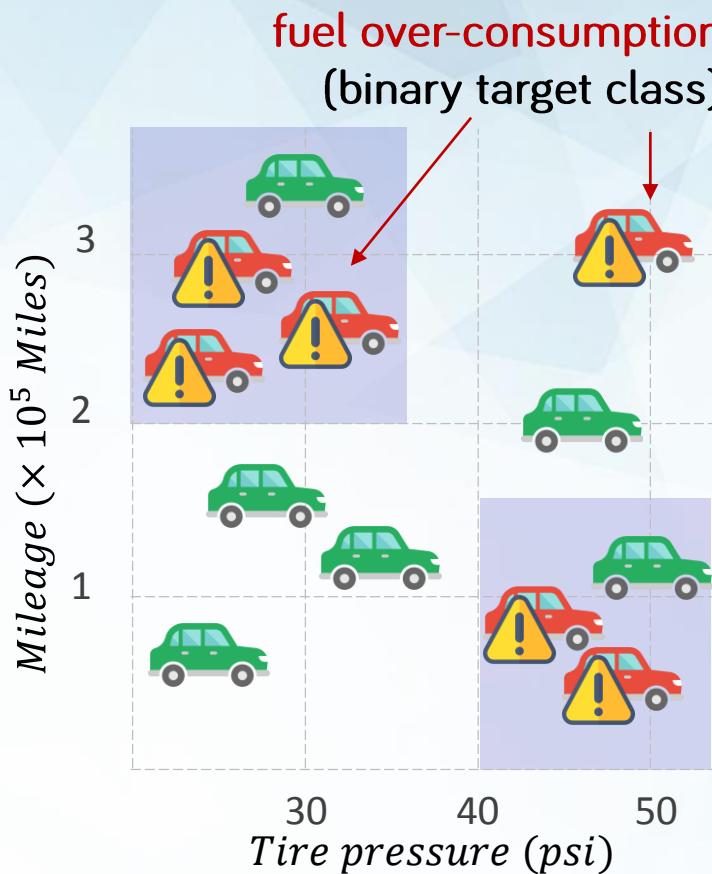


Hypothesis.

Fuel over-consumption is fostered when:

1. *Tire pressure ≤ 35 psi and Milage $\geq 2 \times 10^5$ miles*
2. *Tire pressure ≥ 40 psi and Milage $\leq 1.7 \times 10^5$ miles*

Subgroups - Discriminant Interval patterns



Hypothesis.

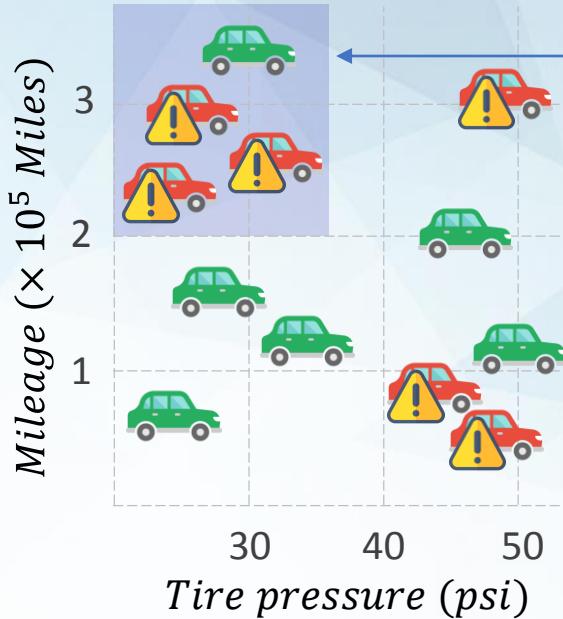
Fuel over-consumption is fostered when:

1. *Tire pressure ≤ 35 psi and Milage $\geq 2 \times 10^5$ miles*
2. *Tire pressure ≥ 40 psi and Milage $\leq 1.7 \times 10^5$ miles*

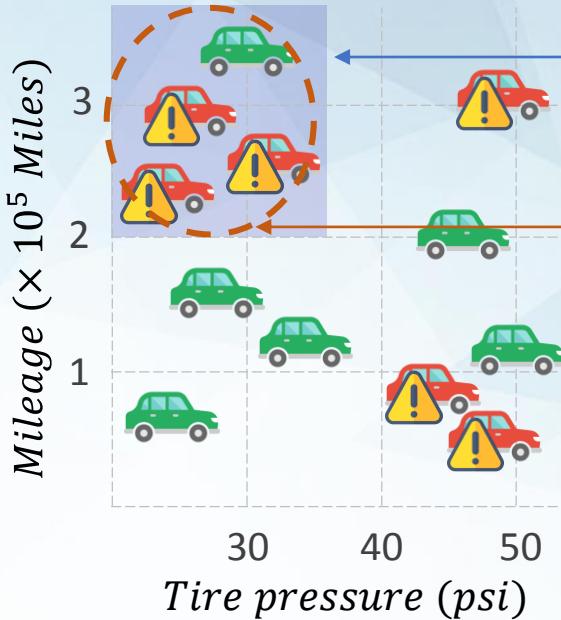
Subgroups - Discriminant Interval patterns

Subgroup discovery. The task of discovering patterns that accurately discriminate the target class from others.

*We are interested only in datasets where all the attributes are numerical.



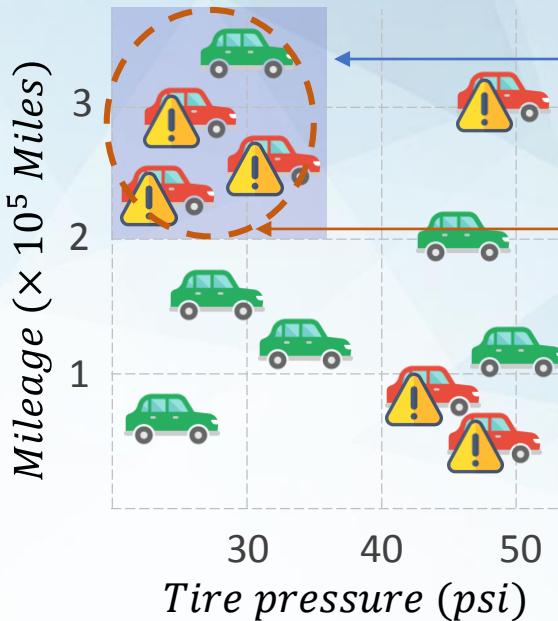
Interval pattern intent. Conjunction of interval restrictions:
E.g. $x \in [0, 35]$ and $y \in [2, +\infty)$



Interval pattern intent. Conjunction of interval restrictions:

E.g. $x \in [0, 35]$ and $y \in [2, +\infty)$

Interval pattern extent. Set of instances for which the pattern hold.

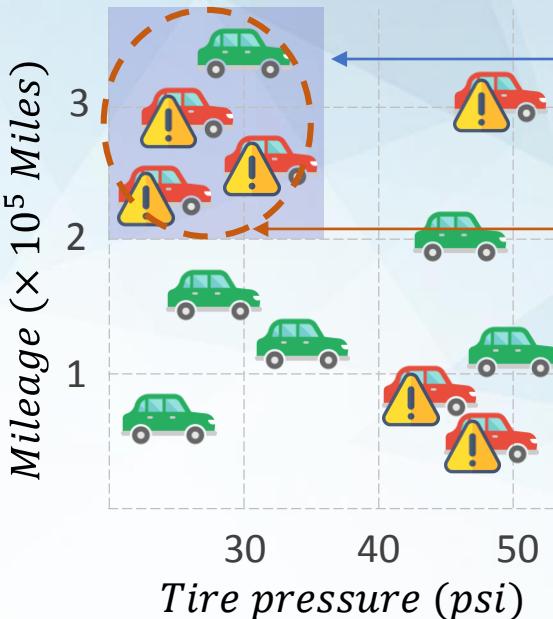


Interval pattern intent. Conjunction of interval restrictions:

E.g. $x \in [0, 35]$ and $y \in [2, +\infty)$

Interval pattern extent. Set of instances for which the pattern hold.

- A **Quality measure ϕ** evaluates **the discriminative power of patterns**.



Interval pattern intent. Conjunction of interval restrictions:

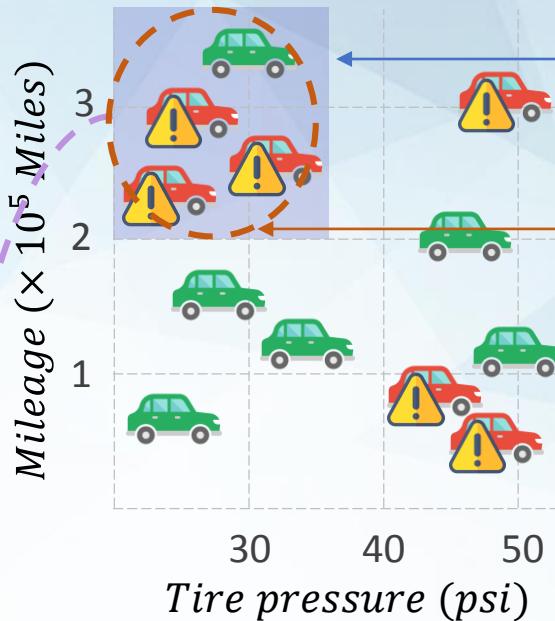
E.g. $x \in [0, 35]$ and $y \in [2, +\infty)$

Interval pattern extent. Set of instances for which the pattern hold.

- A **Quality measure** ϕ evaluates the **discriminative power** of patterns.
- We use **measures increasing** with tpr^* and **decreasing** with the fpr^* .

* True Positive Rate

* False Positive Rate



Interval pattern intent. Conjunction of interval restrictions:

E.g. $x \in [0, 35]$ and $y \in [2, +\infty)$

Interval pattern extent. Set of instances for which the pattern hold.

- A **Quality measure ϕ** evaluates **the discriminative power of patterns**.
- We use **measures increasing with tpr^*** and **decreasing with the fpr^*** .

* True Positive Rate

* False Positive Rate

$$\phi(p) = \frac{3}{6} - \frac{1}{6} = 0.33$$

$$\phi(\text{pattern}) = tpr(\text{pattern}) - fpr(\text{pattern})$$

***Informedness** is order-equivalent to the well-known **WRAcc measure**.



Exhaustive Search

Generally **unfeasible**.

e.g. A 200×3 dataset has up to **8 Billions Interval Patterns**.

M. Kaytoue, S. O. Kuznetsov, A. Napoli.
Revisiting Numerical Pattern Mining with FCA. In IJCAI 2011.



Exhaustive Search

Generally **unfeasible**.

e.g. A 200×3 dataset has up to **8 Billions Interval Patterns**.

M. Kaytoue, S. O. Kuznetsov, A. Napoli.
Revisiting Numerical Pattern Mining with FCA. In IJCAI 2011.

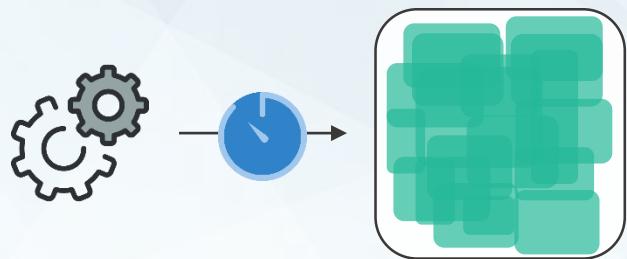


Prior Discretization

Loss of information.

e.g. overlapping intervals are ignored. Look for [1,2] and [2,3] patterns but **ignore [1,3]**.

M. Artzmueller, F. Puppe. **SD-Map – A Fast Algorithm for Exhaustive Subgroup Discovery**. In PKDD 2006.

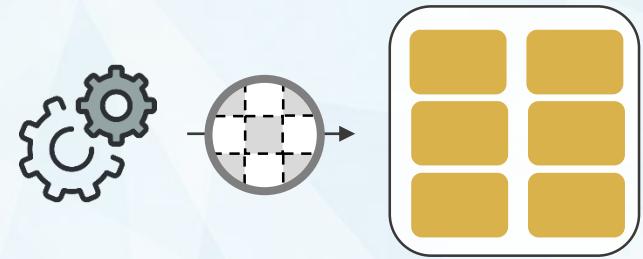


Exhaustive Search

Generally **unfeasible**.

e.g. A 200×3 dataset has up to **8 Billions Interval Patterns**.

M. Kaytoue, S. O. Kuznetsov, A. Napoli.
Revisiting Numerical Pattern Mining with FCA. In IJCAI 2011.

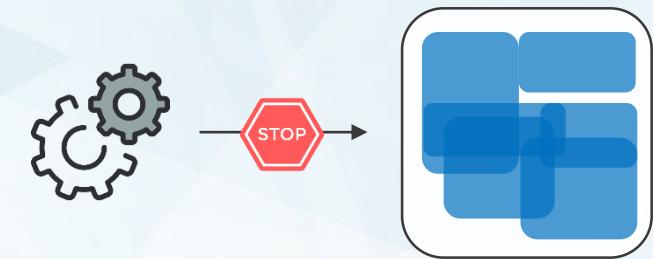


Prior Discretization

Loss of information.

e.g. overlapping intervals are ignored. Look for [1,2] and [2,3] patterns but **ignore** [1,3].

M. Artzmueller, F. Puppe. **SD-Map – A Fast Algorithm for Exhaustive Subgroup Discovery**. In PKDD 2006.

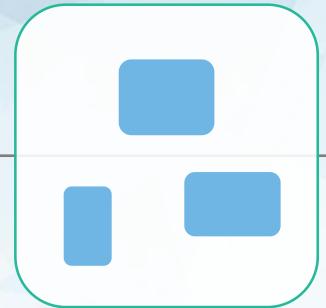


Heuristics and sampling

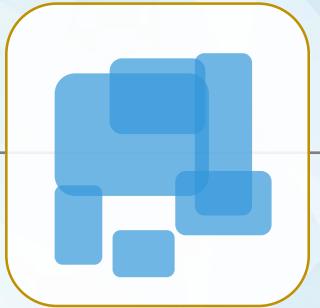
No guarantee on the outputted patterns **at the end** or upon **interruption**

G. Bosc, J.F. Boulicaut, C. Raïssi, M. Kaytoue. **Anytime discovery of a diverse set of patterns with MCTS**. In DMKD 2018.

Algorithm searching for discriminant interval patterns ...



Few patterns
after 5"

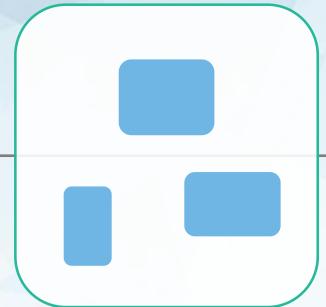


More patterns
after 1'

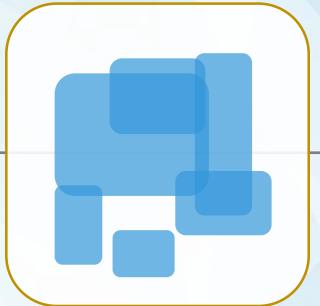


Time

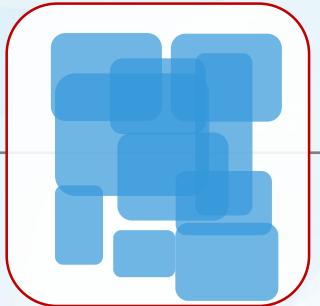
Algorithm searching for discriminant interval patterns ...



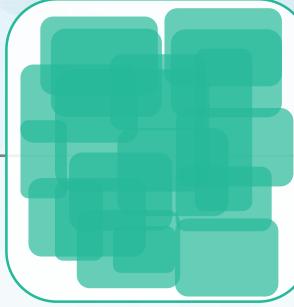
Few patterns
after 5"



More patterns
after 1'



Many patterns
after 20'

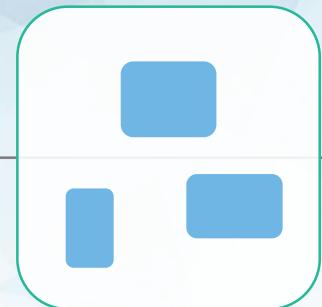


Some or All
patterns after
the End

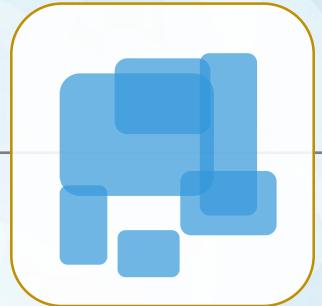


Time

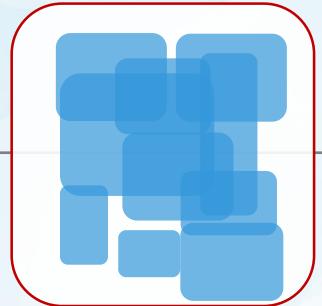
Algorithm searching for discriminant interval patterns ...



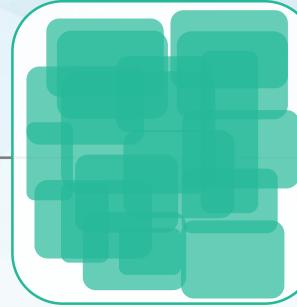
Few patterns
after 5"



More patterns
after 1'



Many patterns
after 20'



Some or All
patterns after
the End



Time



All state-of-the-art subgroup discovery algorithms provide **no guarantees** on the **outputted discriminant interval patterns** when **interrupted during execution**



Propose the new algorithm

Refine&Mine



Anytime
Algorithm

Provide patterns **anytime**
with **two guarantees** on them





Propose the new algorithm

Refine&Mine



Anytime Algorithm

Provide patterns **anytime** with **two guarantees** on them



Guarantee 1: Accuracy

Bound the quality difference between the **best possible pattern** and the best **already found pattern**



Propose the new algorithm

Refine&Mine



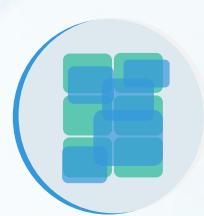
Anytime Algorithm

Provide patterns **anytime** with **two guarantees** on them
* The **more time** it is given, the lower are **the bounds**.



Guarantee 1: Accuracy

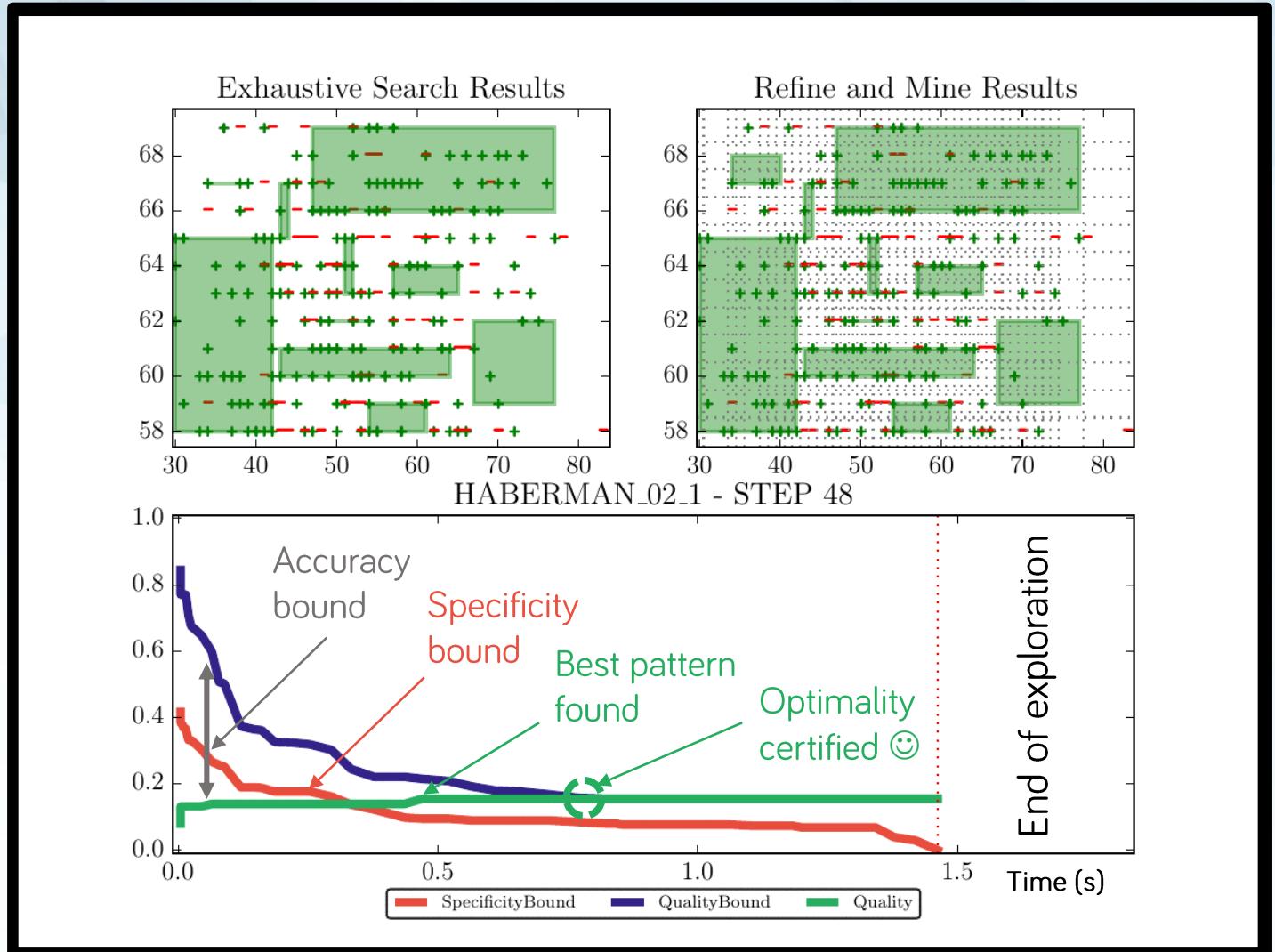
Bound the quality difference between the **best possible pattern** and the best **already found pattern**

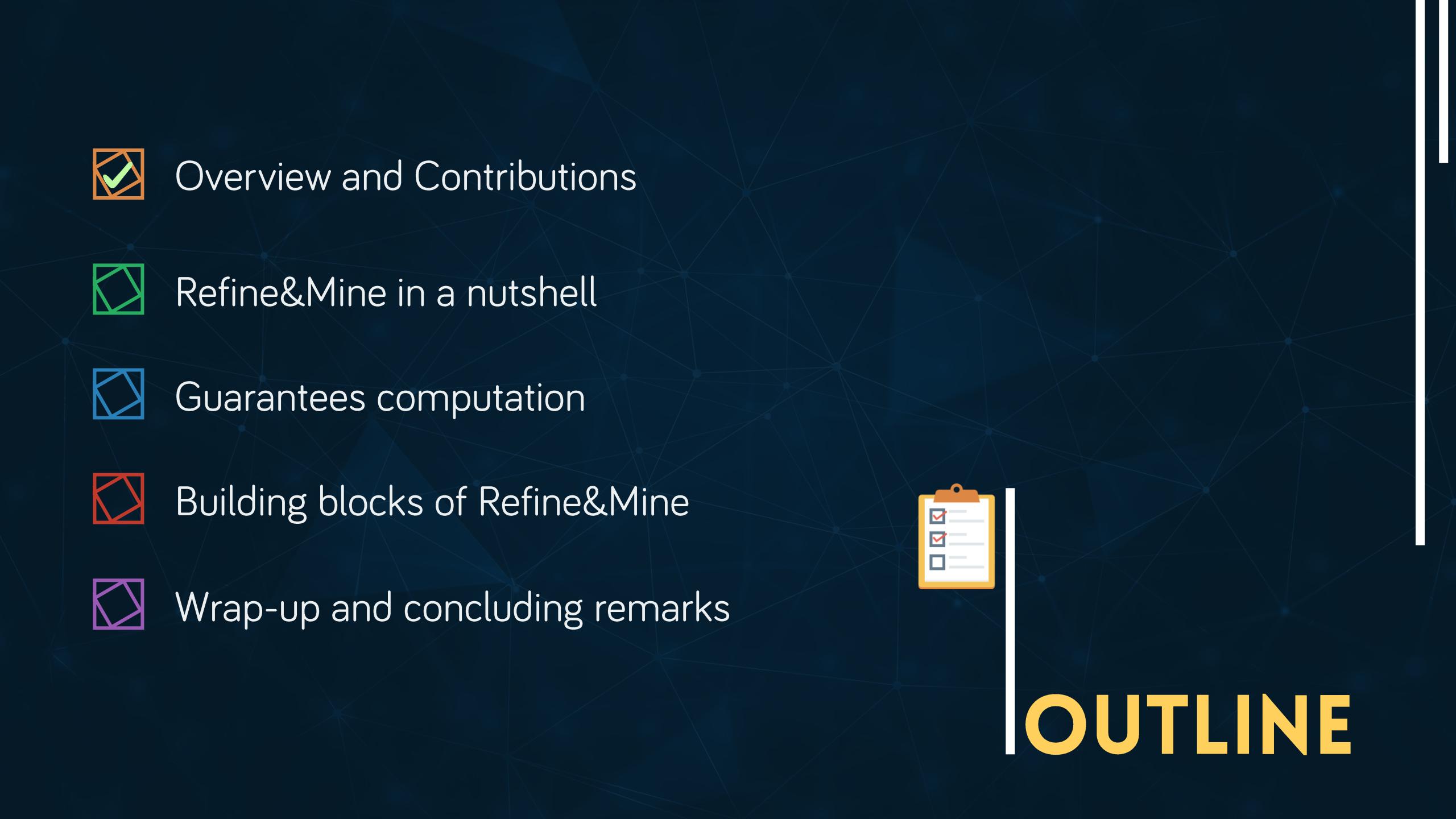


Guarantee 2: Specificity

Bound the distance of **found patterns** to **ground truth patterns** to ensure diversity of found patterns.

Contributions (2)



- 
-  Overview and Contributions
 -  Refine&Mine in a nutshell
 -  Guarantees computation
 -  Building blocks of Refine&Mine
 -  Wrap-up and concluding remarks



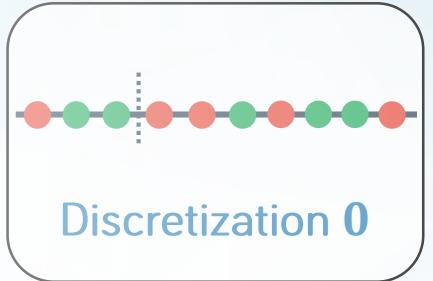
OUTLINE



| **REFINE&MINE**

Refine&Mine

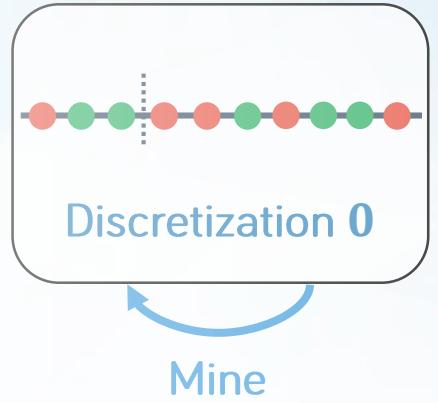
Sequence of discretization **refinements** and mining patterns within the **discretization**.



Discretization 0

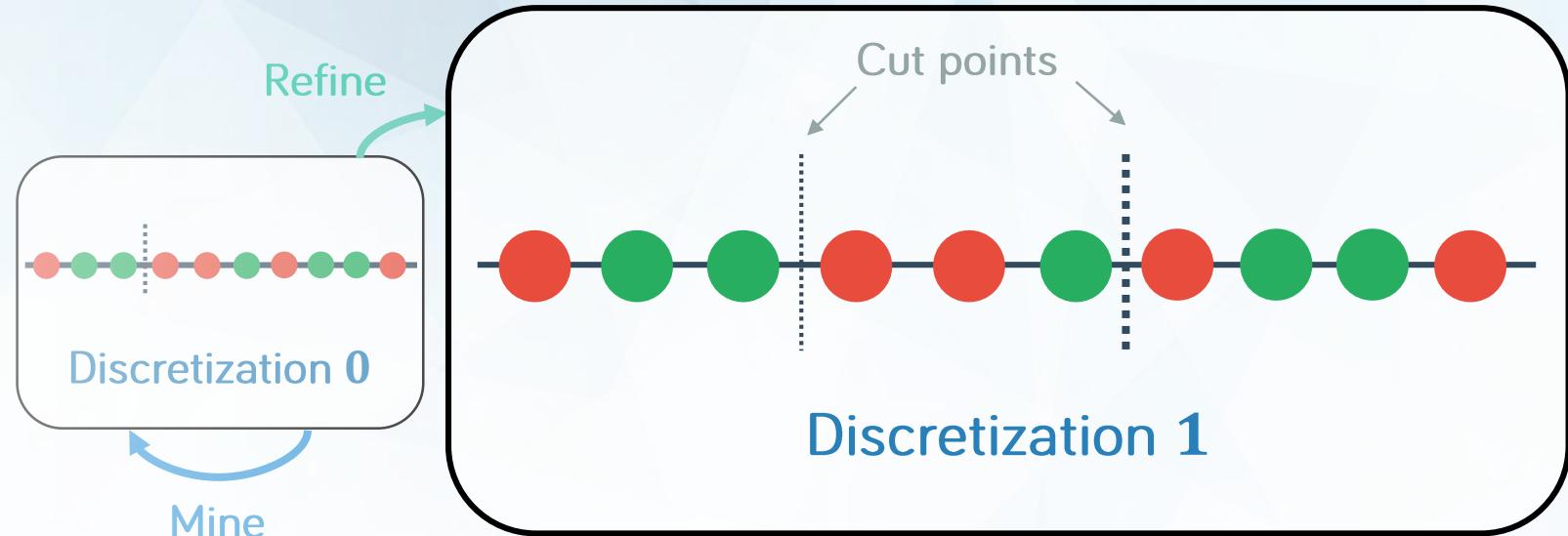
Refine&Mine

Sequence of discretization **refinements** and mining patterns within the **discretization**.



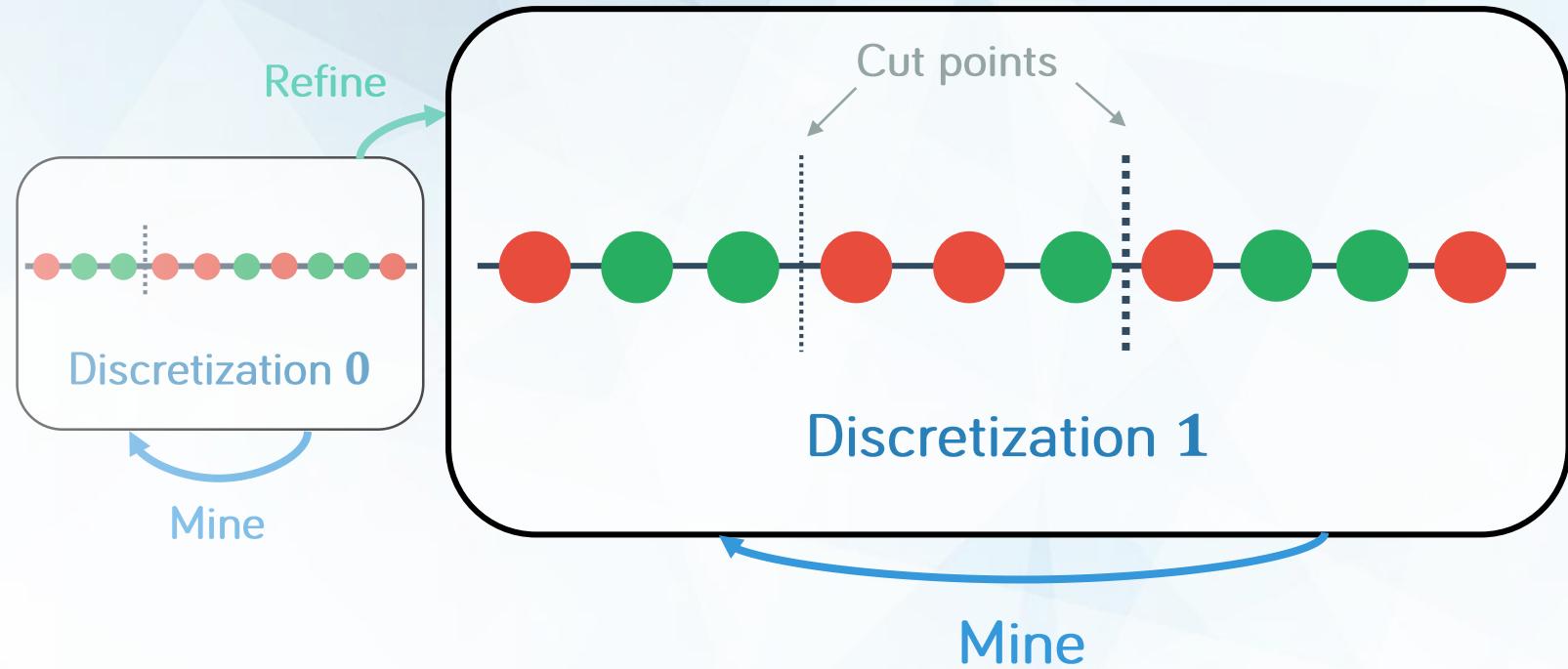
Refine&Mine

Sequence of discretization **refinements** and mining patterns within the **discretization**.



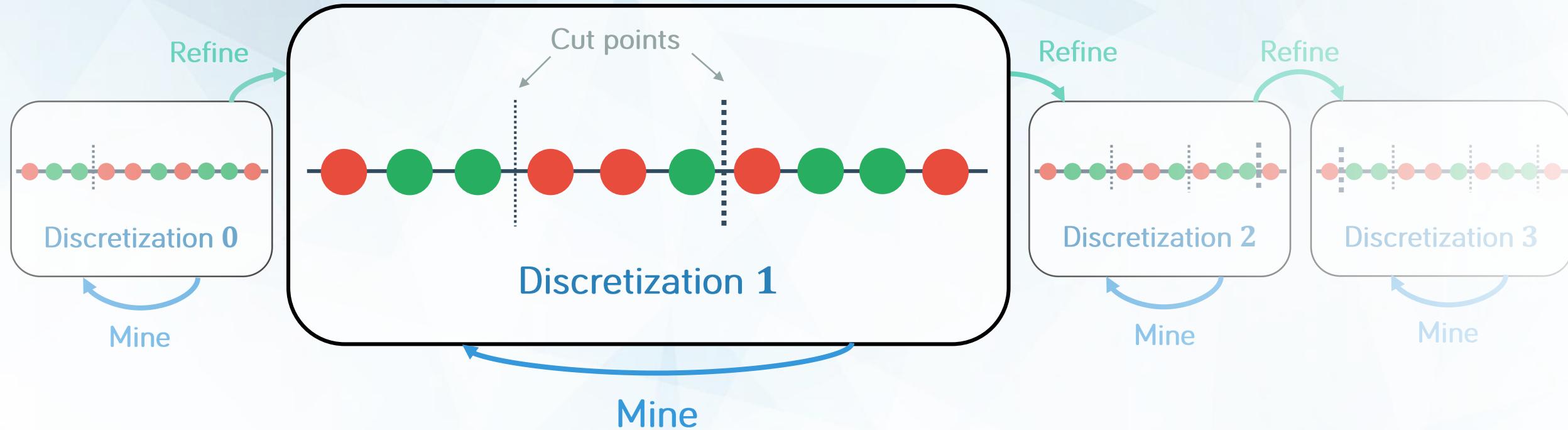
Refine&Mine

Sequence of discretization **refinements** and mining patterns within the **discretization**.



Refine&Mine

Sequence of discretization **refinements** and mining patterns within the **discretization**.

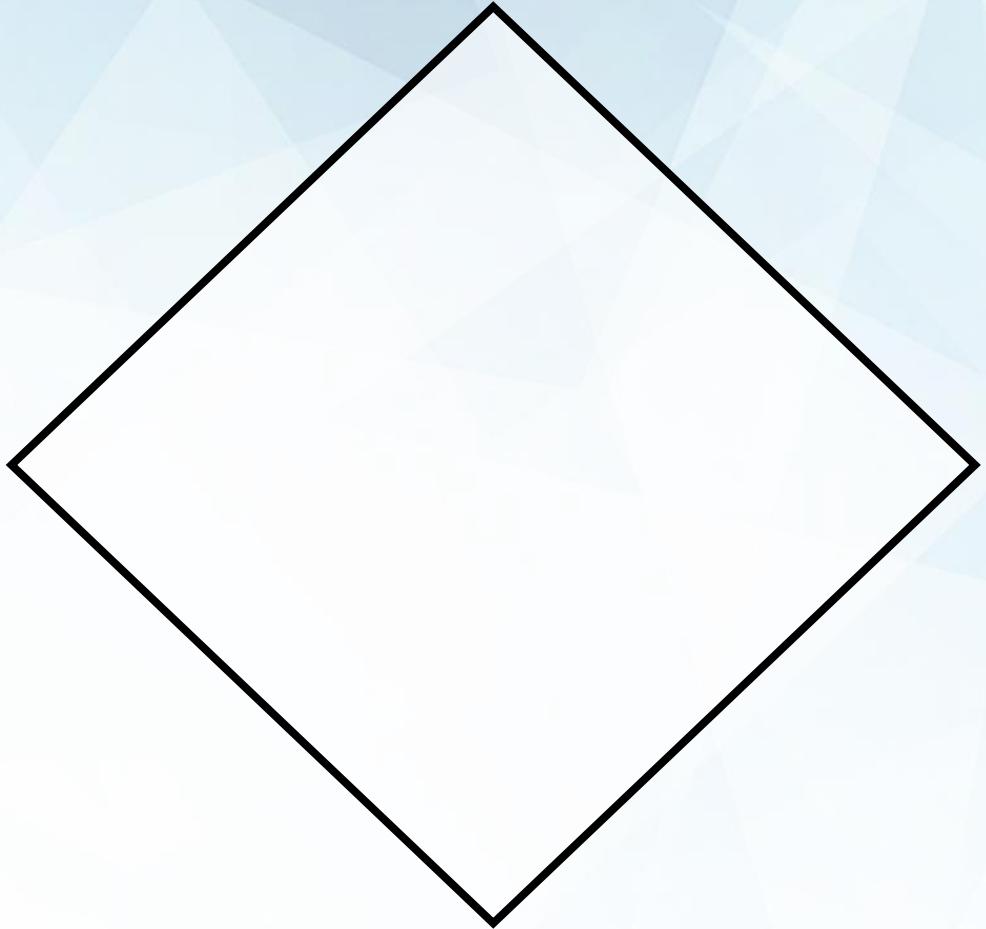


Refine&Mine

Sequence of discretization **refinements** and mining patterns within the **discretization**.

REFINE&MINE: In a Nutshell (2)

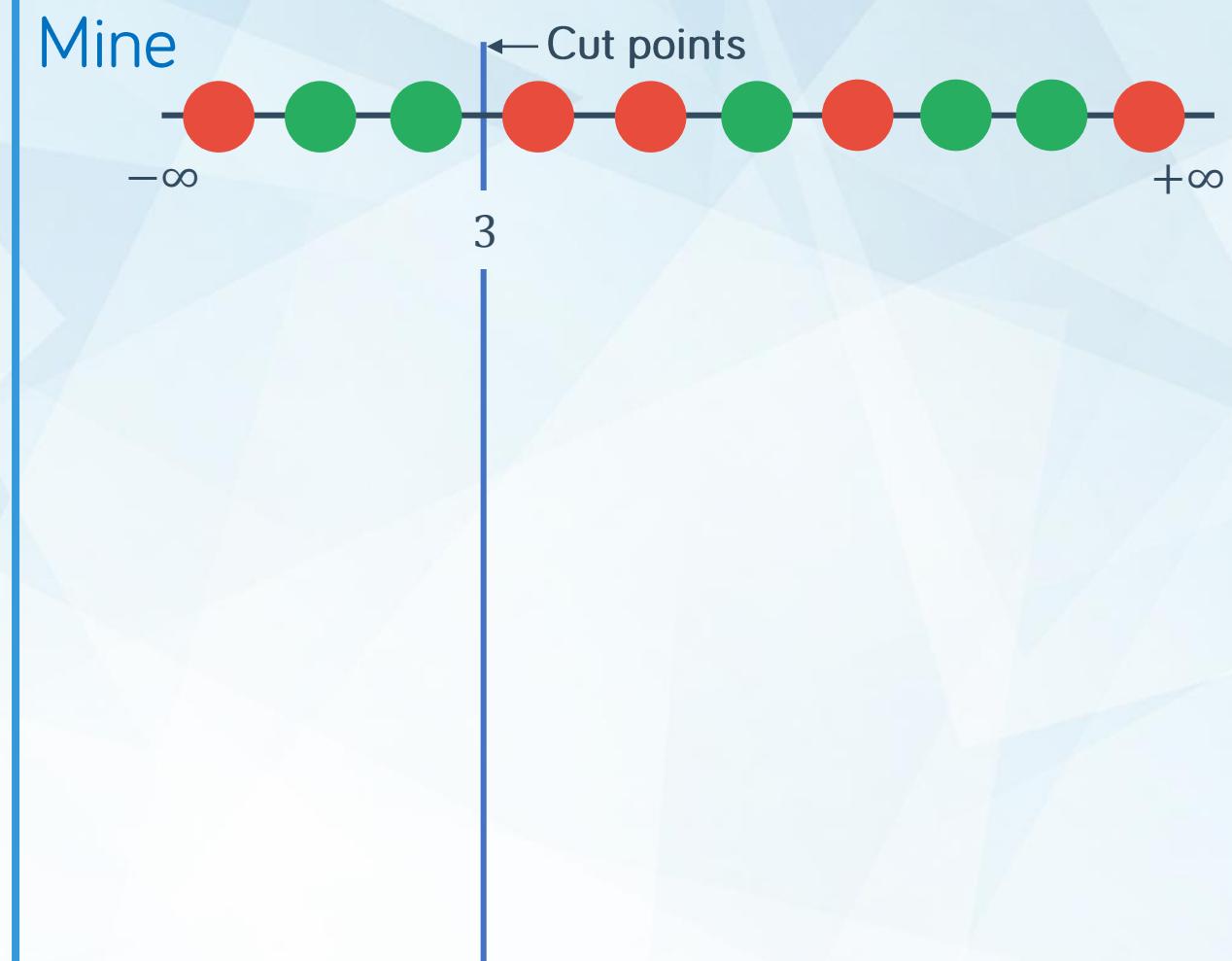
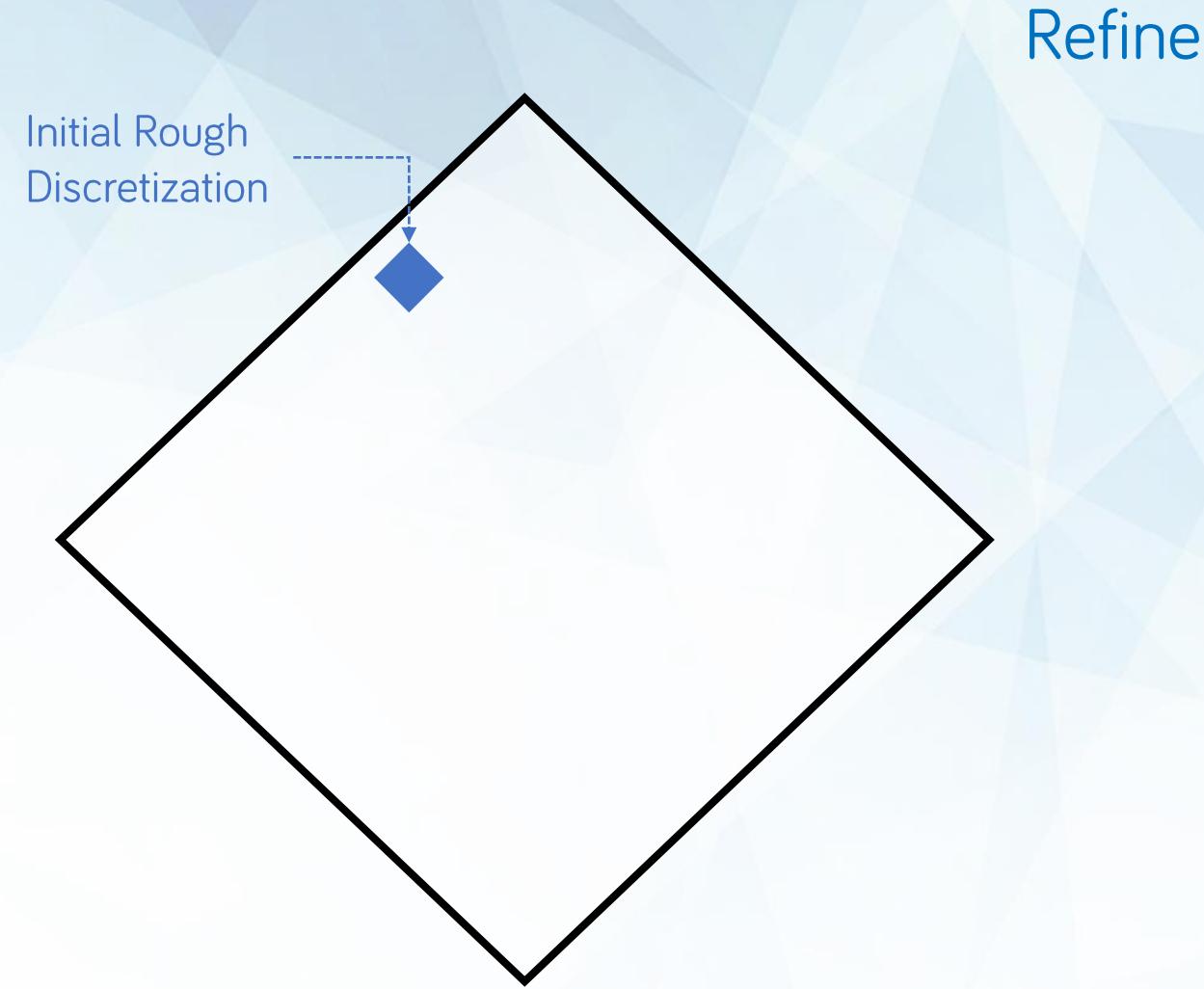
Refine



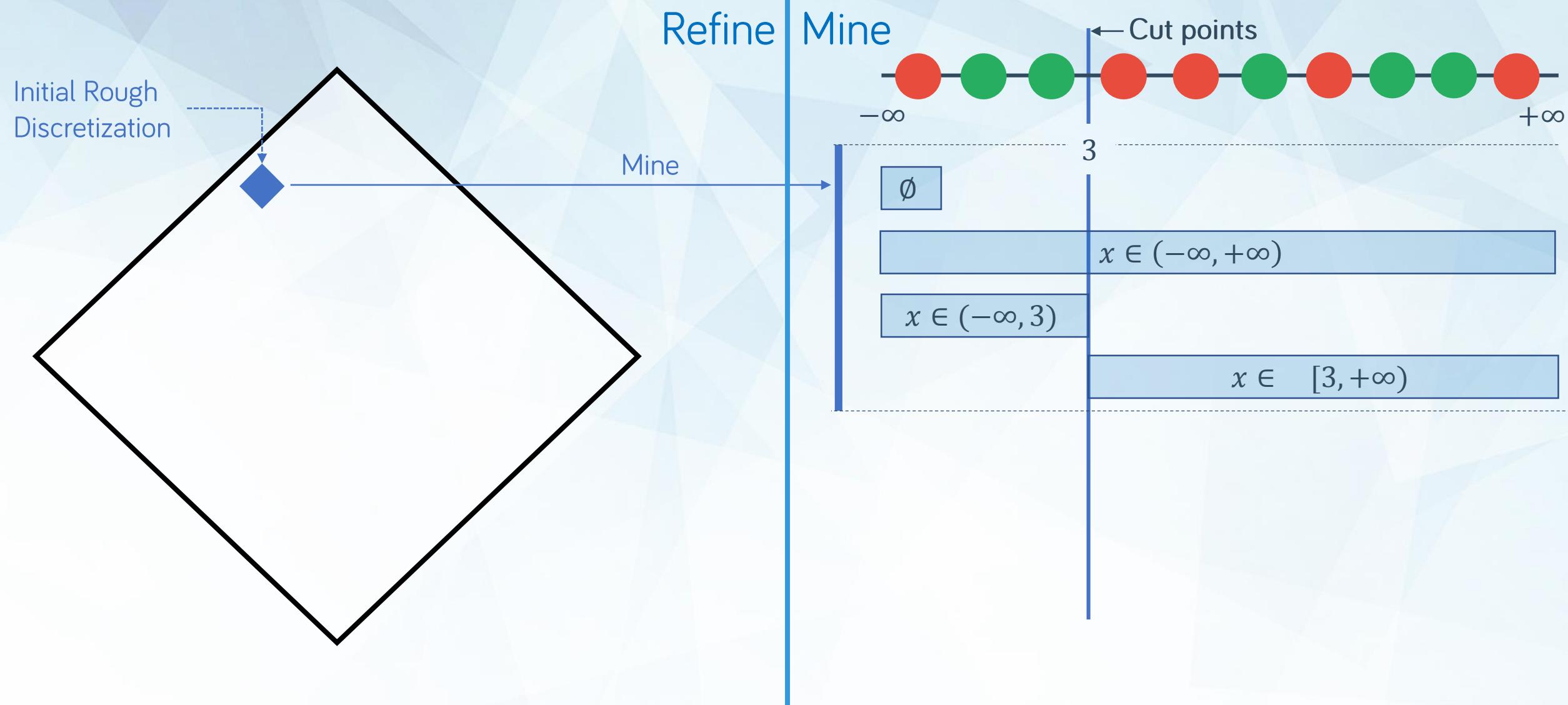
Mine



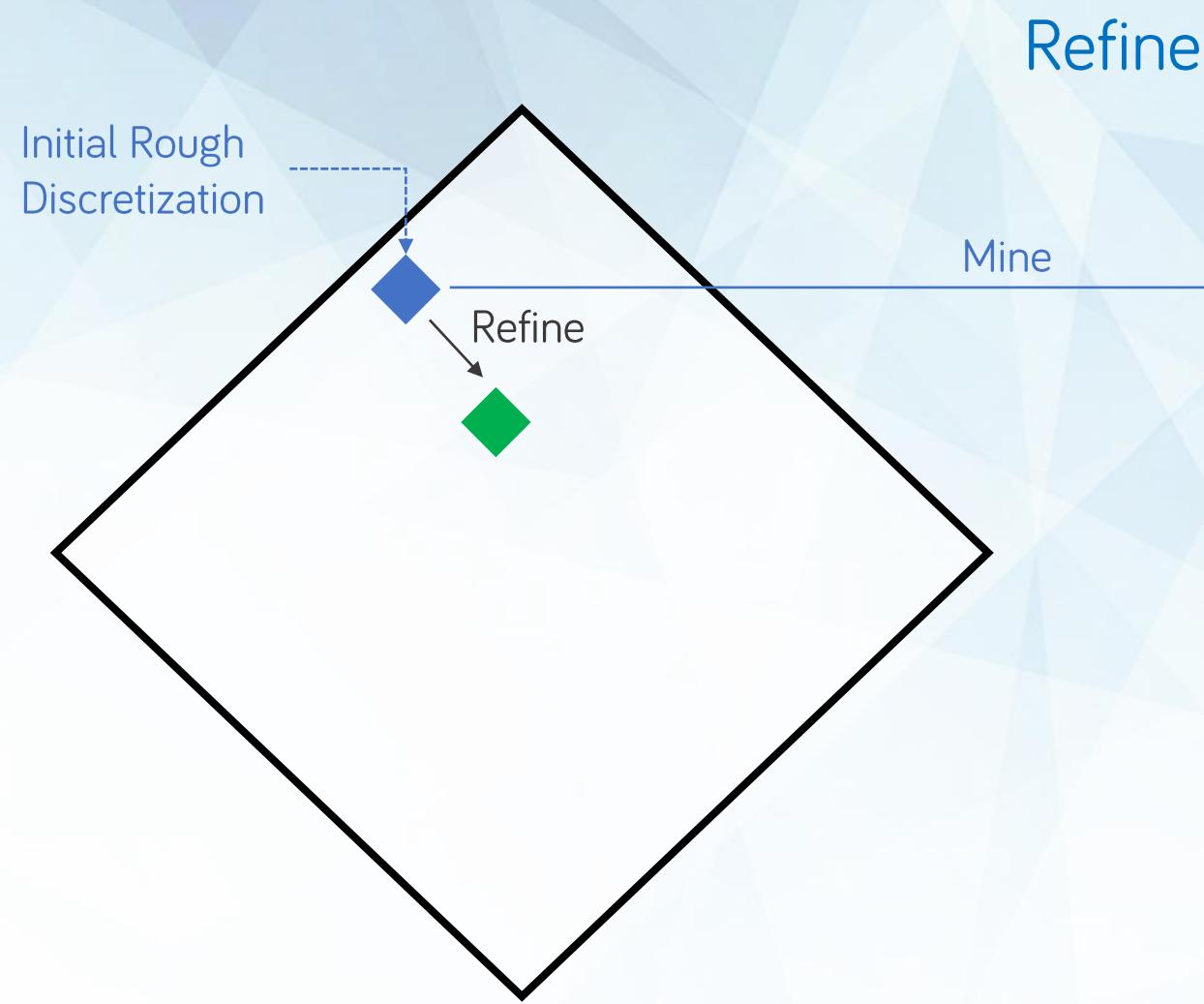
REFINE&MINE: In a Nutshell (2)



REFINE&MINE: In a Nutshell (2)

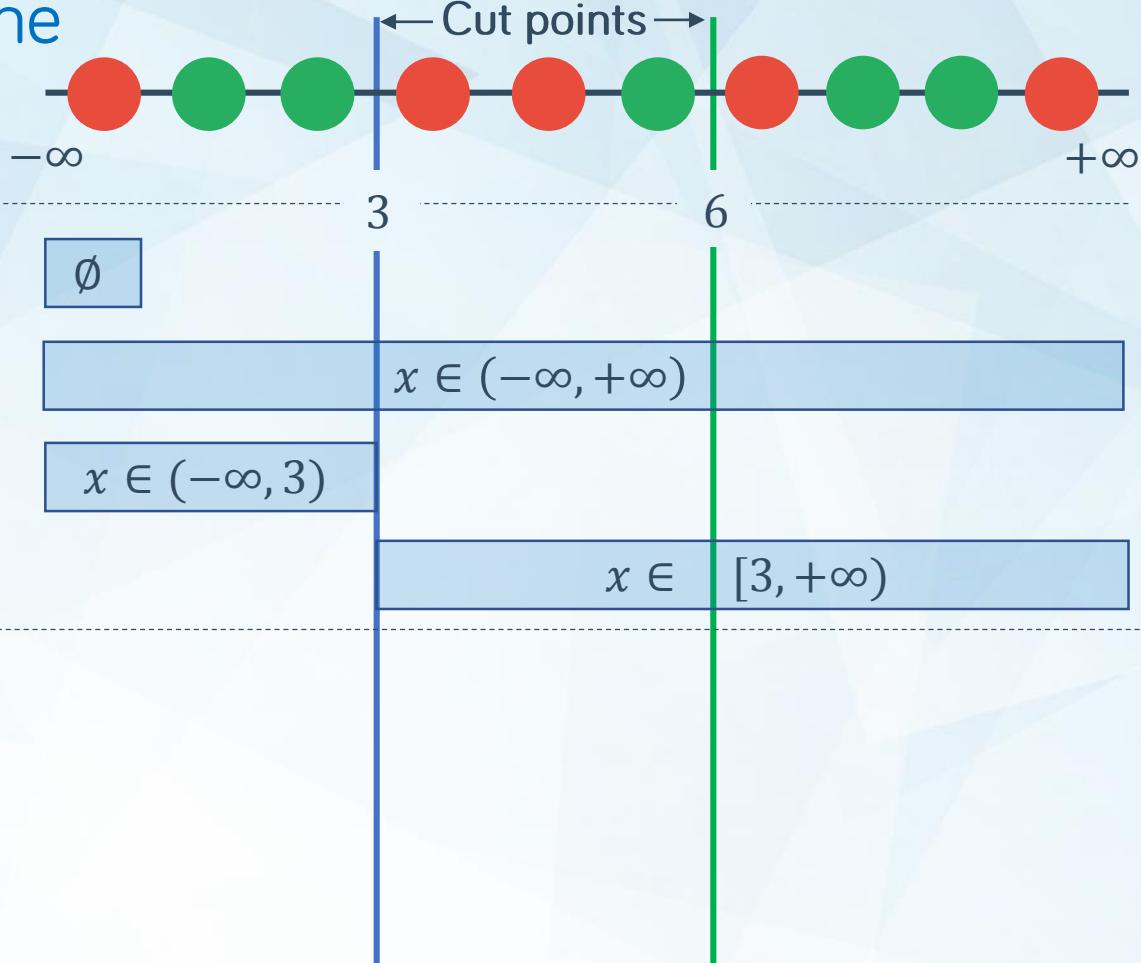


REFINE&MINE: In a Nutshell (2)

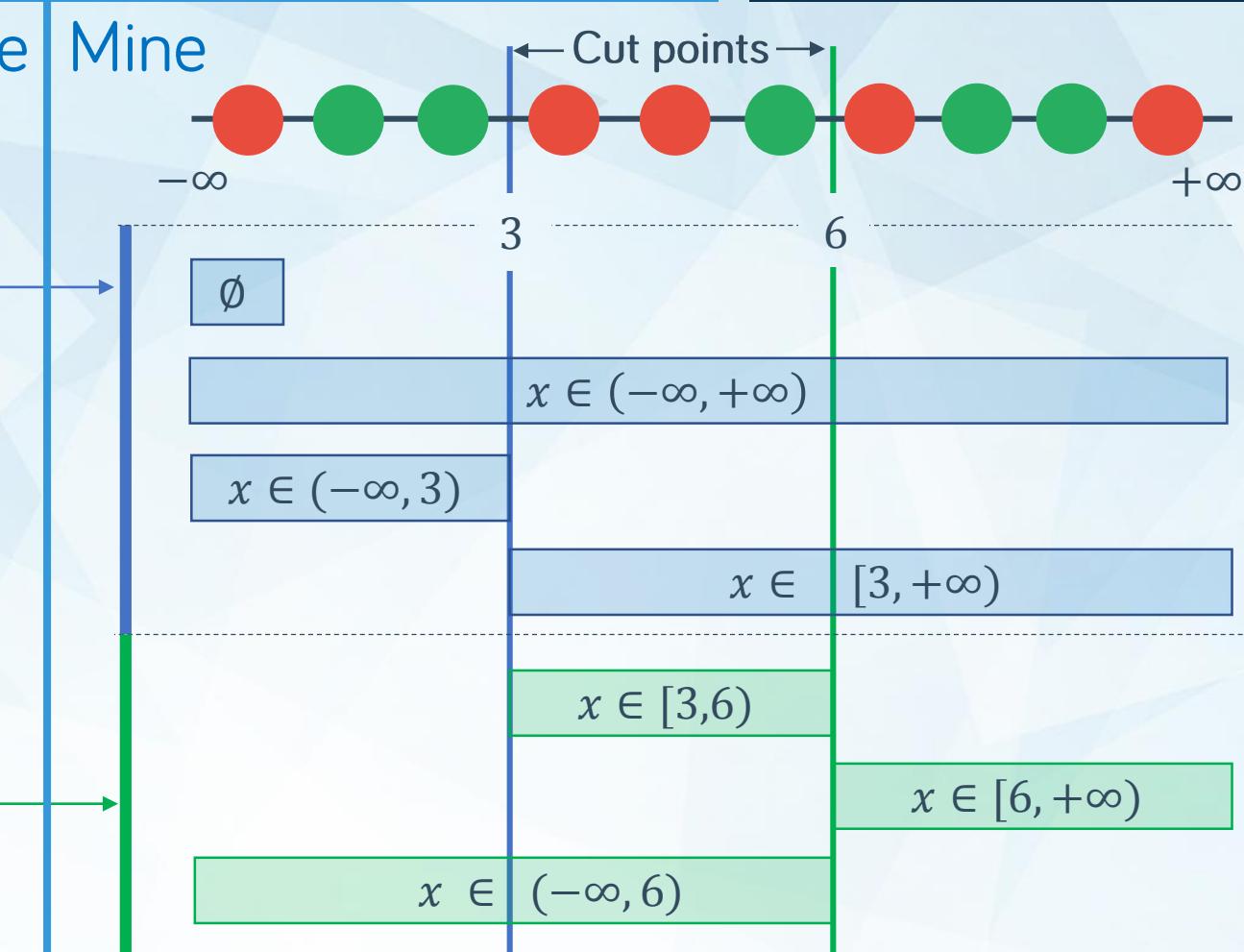
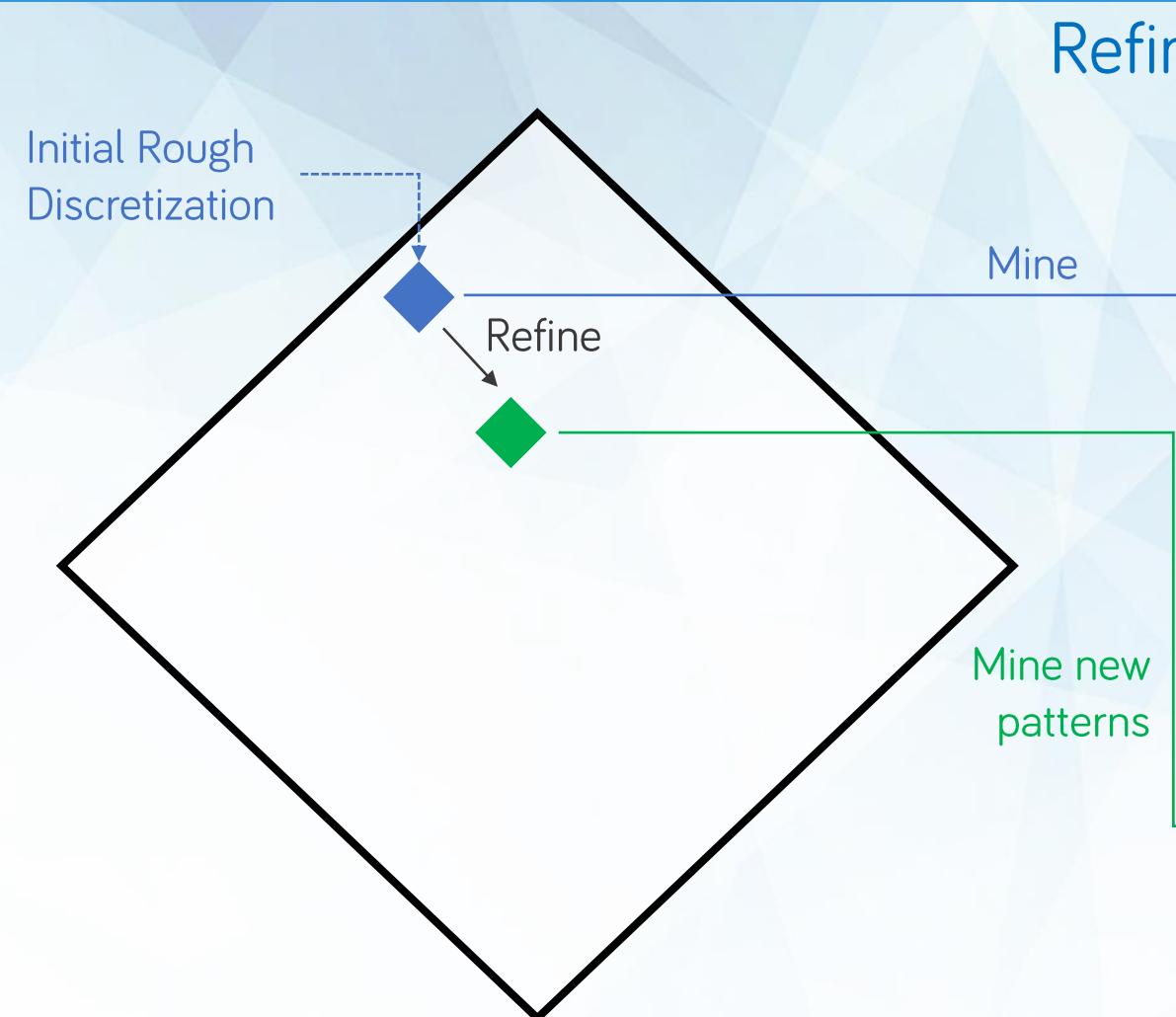


Refine

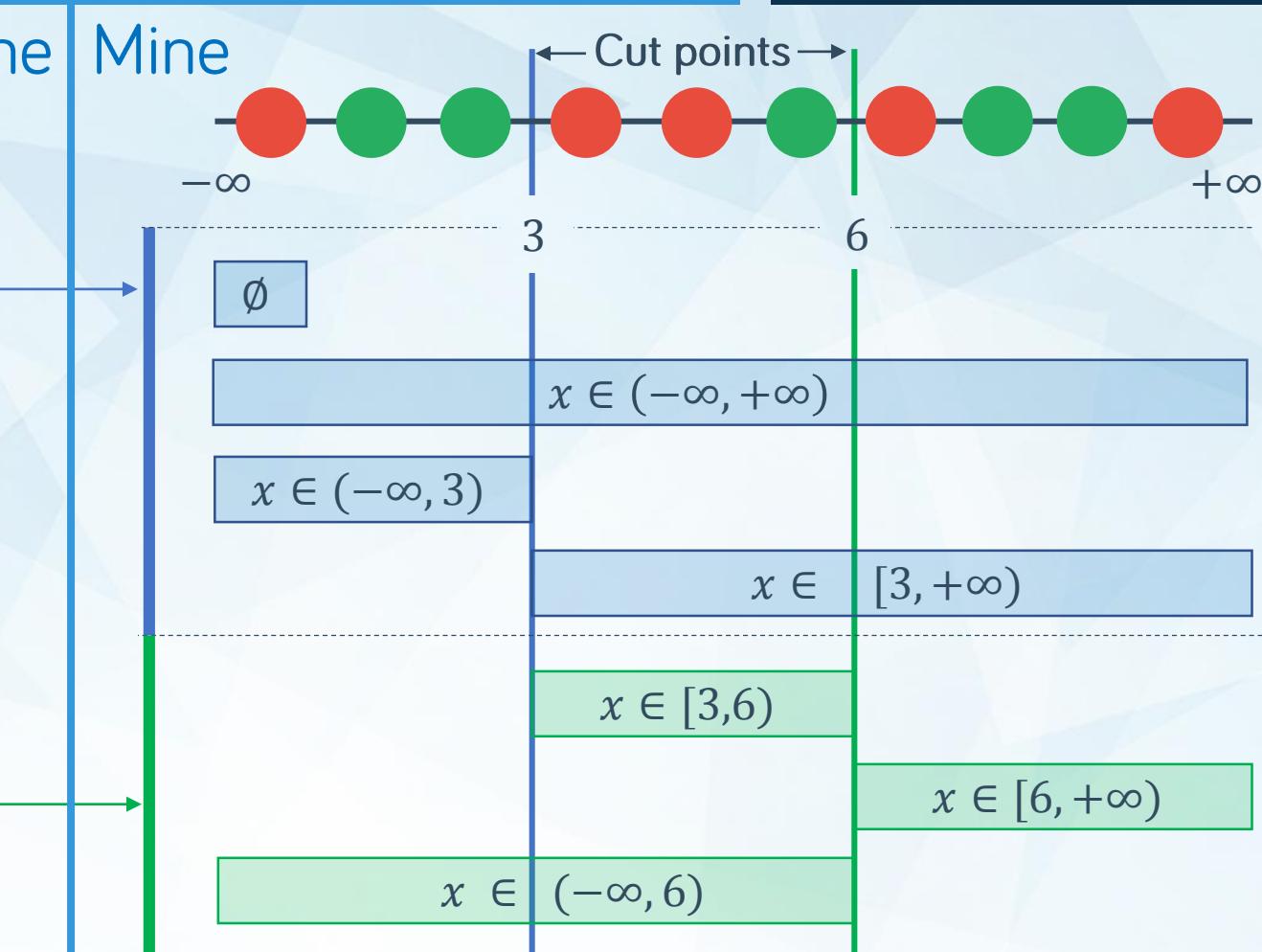
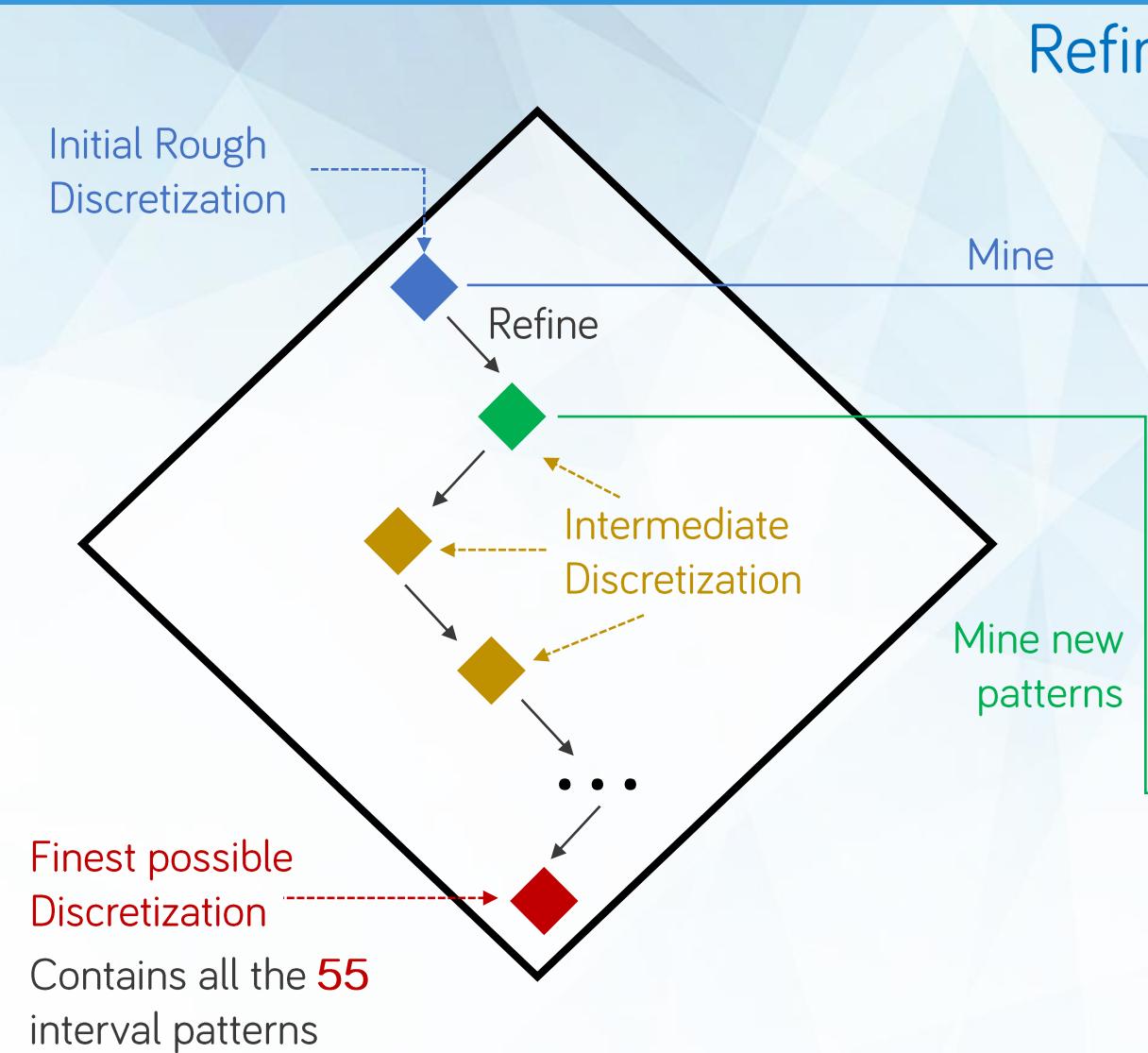
Mine



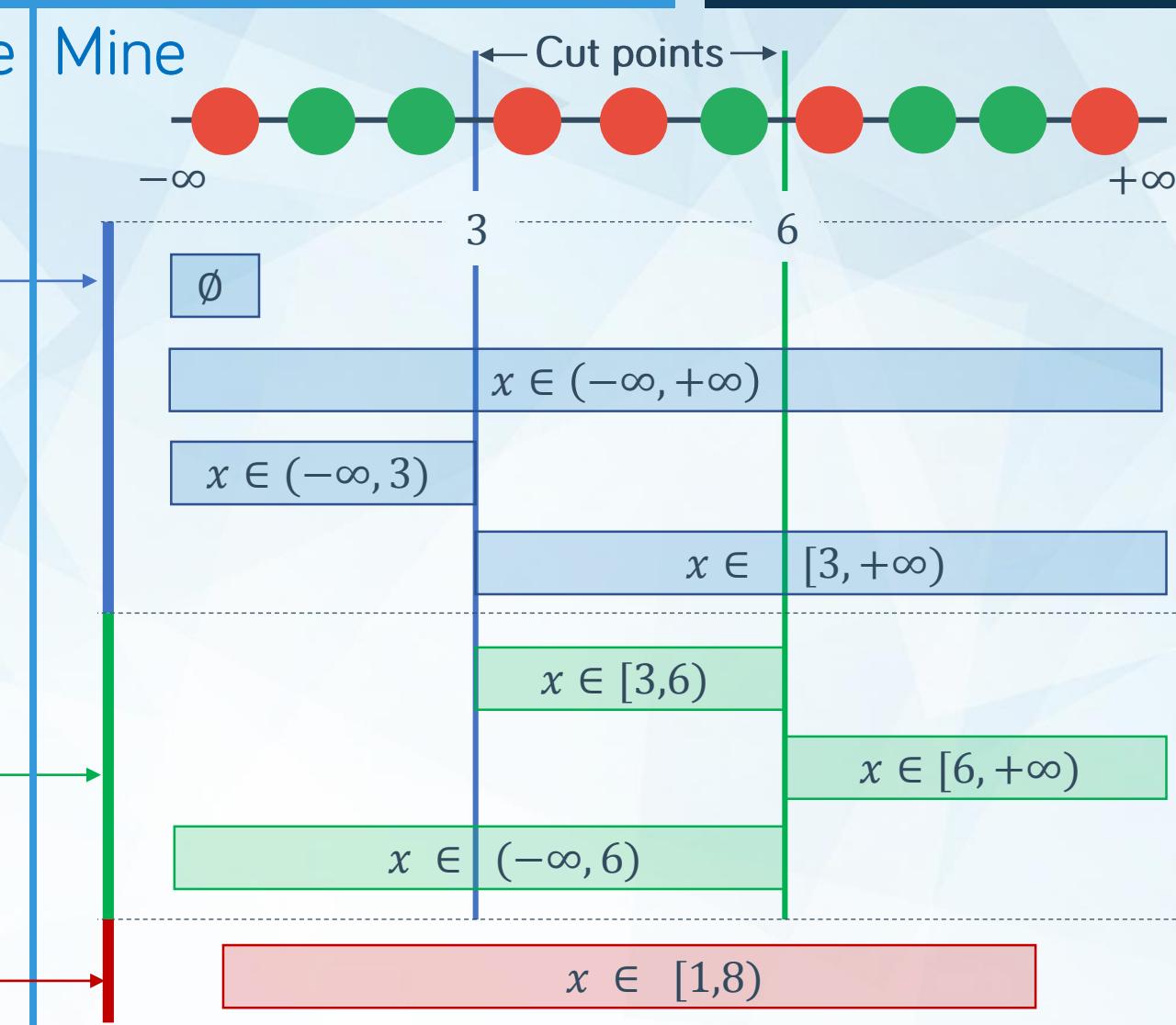
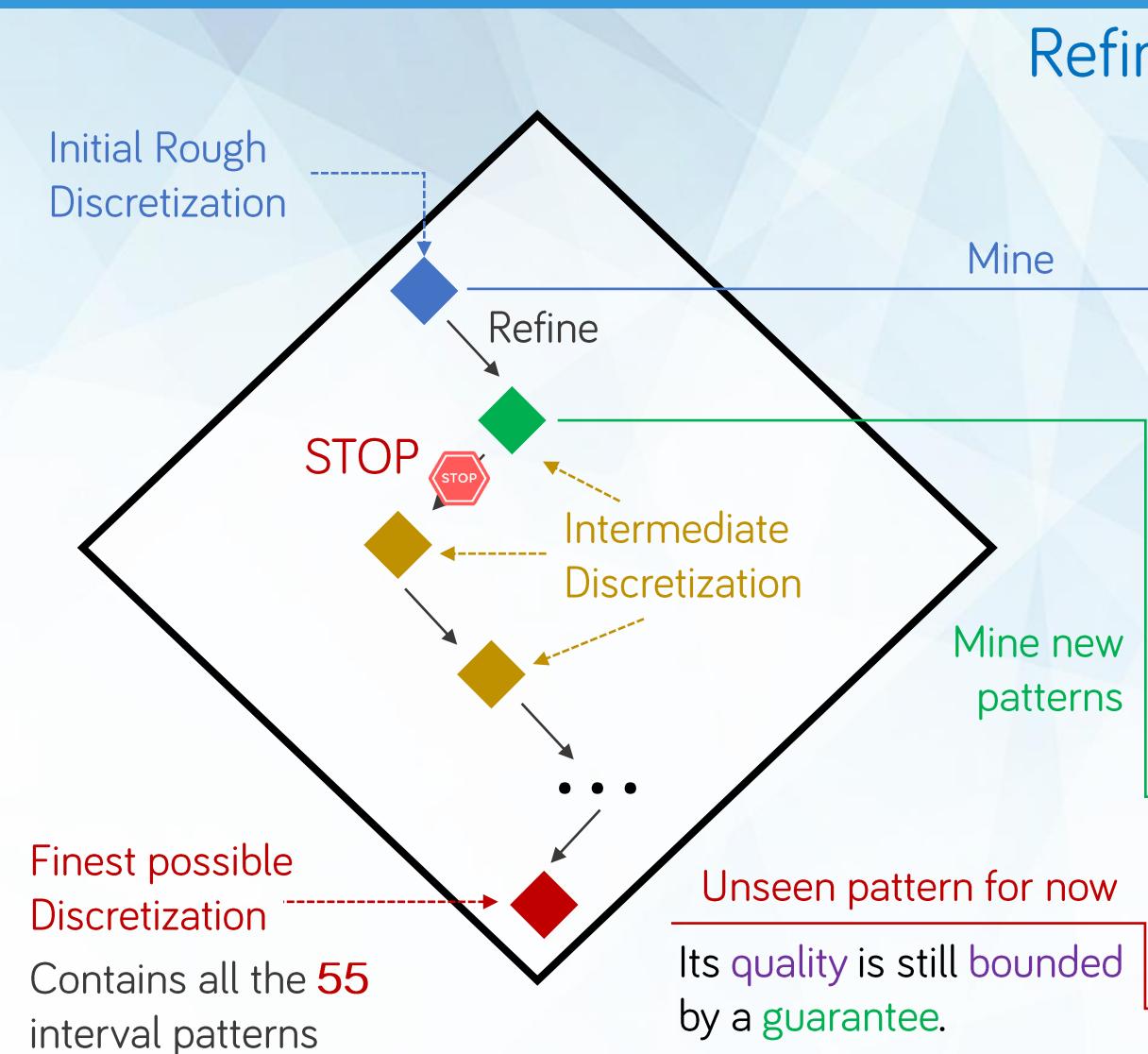
REFINE&MINE: In a Nutshell (2)



REFINE&MINE: In a Nutshell (2)



REFINE&MINE: In a Nutshell (2)



Guarantees provided by Refine&Mine

e.g. Informedness



$$\text{Accuracy} = \phi(\text{best pattern in the dataset}) - \phi(\text{best found pattern})$$

e.g. Informedness

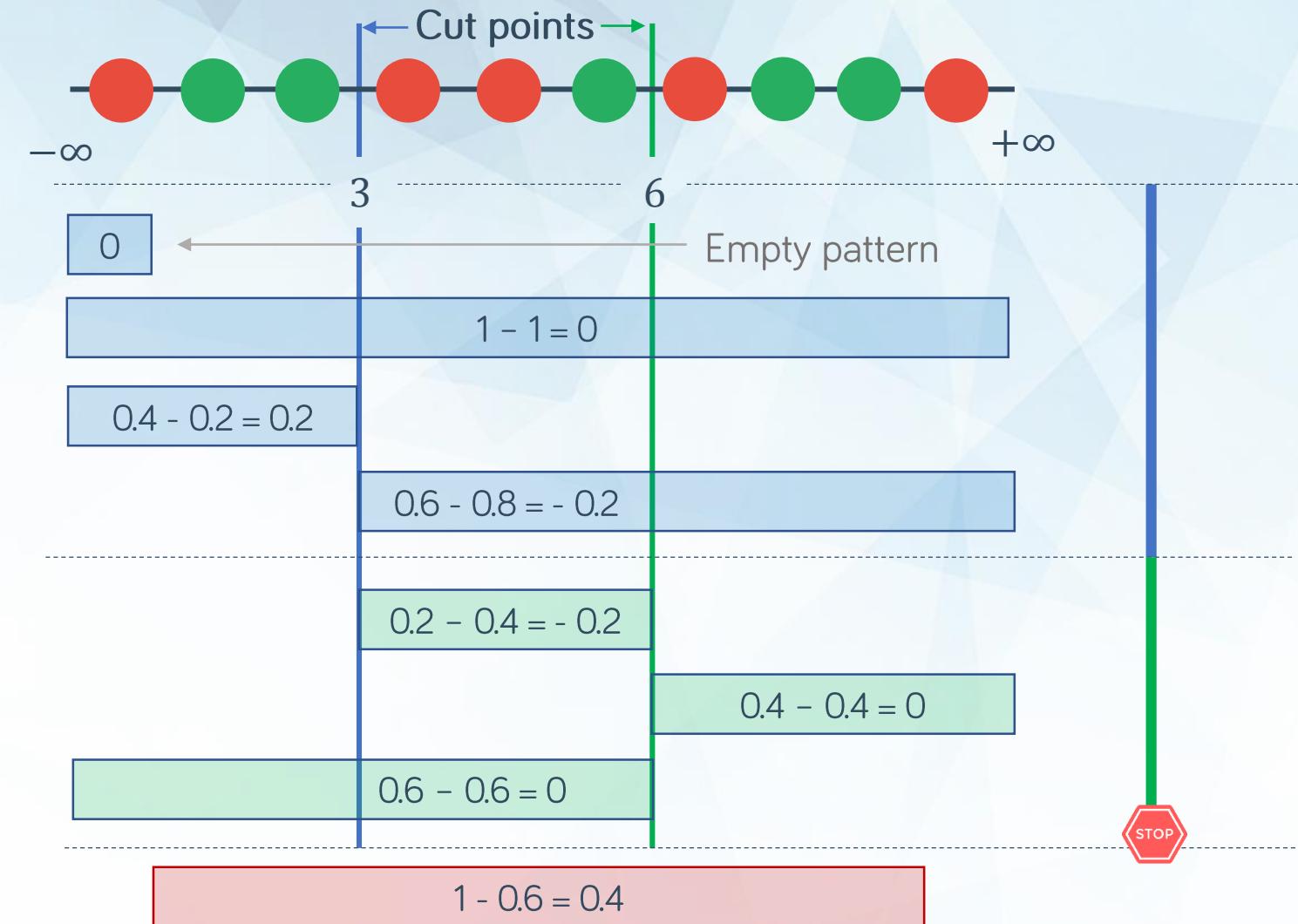


$$\text{Accuracy} = \phi(\text{best pattern in the dataset}) - \phi(\text{best found pattern})$$

$$\text{Accuracy} \leq \text{Accuracy bound}$$

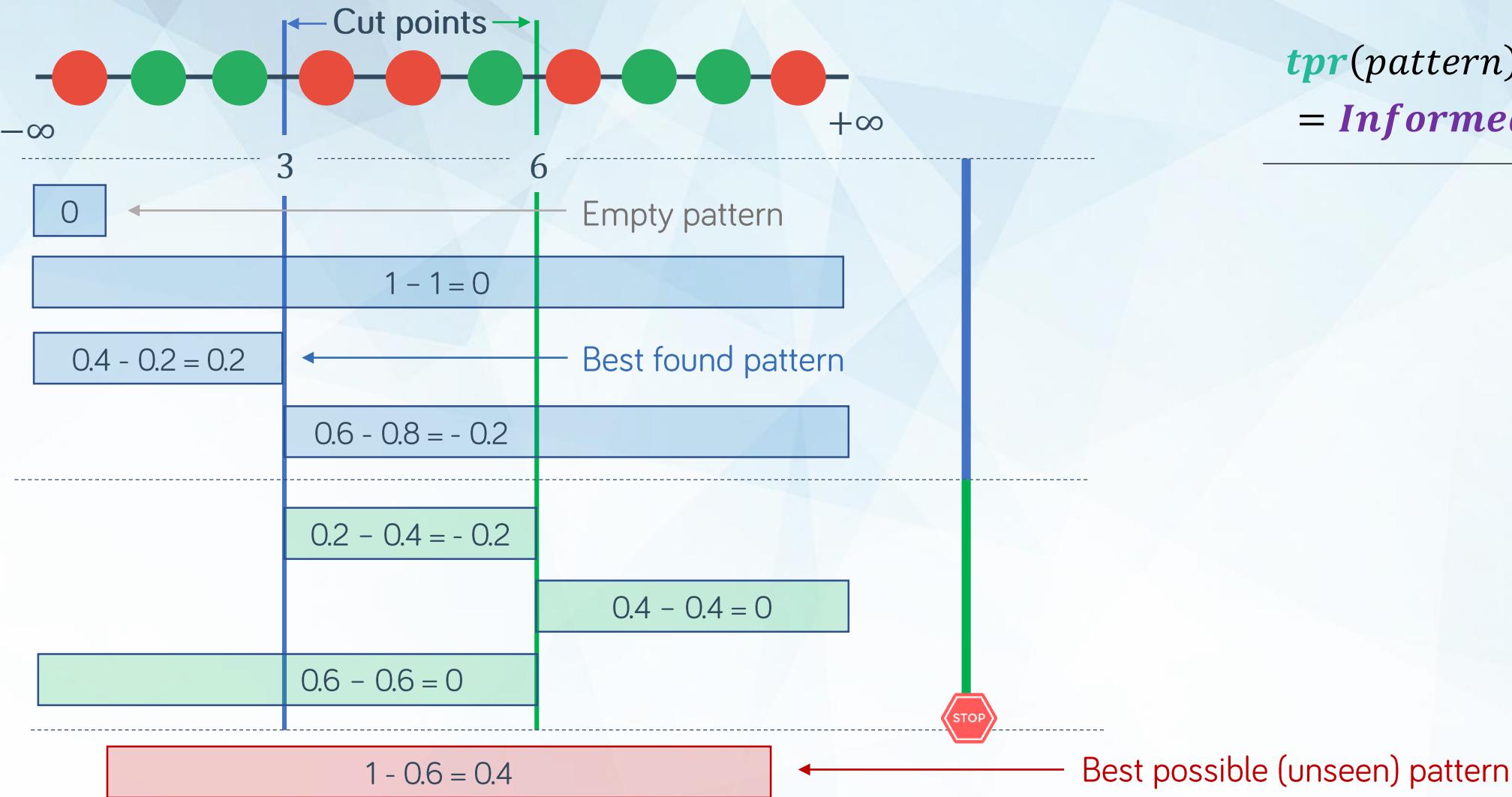


Refine&Mine computes this bound
using only already found patterns.



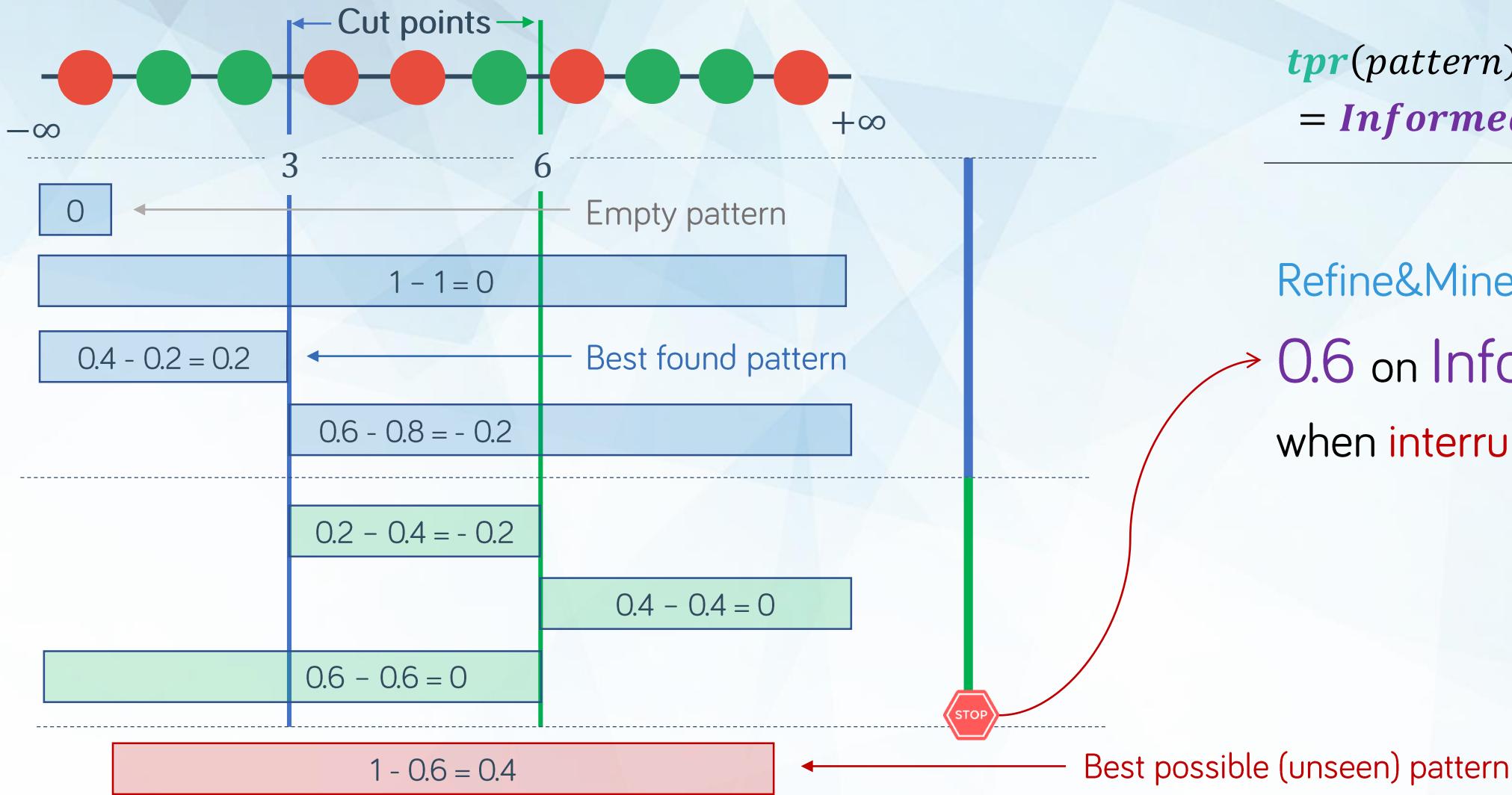
$tpr(\text{pattern}) - fpr(\text{pattern})$
 $= \text{Informedness}(\text{pattern})$

Accuracy Guarantee – Introduction



$$tpr(\text{pattern}) - fpr(\text{pattern}) = \text{Informedness}(\text{pattern})$$

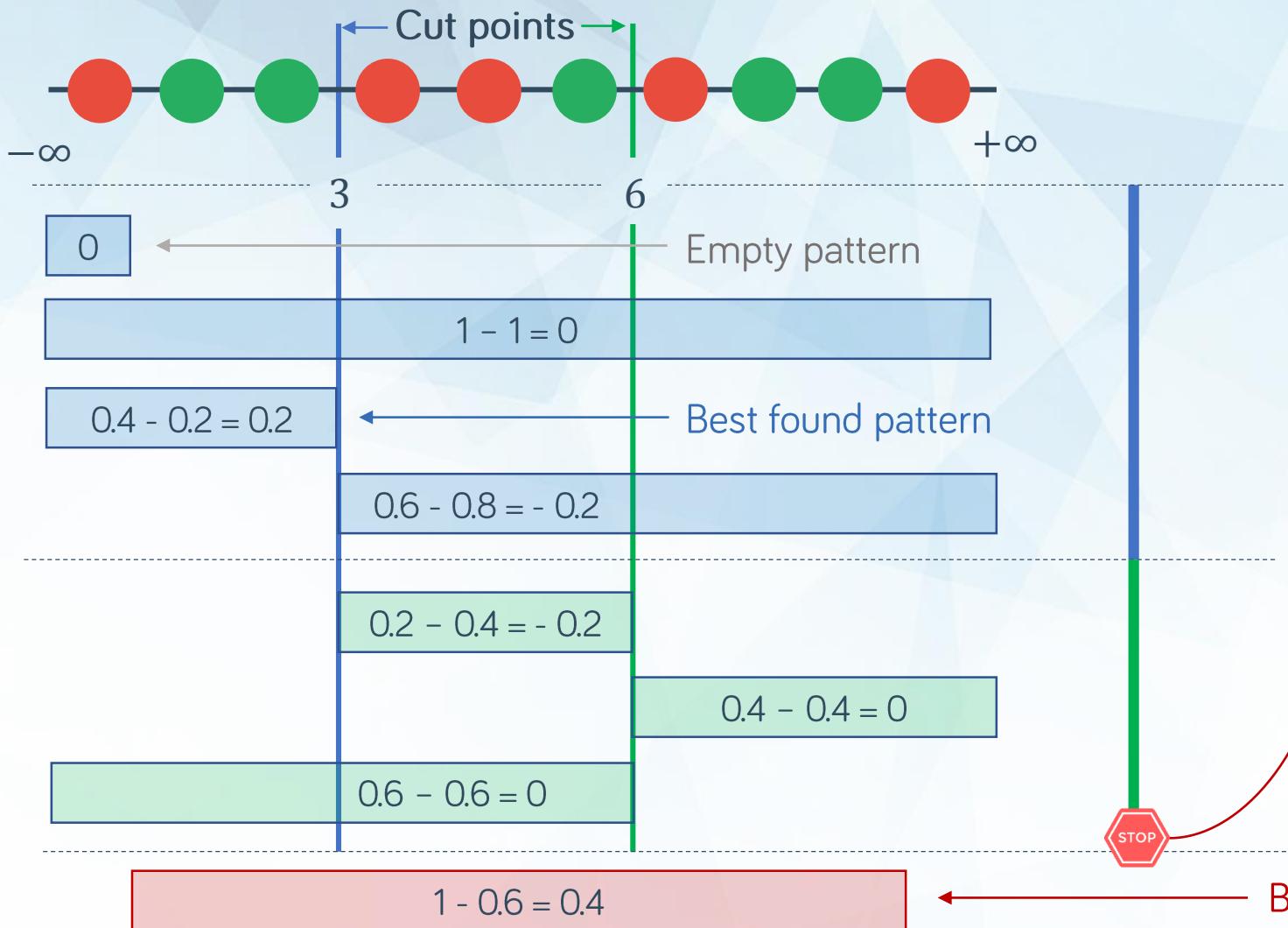
Accuracy Guarantee – Introduction



$$tpr(\text{pattern}) - fpr(\text{pattern}) = \text{Informedness}(\text{pattern})$$

Refine&Mine gives a bound of 0.6 on Informedness when interrupted.

Accuracy Guarantee – Introduction

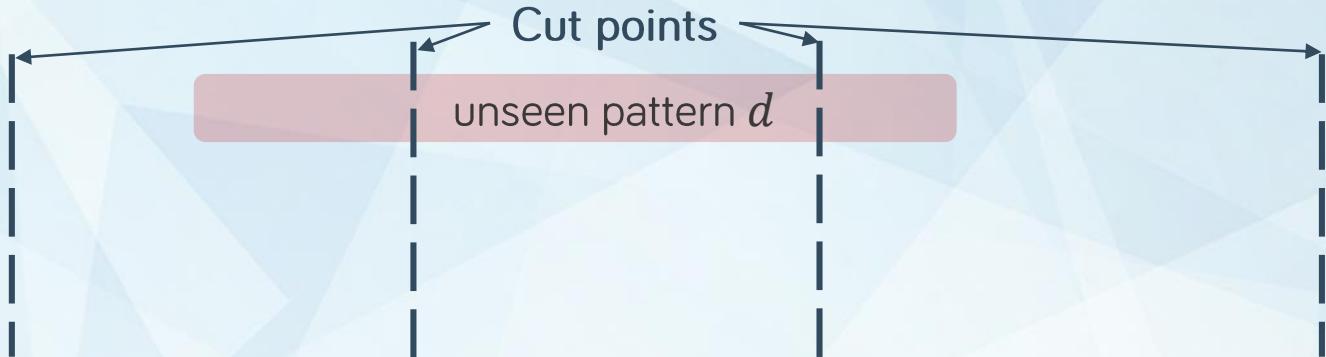


$$tpr(\text{pattern}) - fpr(\text{pattern}) = \text{Informedness}(\text{pattern})$$

Refine&Mine gives a bound of 0.6 on Informedness when interrupted.



How?

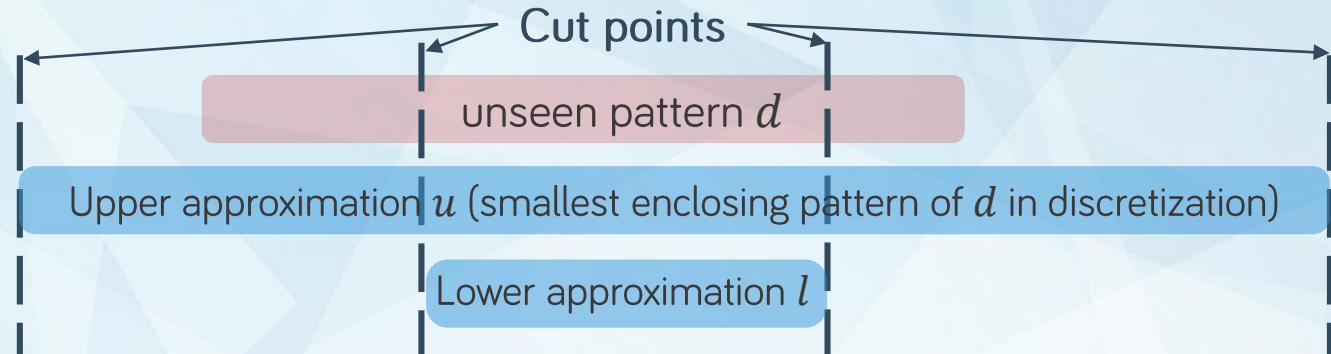


Z. Pawlak. [Rough sets](#). In International Journal of Parallel Programming 1982.

A. Tarski. [A lattice-theoretical fixpoint theorem and its applications](#). In Pacific journal of Mathematics 1955.



Any **unseen interval pattern**
can be encapsulated between
two found patterns.



Z. Pawlak. **Rough sets**. In International Journal of Parallel Programming 1982.

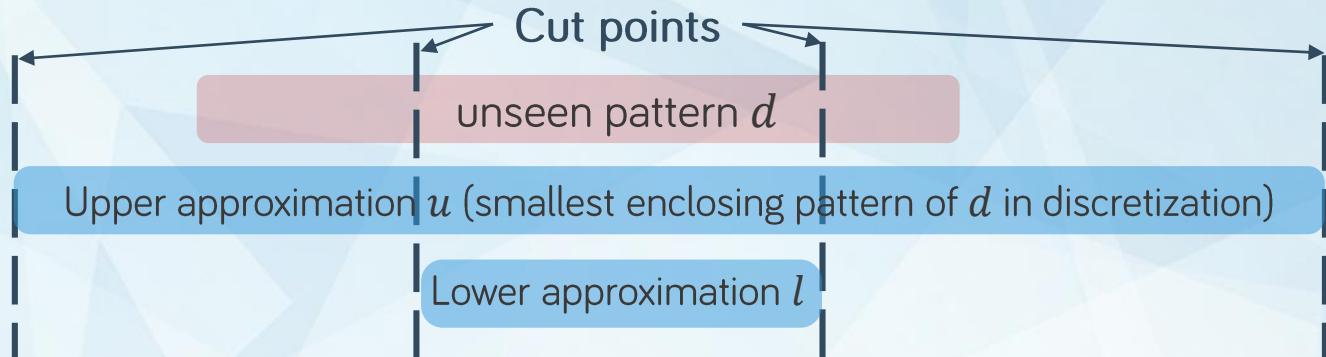
A. Tarski. **A lattice-theoretical fixpoint theorem and its applications**. In Pacific journal of Mathematics 1955.



Any **unseen interval pattern** can be encapsulated between **two found patterns**.



tpr and **fpr** measures are monotonous.



- $\text{tpr}(d) \leq \text{tpr}(u)$
- $\text{fpr}(d) \geq \text{fpr}(l)$ $\Rightarrow \text{informedness}(d) \leq \text{tpr}(u) - \text{fpr}(l)$



Z. Pawlak. **Rough sets**. In International Journal of Parallel Programming 1982.

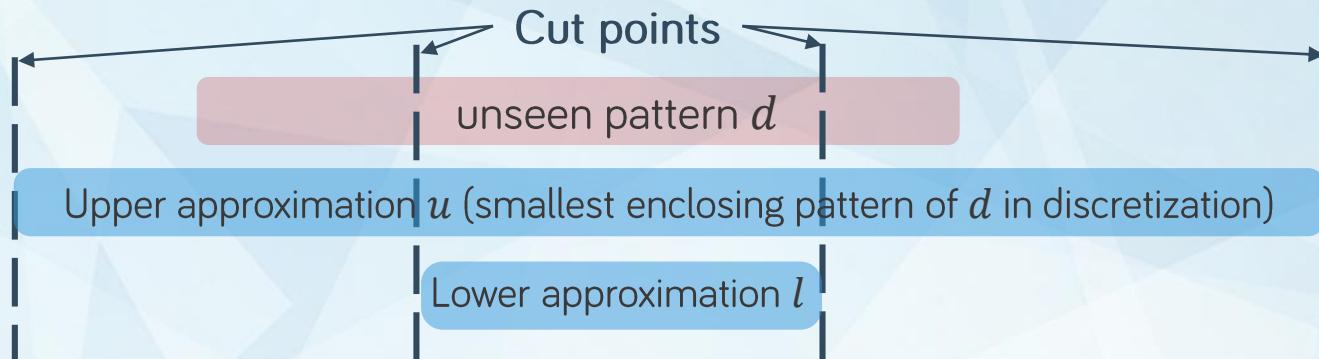
A. Tarski. **A lattice-theoretical fixpoint theorem and its applications**. In Pacific journal of Mathematics 1955.



Any **unseen interval pattern**
can be encapsulated between
two found patterns.



tpr and **fpr** measures are
monotonous.



- $\text{tpr}(d) \leq \text{tpr}(u)$
- $\text{fpr}(d) \geq \text{fpr}(l)$ $\Rightarrow \text{informedness}(d) \leq \text{tpr}(u) - \text{fpr}(l)$

Now, we gave a bound on the **quality** of **an unseen pattern**
using its **two already found patterns**.



Z. Pawlak. **Rough sets**. In International Journal of Parallel Programming 1982.

A. Tarski. **A lattice-theoretical fixpoint theorem and its applications**. In Pacific journal of Mathematics 1955.



Q. So, do we need to look for pair of patterns in the discretization search space to give bound the quality measure of all unseen patterns?



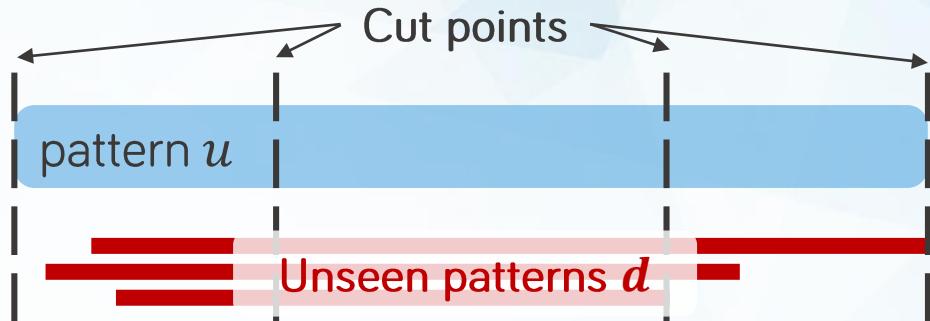
Q. So, do we need to look for pair of patterns in the discretization search space to give bound the quality measure of all unseen patterns?

A. No, we use the notion of core



Q. So, do we need to look for pair of patterns in the discretization search space to give bound the quality measure of all unseen patterns?

A. No, we use the notion of core

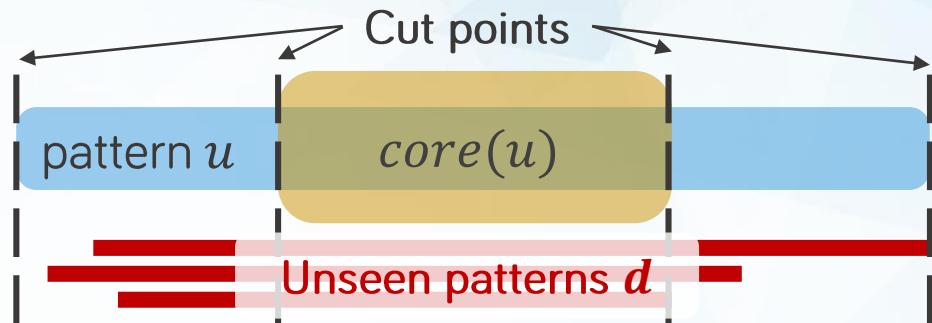


- Take all unseen patterns which upper approximation is u .



Q. So, do we need to look for pair of patterns in the discretization search space to give bound the quality measure of all unseen patterns?

A. No, we use the notion of core



- Take all unseen patterns which upper approximation is u .
- The core of a found pattern u the largest possible lower approximation of all these unseen patterns.

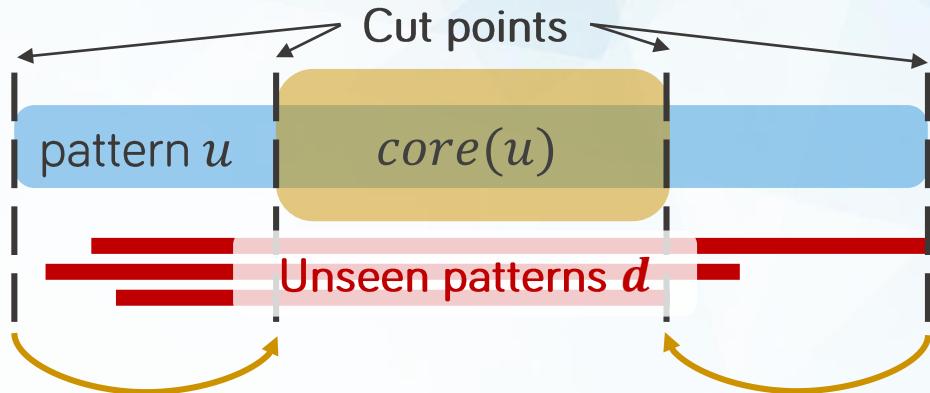
Core(found pattern) \subseteq some unseen patterns \subseteq found pattern

Upper approximation



Q. So, do we need to look for pair of patterns in the discretization search space to give bound the quality measure of all unseen patterns?

A. No, we use the notion of core

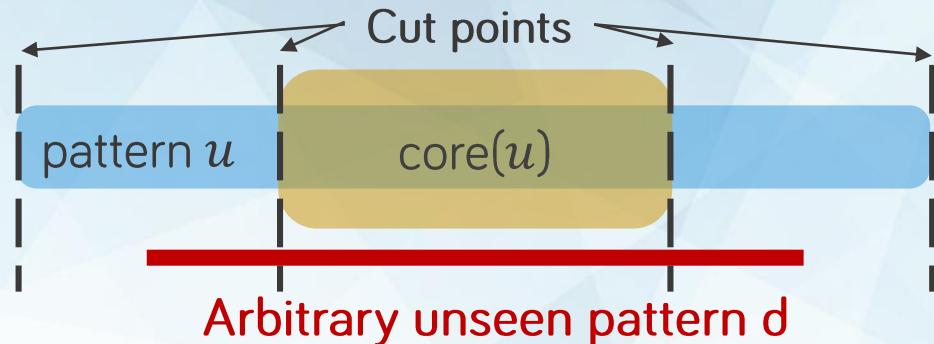


- Take all unseen patterns which upper approximation is u .
- The core of a found pattern u the largest possible lower approximation of all these unseen patterns.

Core(found pattern) \subseteq some unseen patterns \subseteq found pattern



The core computation complexity is $O(d)$ with d is the number of dimensions.

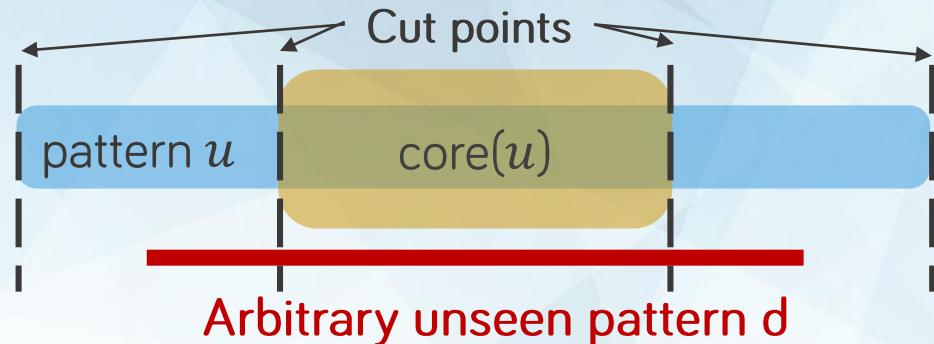


- Core(found pattern) \subseteq some unseen patterns \subseteq found pattern
- Upper approximation

Accuracy =

$$\sup_{d \in all} \phi(d)$$

$$- \sup_{u \in found} \phi(u)$$



- $\text{Core}(\text{found pattern}) \subseteq \text{some unseen patterns} \subseteq \text{found pattern}$

$$\begin{aligned} \text{Accuracy} &= \sup_{d \in \text{all}} \phi(d) - \sup_{u \in \text{found}} \phi(u) \\ &\leq \\ \text{Accuracy} &\leq \sup_{u \in \text{found}} \phi\left(\text{tpr}(u), \text{fpr}(\text{core}(u))\right) - \sup_{u \in \text{found}} \phi(u) \end{aligned}$$

Specificity = $\overrightarrow{distance}$ (All found patterns, All pattern in the dataset)

Specificity = $\overrightarrow{distance}$ (All found patterns, All pattern in the dataset)

Specificity \leq Specificity bound



Same aforementioned notions are used to define the **bound** on specificity

Refine&Mine computes this bound using **only** already found patterns.



Overview and Contributions



Refine&Mine in a nutshell



Guarantees computation



Building blocks of Refine&Mine



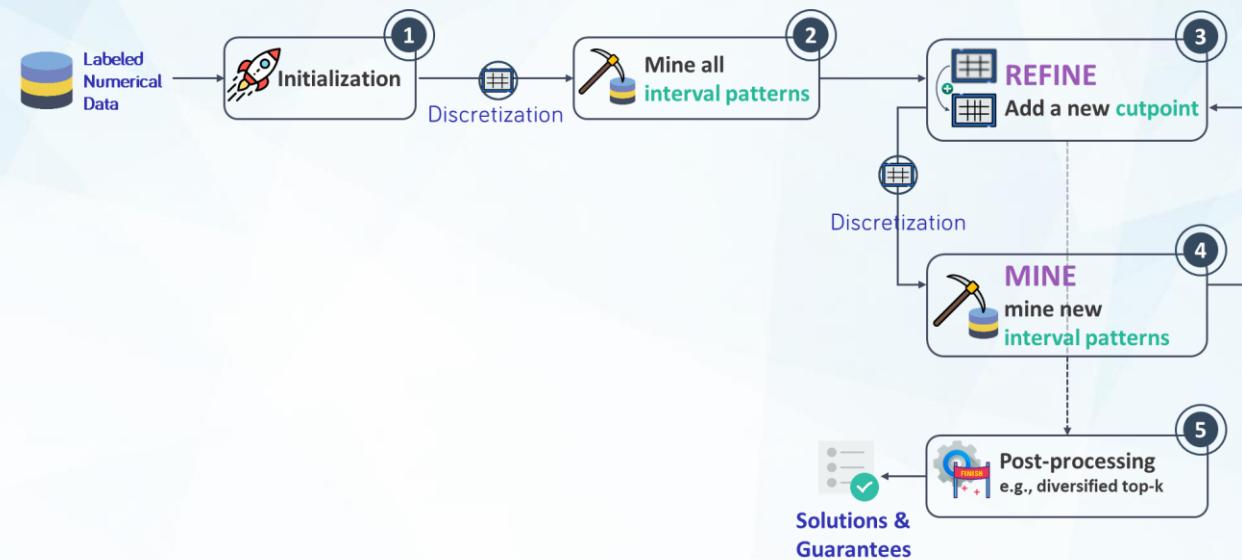
Wrap-up and concluding remarks



OUTLINE

Refine&Mine

Step-by-step





Labeled
Numerical
Data

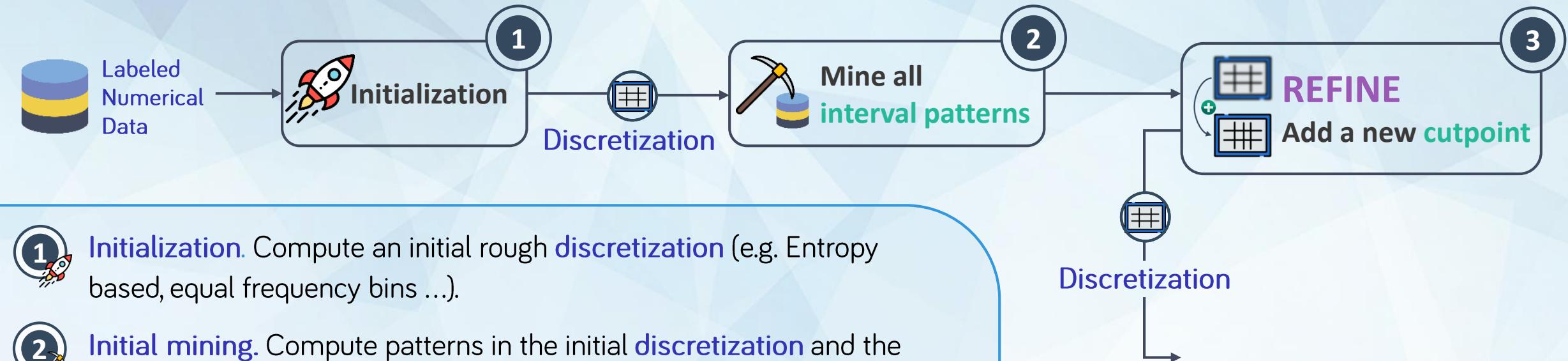


- 1** **Initialization.** Compute an initial rough **discretization** (e.g. Entropy based, equal frequency bins ...).



- 1 **Initialization.** Compute an initial rough **discretization** (e.g. Entropy based, equal frequency bins ...).
- 2 **Initial mining.** Compute patterns in the initial **discretization** and the first guarantees (e.g. use optimizations like closure operator) .

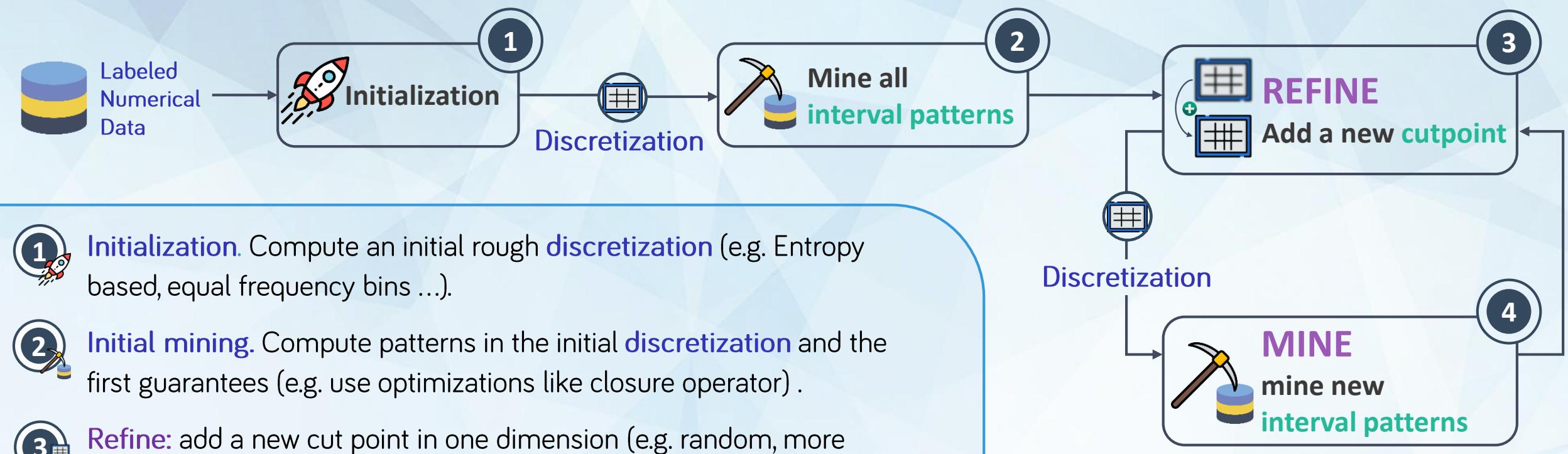
REFINE&MINE: Step-by-step



-  **Initialization.** Compute an initial rough **discretization** (e.g. Entropy based, equal frequency bins ...).
 -  **Initial mining.** Compute patterns in the initial **discretization** and the first guarantees (e.g. use optimizations like closure operator) .
 -  **Refine:** add a new cut point in one dimension (e.g. random, more sophisticated strategy to reduce quickly the accuracy bound).

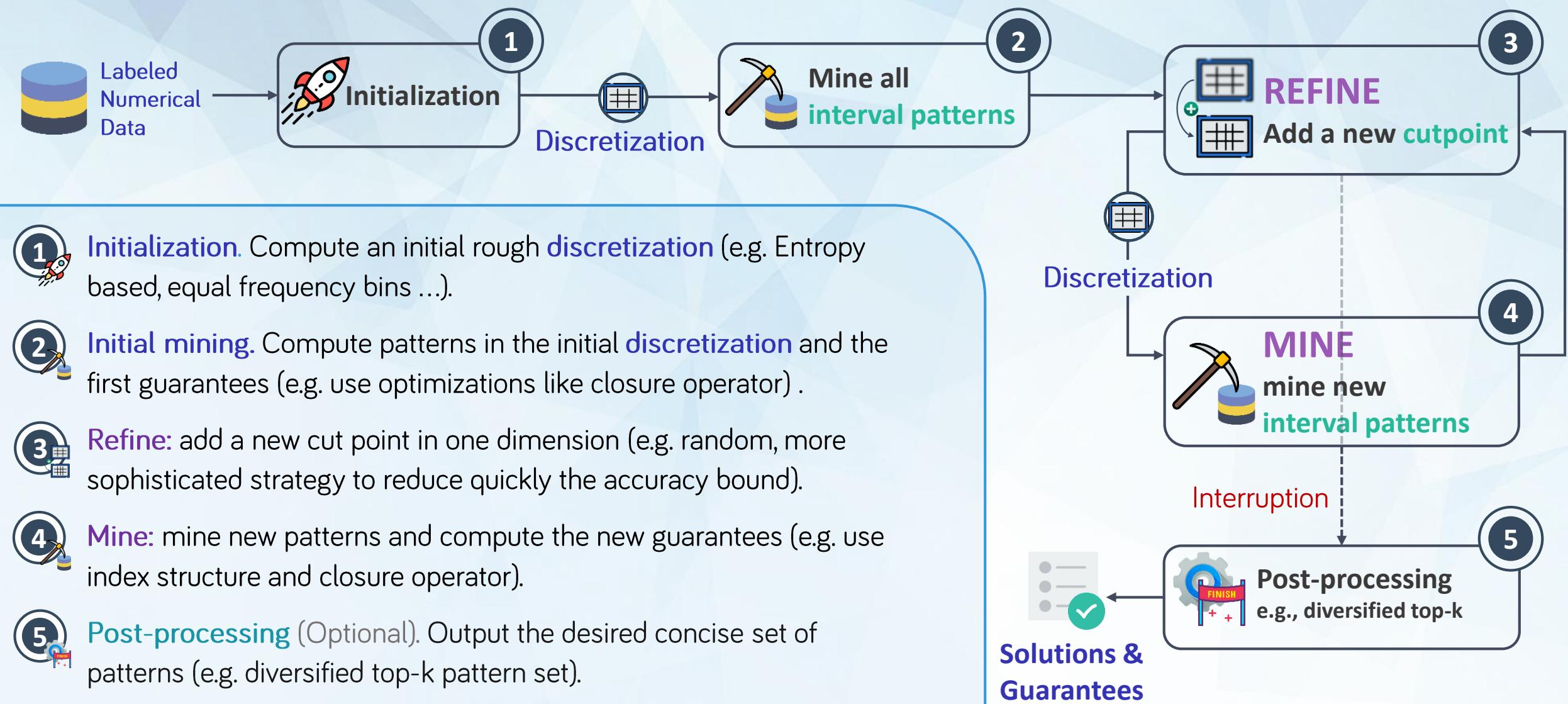


REFINE&MINE: Step-by-step



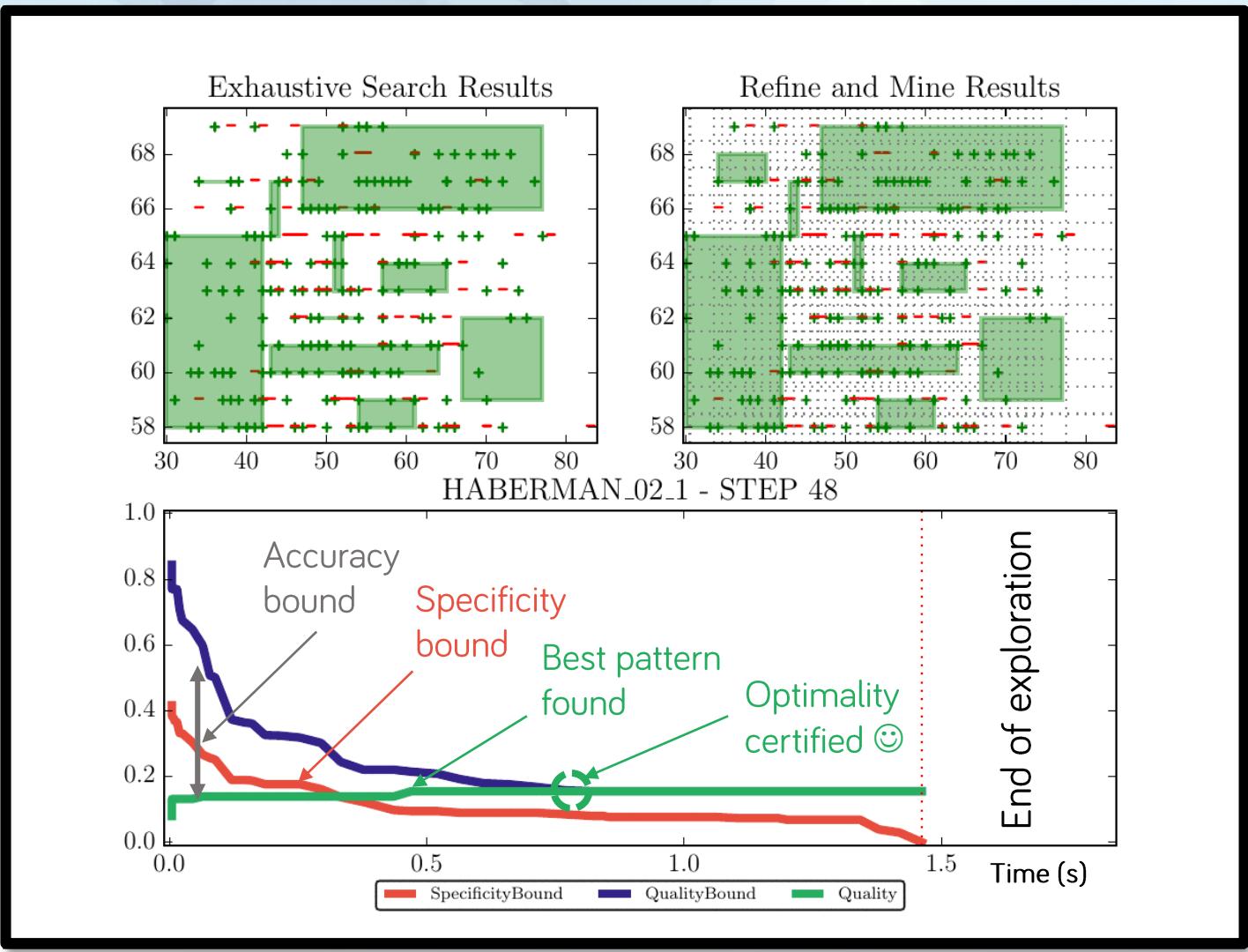
- 1 **Initialization.** Compute an initial rough **discretization** (e.g. Entropy based, equal frequency bins ...).
- 2 **Initial mining.** Compute patterns in the initial **discretization** and the first guarantees (e.g. use optimizations like closure operator).
- 3 **Refine:** add a new cut point in one dimension (e.g. random, more sophisticated strategy to reduce quickly the accuracy bound).
- 4 **Mine:** mine new patterns and compute the new guarantees (e.g. use index structure and closure operator).

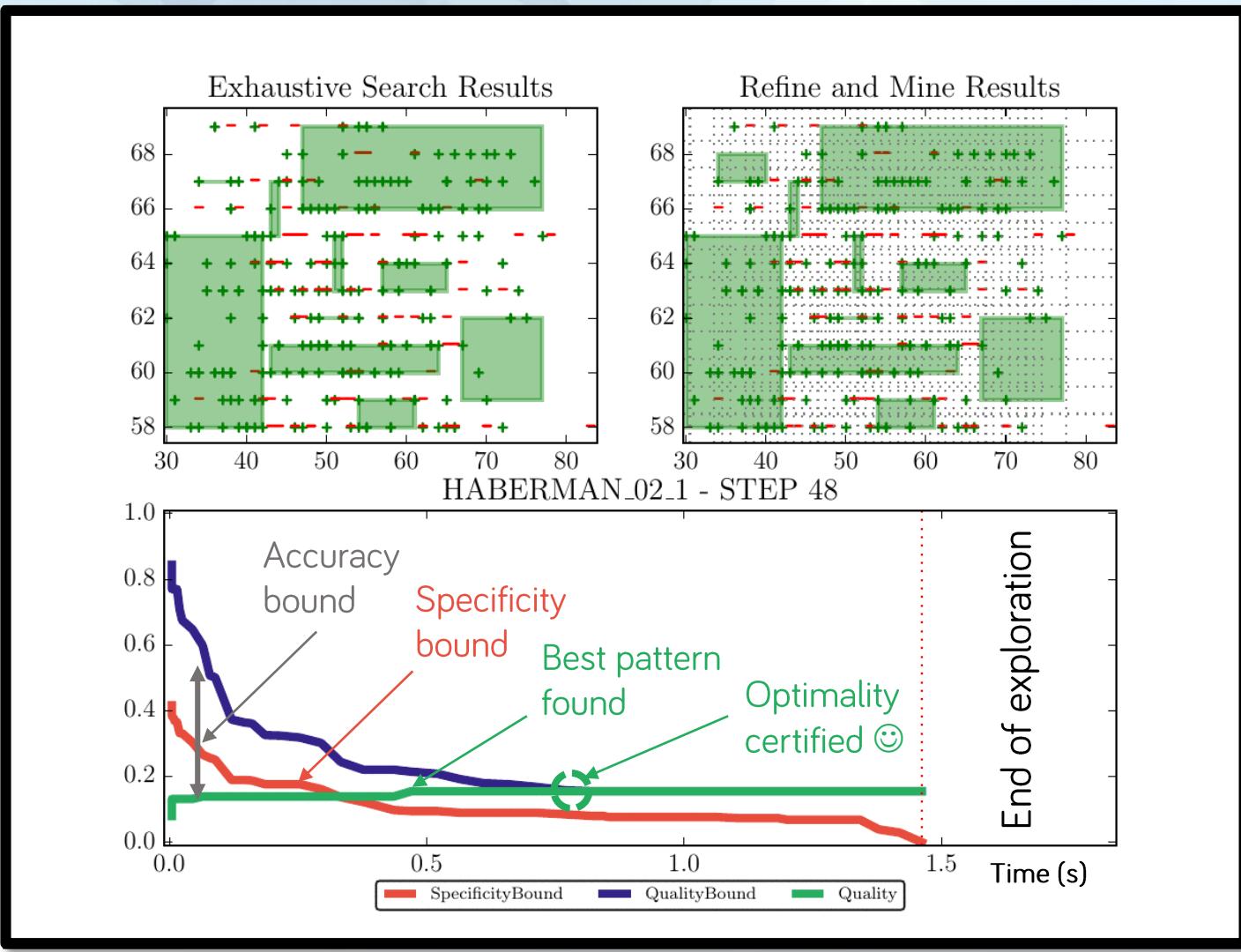
REFINE&MINE: Step-by-step



Refine&Mine

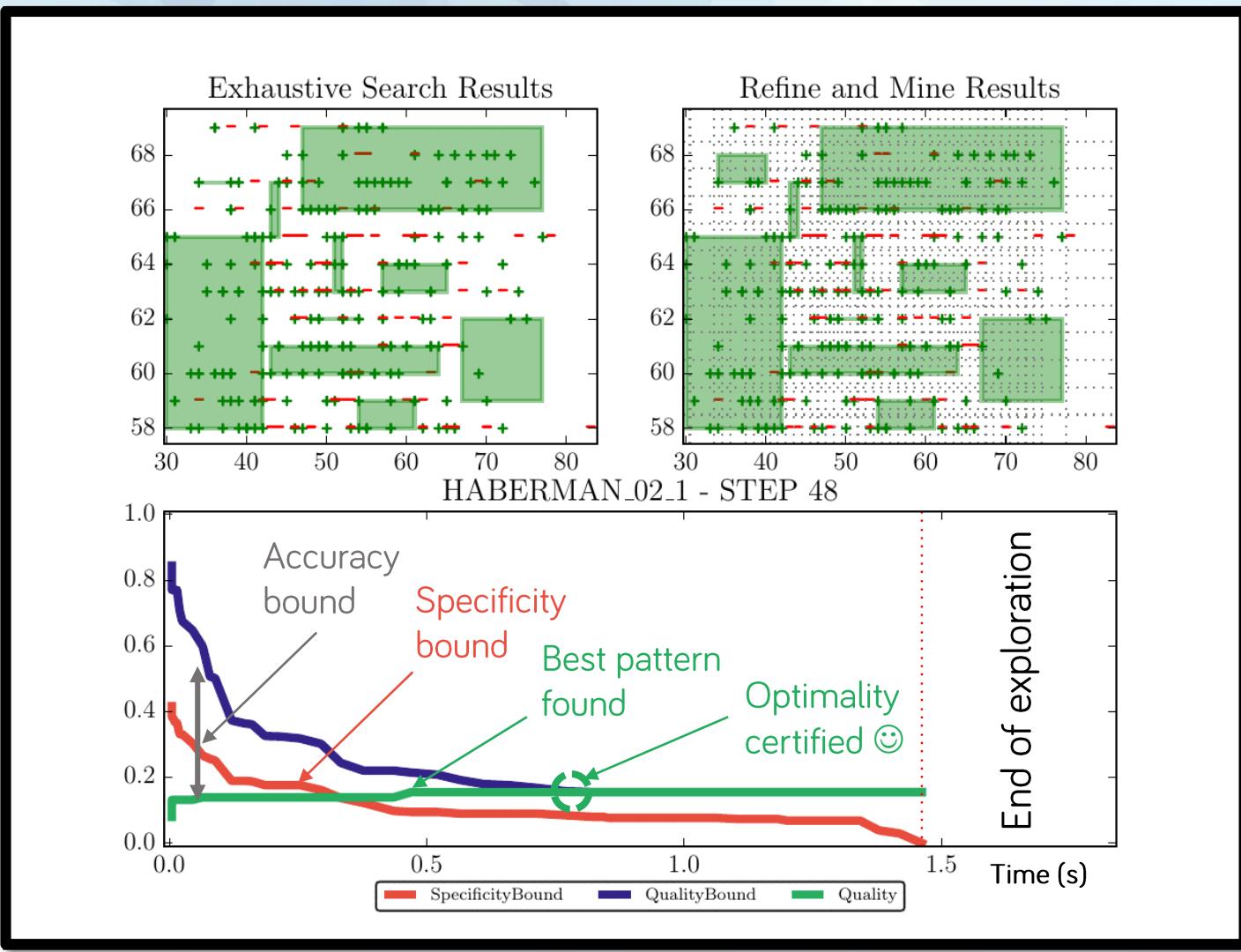
Wrap-up





Empirical Evaluation

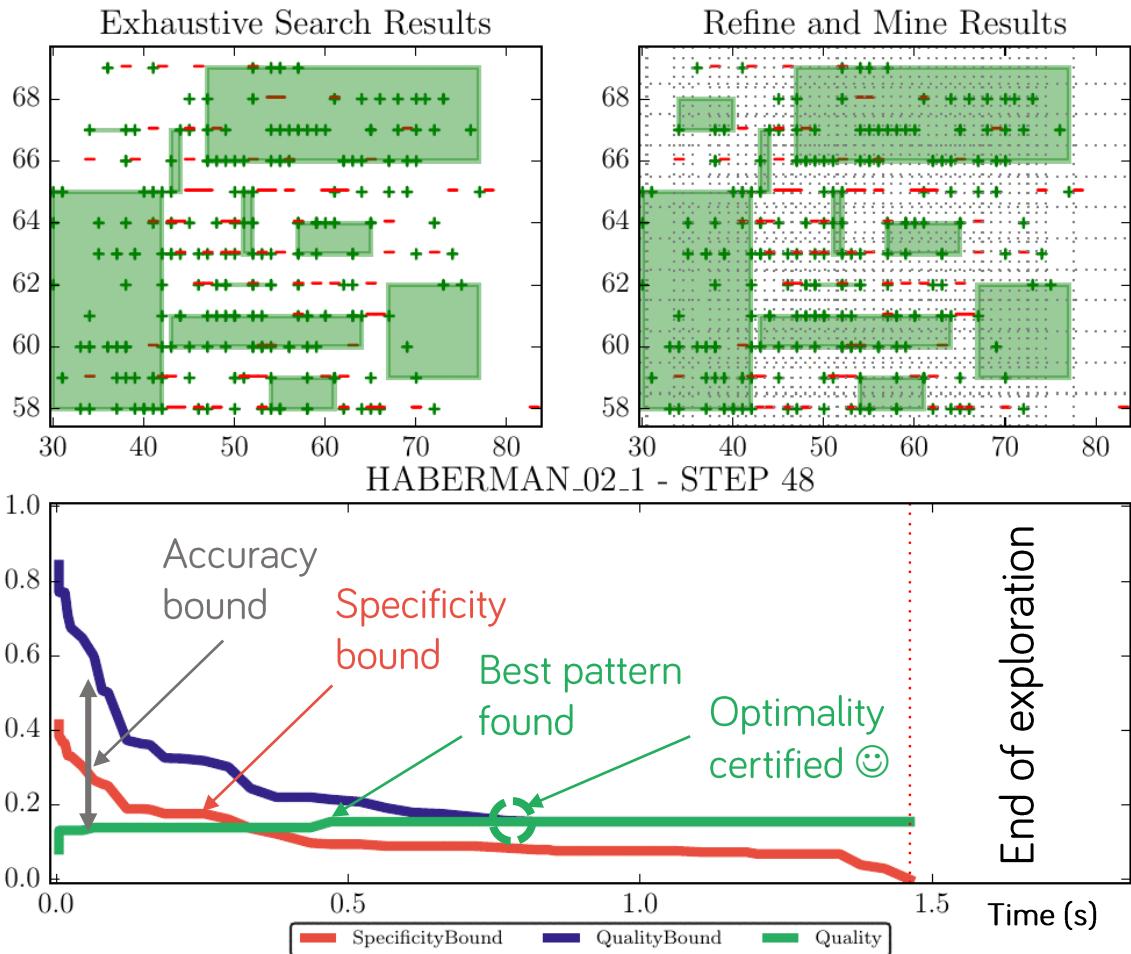
Experiments were performed over 25 UCI Datasets



Empirical Evaluation

Experiments were performed over 25 UCI Datasets

- Refine&Mine achieves to find a 80% ground truth in less than 20% of the time required for the exhaustive search algorithm to finish



Empirical Evaluation

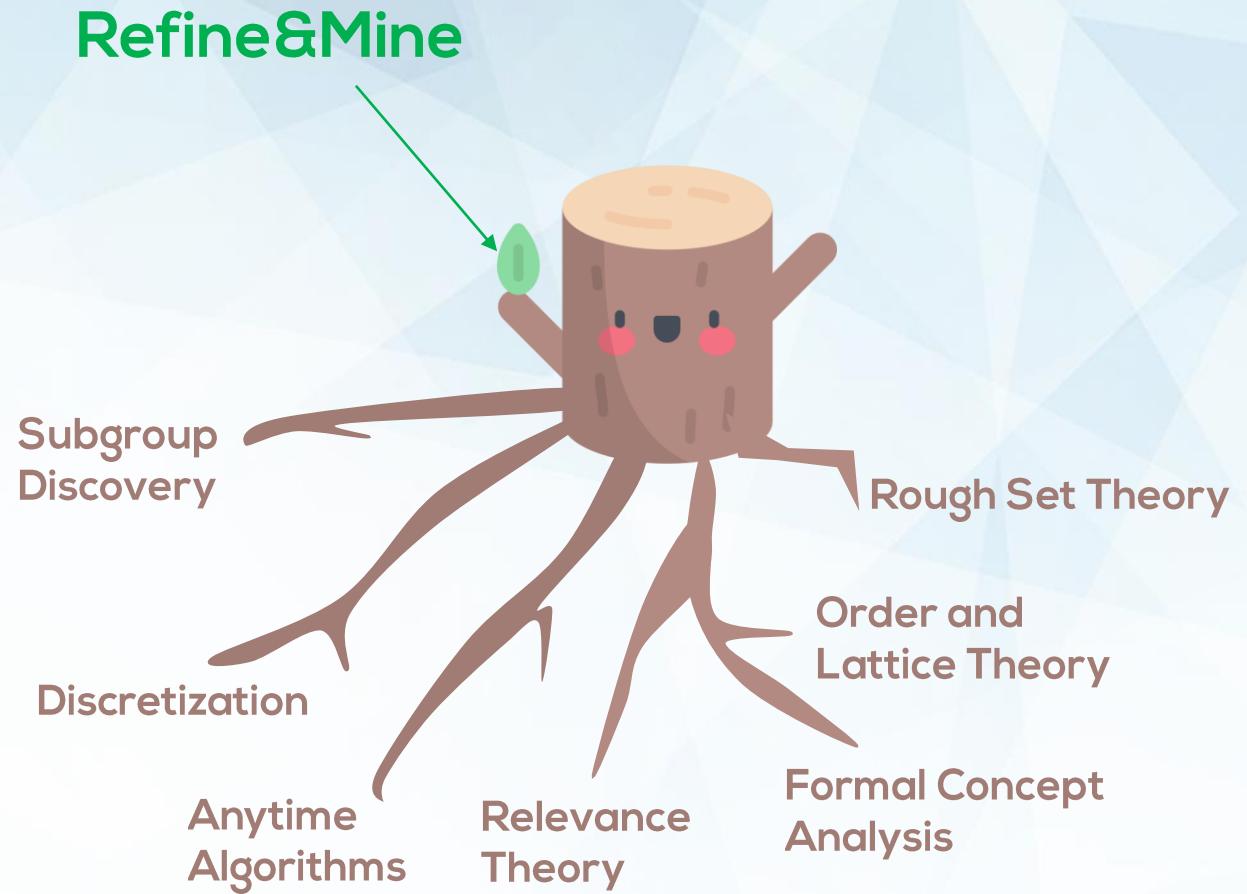
Experiments were performed over 25 UCI Datasets

- Refine&Mine achieves to find a **80% ground truth** in less than **20%** of the time required for the exhaustive search algorithm to finish
- Refine&Mine achieved **better results** both in finding the best pattern and a high diversified results set compared to **MCTS4DM** [*]

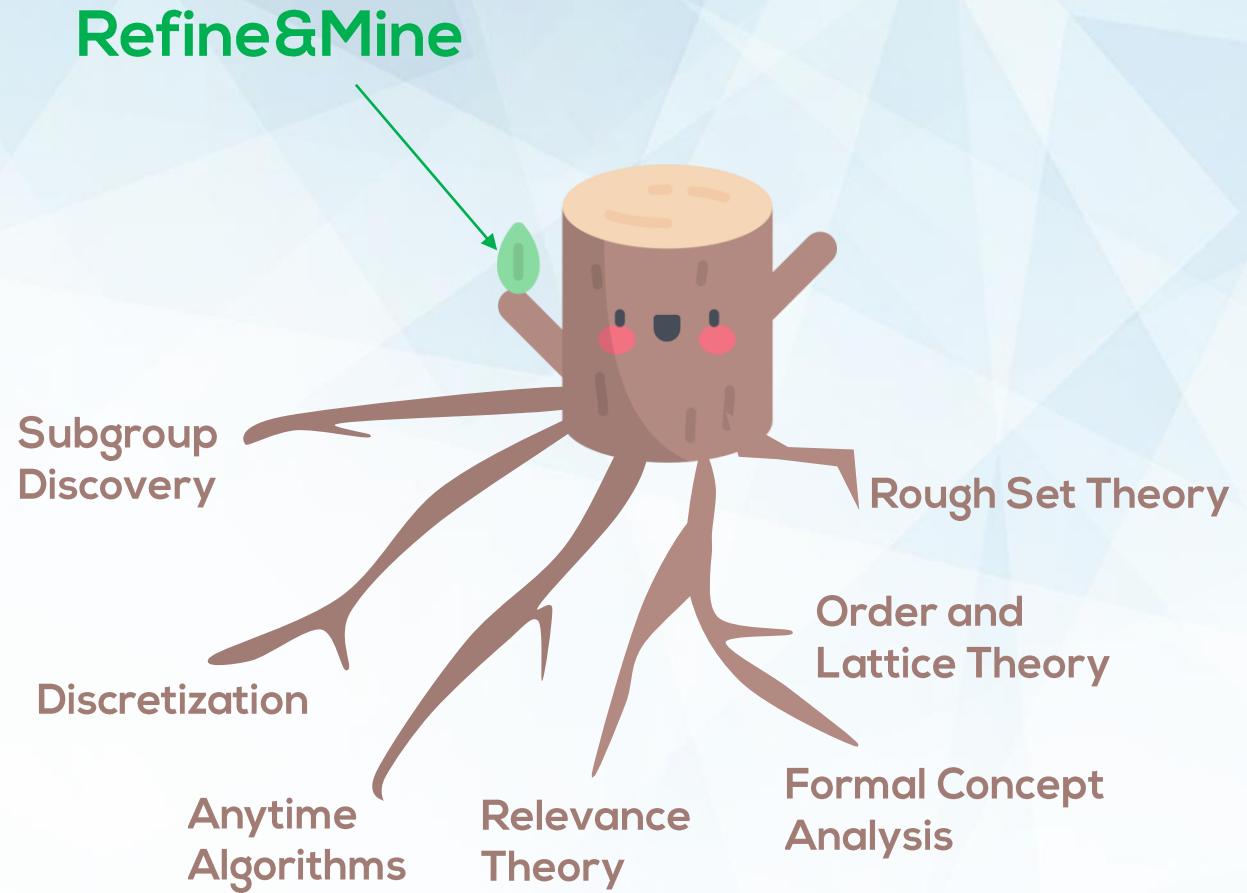
[*] G. Bosc, J.F. Boulicaut, C. Raïssi, M. Kaytoue. *Anytime discovery of a diverse set of patterns with MCTS*. In DMKD 2018.



CONCLUSION & PERSPECTIVES



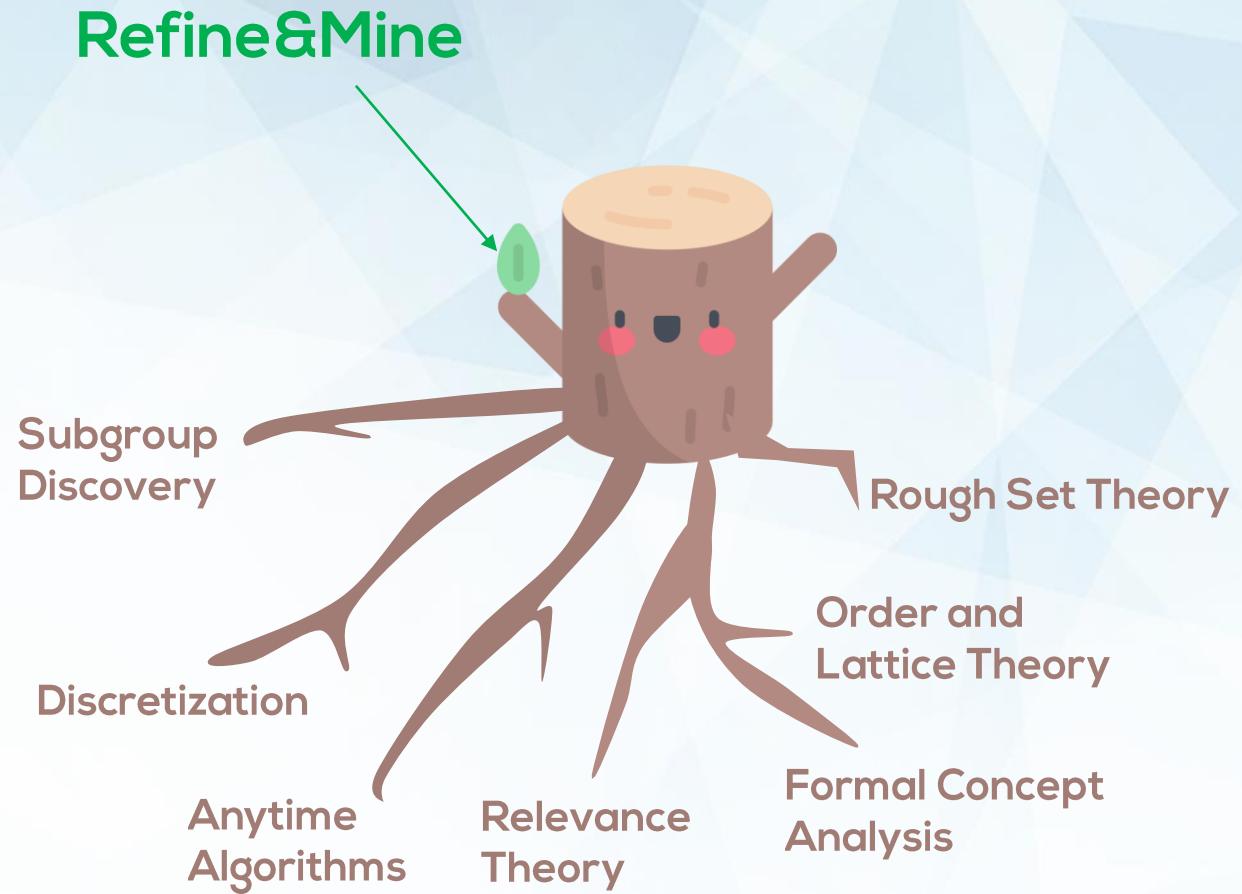
Proposed the new algorithm **Refine&Mine**



Proposed the new algorithm **Refine&Mine**



Anytime Algorithm



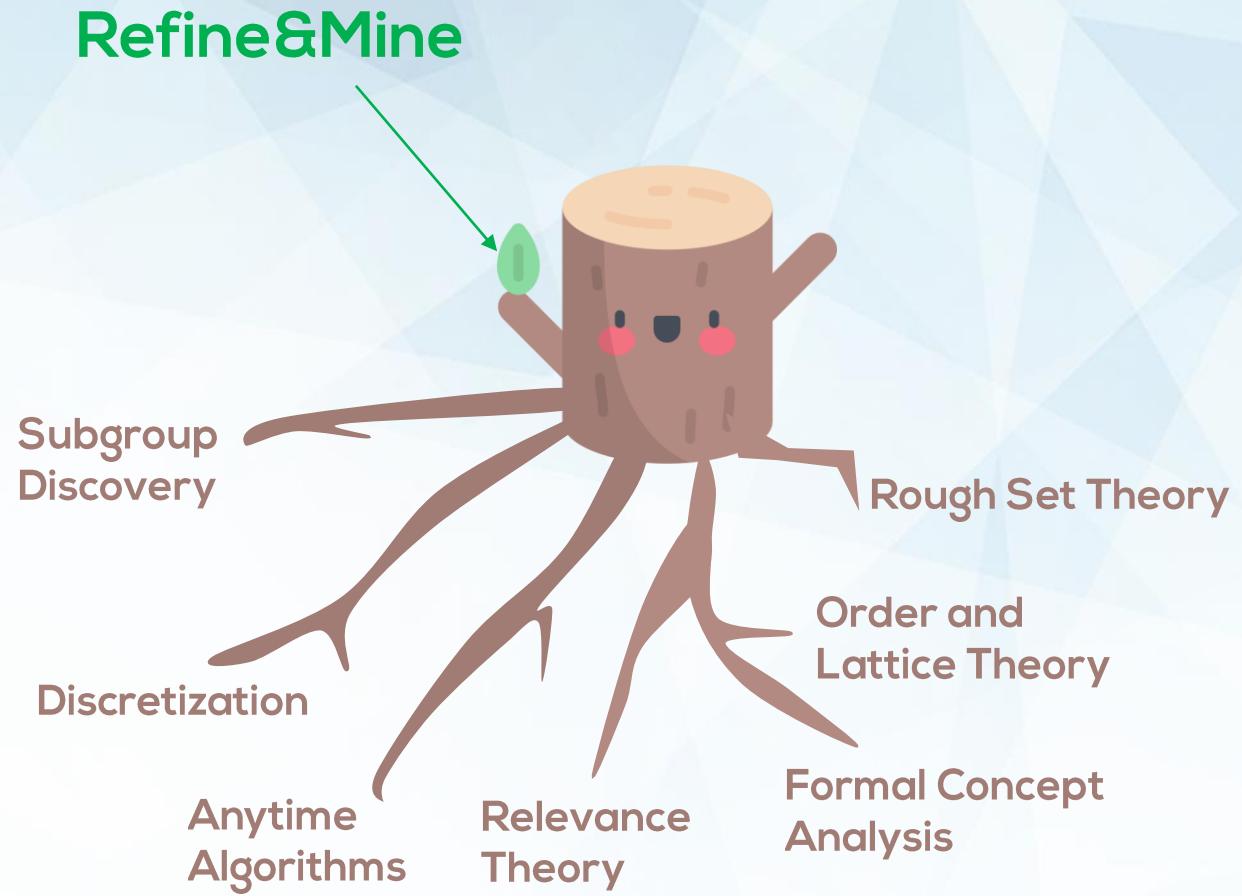
Proposed the new algorithm **Refine&Mine**



Anytime Algorithm



Guarantees on the found patterns anytime



Proposed the new algorithm **Refine&Mine**



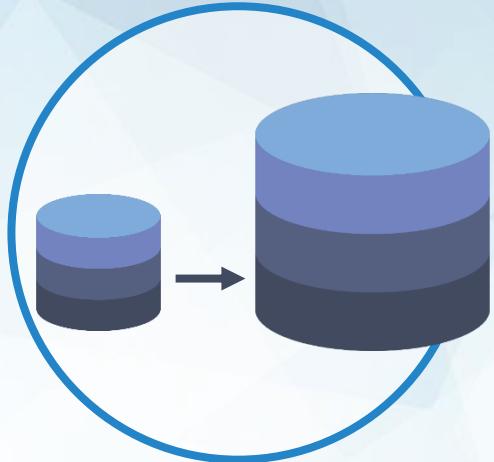
Anytime Algorithm



Guarantees on the found patterns anytime

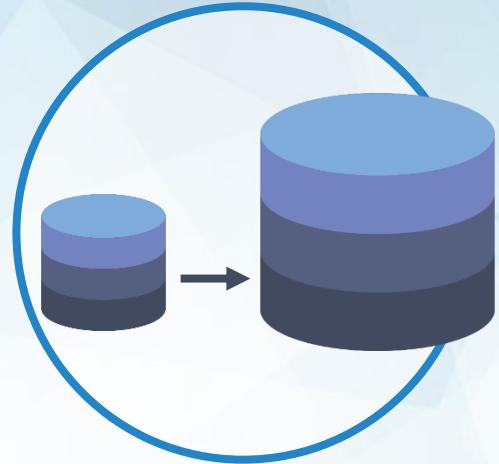


It performs an exhaustive search If given enough time



Handle **High**
dimensional data.

Limit. We need **2 cut points in each dimension** to compute the first guarantees:
 d attributes $\Leftrightarrow 6^d$ patterns



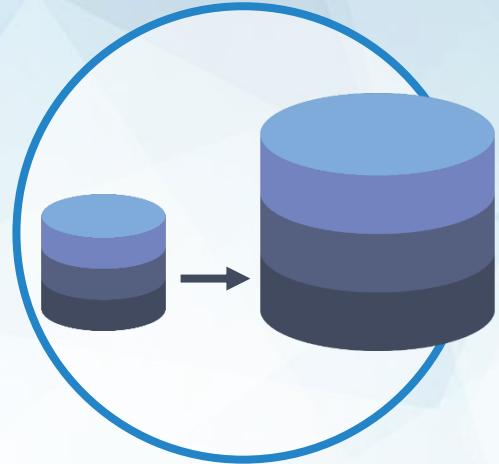
Handle **High**
dimensional data.

Limit. We need **2 cut points in each dimension** to compute the first guarantees:
 d attributes $\Leftrightarrow 6^d$ patterns



Handle datasets with
hybrid attributes.

e.g. binary, nominal, numerical,
sequences, graphs.



Handle **High** dimensional data.

Limit. We need **2 cut points** in **each dimension** to compute the first guarantees:
 d attributes $\Leftrightarrow 6^d$ patterns



Handle datasets with **hybrid** attributes.

e.g. binary, nominal, numerical, sequences, graphs.



Provide **guarantees** on **pattern sets** rather than **all patterns**.

GRAÇIAS

AČIŪ

DĚKUJI

THANKS FOR YOUR TIME



?

QUESTIONS

Contact : adnene.belfodil@insa-lyon.fr & aimene.belfodil@insa-lyon.fr

Materials: <https://github.com/Adnene93/RefineAndMine>

Looking forward to meeting up with you during the poster session :-)

Ευχαριστώ

MULUMIRI

TACK

KIITOS

OBRIGADO

HVALA

BEDANKT

DANKE

KÖSZ

MERCI

AITÄH

MERCY

GO RAIBH MAITH AGAT

AČIŪ

GO RAIBH MAITH AGAT

DI

GRAZIE

DIES

GRAZIE

OBIGADO

KIITOS

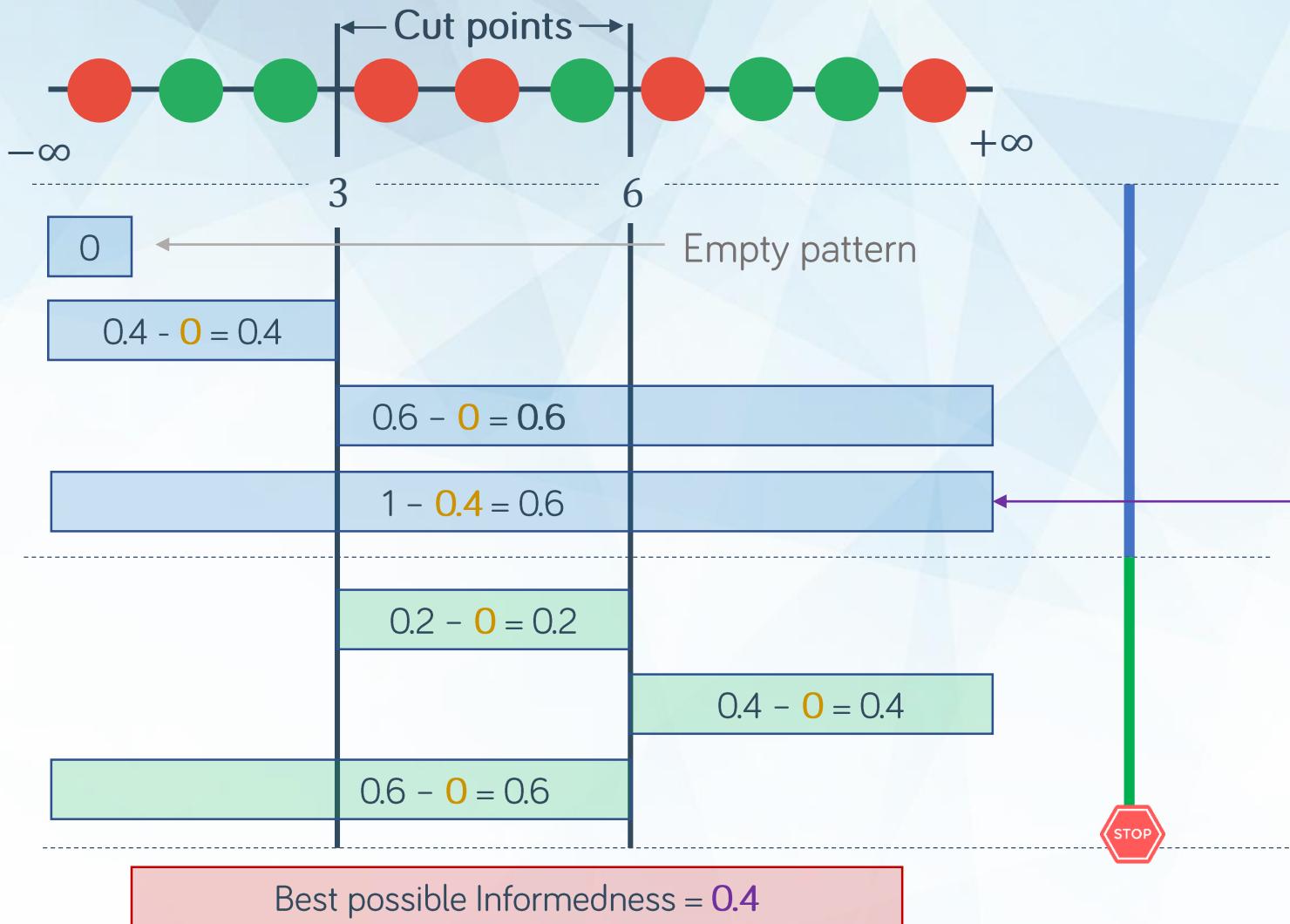
HVALA

MERCI



ADDITIONAL DETAILS

Accuracy Guarantee – Example



$$tpr(\text{pattern}) - fpr(\text{core}(\text{pattern})) = \text{informedness bound}(\text{pattern})$$

- Best found Informedness = 0.2
 - Best unseen Informedness = 0.4
 - Highest Informedness bound = 0.6
-
- Accuracy = 0.4 - 0.2 = 0.2
 - Accuracy bound = 0.6 - 0.2 = 0.4

First guarantee

Specificity = $\overrightarrow{distance}$ (All found patterns, All pattern in the dataset)

Specificity \leq Specificity bound

Refine&Mine computes this bound
using only already found patterns.

$$\text{distance}(c, d) = \frac{|c \Delta d|}{n} \in [0,1]$$

Pattern extents Symmetric difference Dataset size

The set of found patterns is compared to the set of ground truth patterns in this direction as follow:

- Each possible pattern is compared to its nearest found pattern.
- The top nearest distance is then the specificity

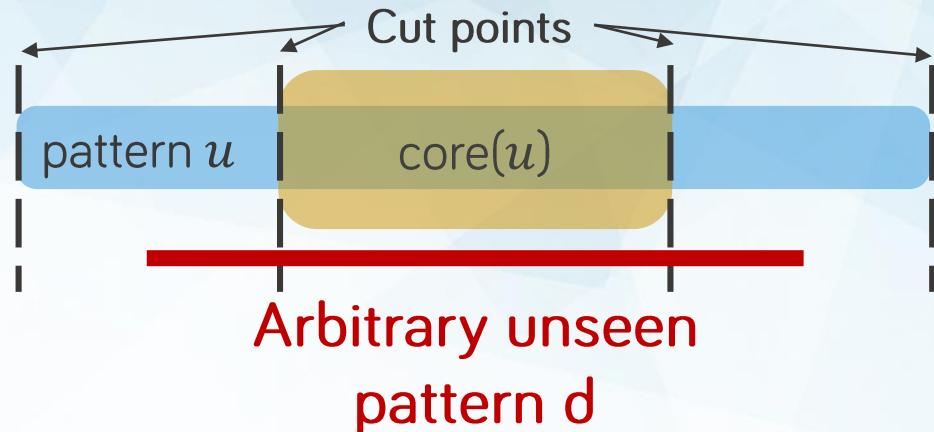
$$\overrightarrow{\text{distance}}(\text{Found}, \text{All}) = \sup_{d \in \text{all}} \inf_{c \in \text{found}} \text{distance}(c, d)$$

Specificity = Directed Hausdorff distance

Distance to the nearest found pattern

*The lower the specificity, the better found patterns are diversified

Specificity = 0 \Leftrightarrow Found Patterns \supseteq All Patterns.



$$\text{Specificity} = \sup_{d \in \text{all}} \inf_{c \in \text{found}} |c \Delta d| / n$$

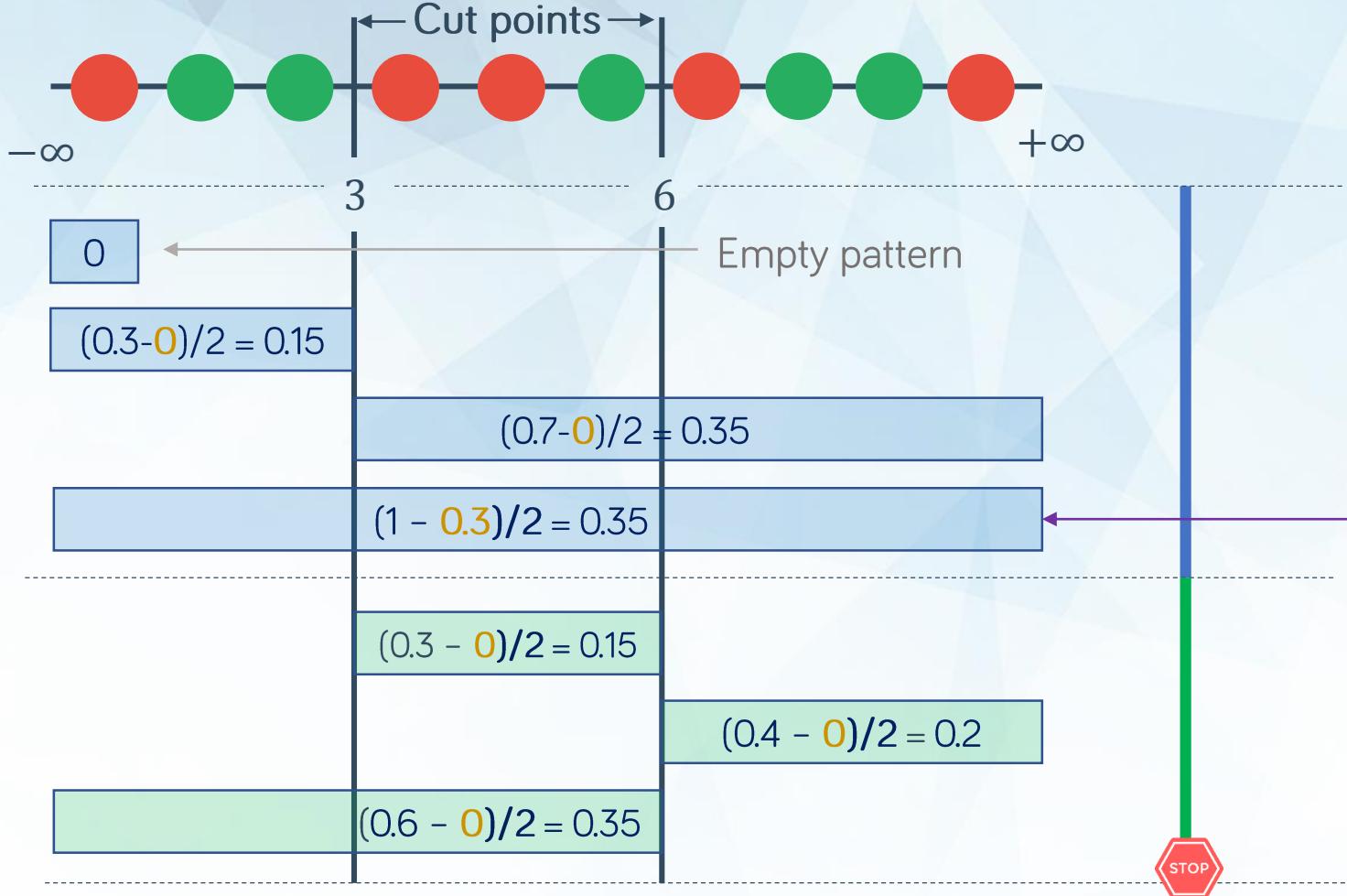
- $\text{Core}(\text{found pattern}) \subseteq \text{some unseen patterns} \subseteq \text{found pattern}$
 - $\inf_{d \in \text{all}} \text{distance}(c, d) \leq \text{distance}(u, d) = \frac{|u| - |d|}{n}$
 - $\inf_{d \in \text{all}} \text{distance}(c, d) \leq \text{distance}(\text{core}(u), d) = \frac{|d| - |\text{core}(u)|}{n}$
- Upper approximation



$$\text{Specificity} \leq \sup_{u \in \text{found}} \frac{|u| - |\text{core}(u)|}{2 \cdot n}$$

*This is a simplified formula.

Specificity Guarantee – Example



$$(|\text{pattern}| - |\text{core}(\text{pattern})|)/2 \cdot 10$$

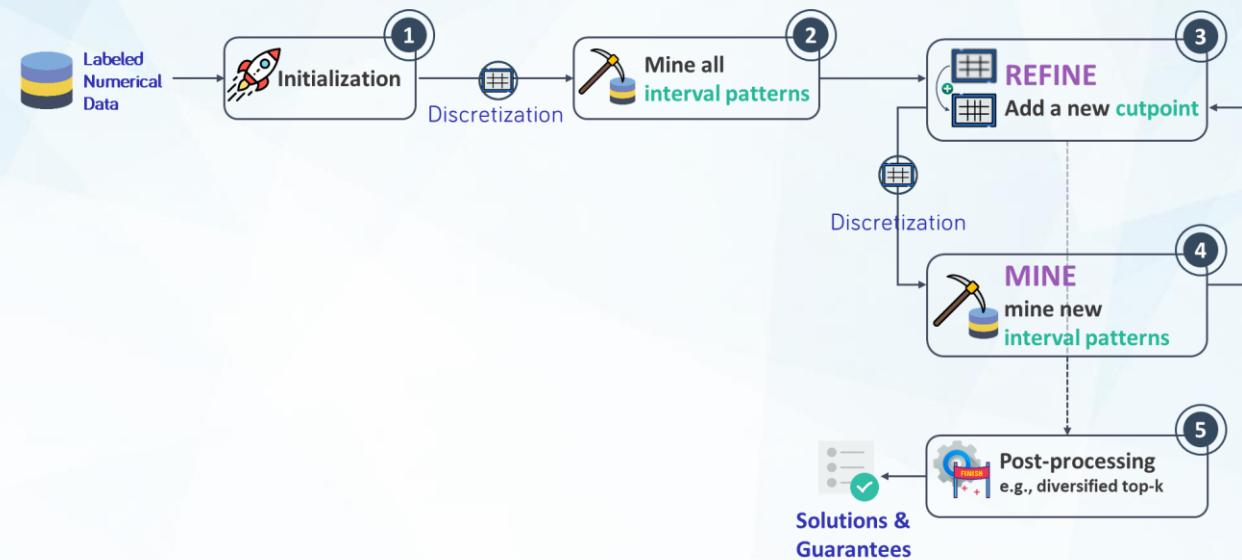
= **nearest distance bound**(pattern)

- Specificity bound = 0.35

Second
guarantee

Refine&Mine

Step-by-step





1 Initialization



Compute an initial discretization [10,11]. e.g. Entropy based discretization, frequency based discretization...

[10] U.M. Fayyad, K.B. Irani. **Multi-Interval discretization of continuous-valued attributes for classification learnings.** In IJCAI 1993.

[11] Y. Yang, G.I. Webb, X., Wu: **Descretization methods.** In DMKD 2010.



2 Mine

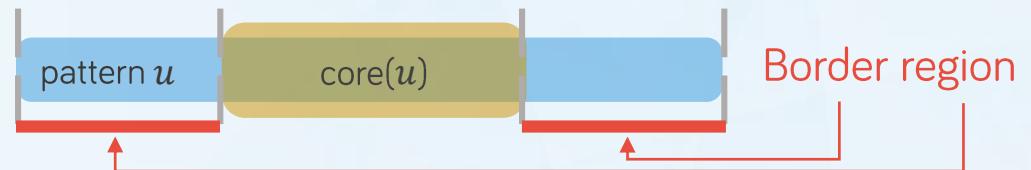
- Compute all **interval patterns** using **MinIntChange** algorithm [6] considering only the input-discretization selected points
- Compute the first guarantees.

[6] M. Kaytoue, Sergei O. Kuznetsov, Amedeo Napoli. **Revisiting Numerical Pattern Mining with Formal Concept Analysis**. In IJCAI 2011.



3 Refine

- Add one **cutpoint** in one dimension generating a new discretization, this can be done by:
 - Adding a **random cut-point** on a **random dimension**.
 - Adding a cut-point belonging to the **border region** of the pattern yielding the **highest accuracy bound**.

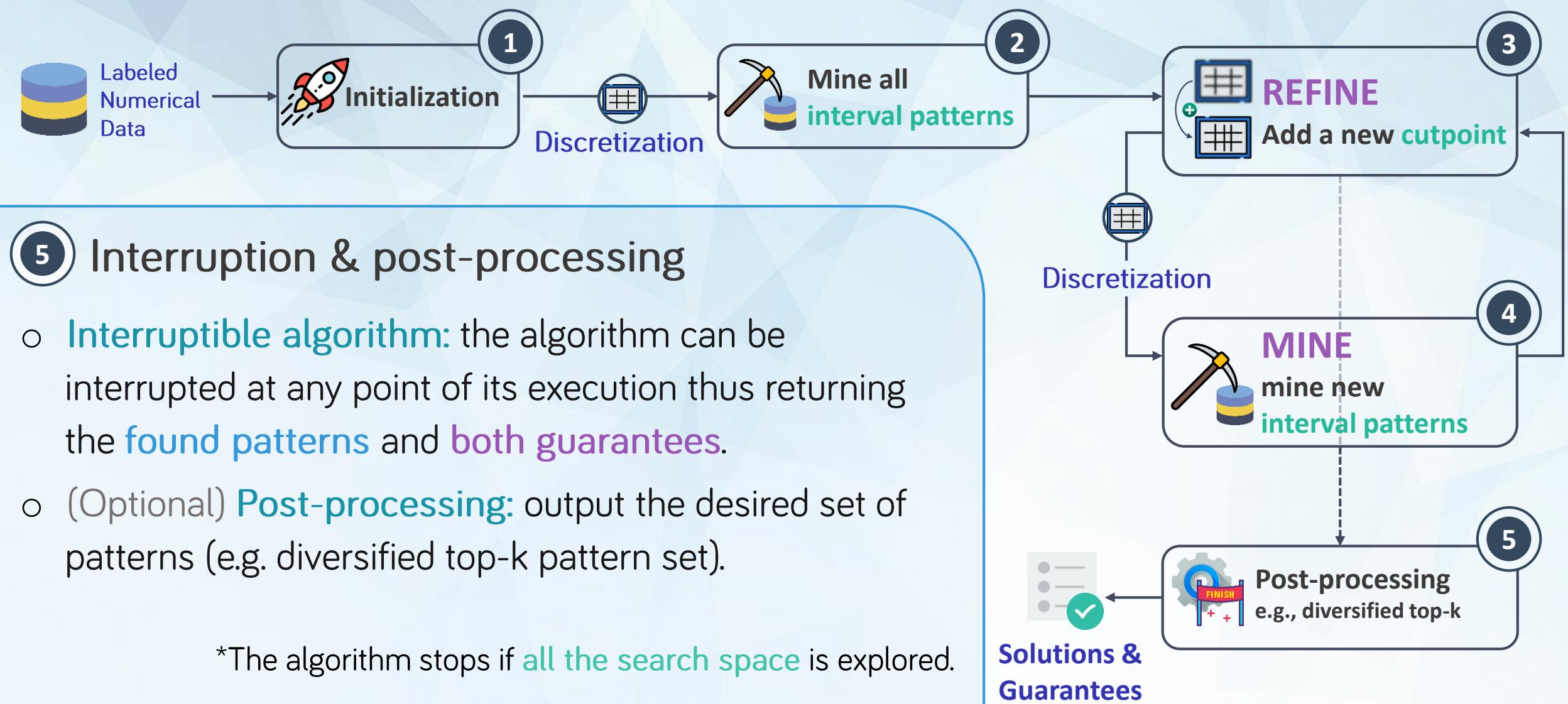


REFINE&MINE: Step-by-step (4)



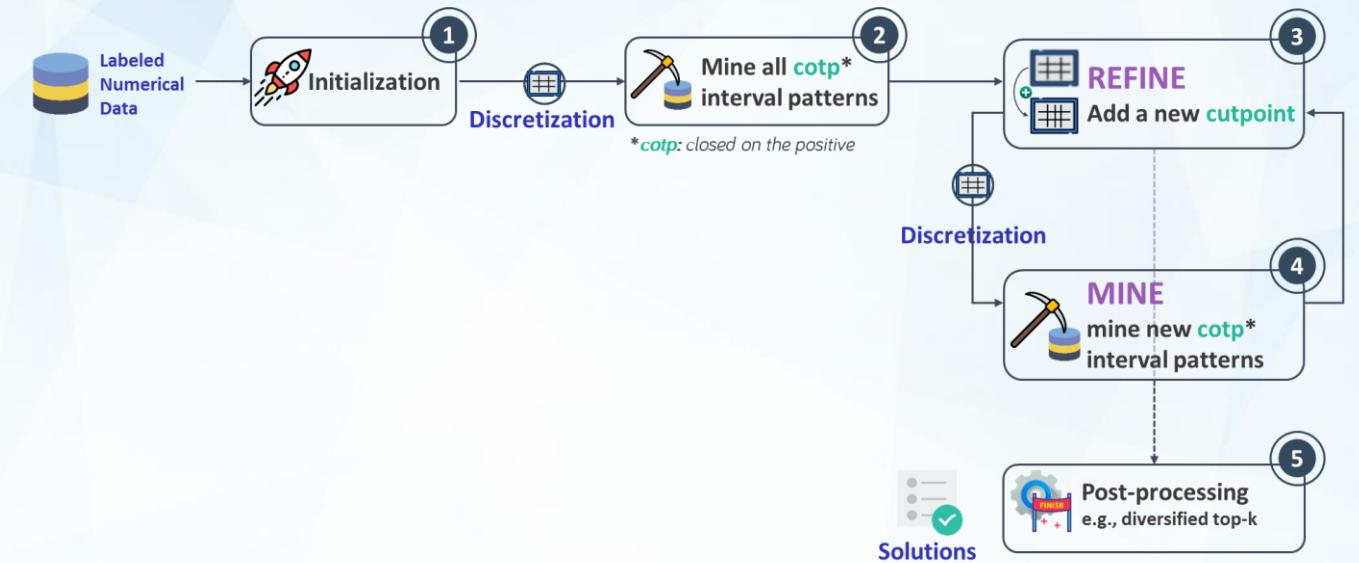
- Mine new interval patterns using [6]
by only visiting those whose left or right bound on the split dimension is the new selected cut-point.
 - Update both guarantees
An index structure is used to update quickly the guarantees

[6] M. Kaytoue, Sergei O. Kuznetsov, Amedeo Napoli. [Revisiting Numerical Pattern Mining with Formal Concept Analysis](#). In IJCAI 2011.



Refine&Mine

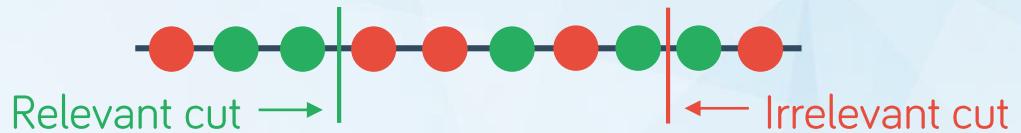
Step-by-step





1 Initialization

- Compute **relevant cuts** [10] (cuts separating two classes).



- Compute an initial discretization [11].

[10] U.M. Fayyad, K.B. Irani. **Multi-Interval discretization of continuous-valued attributes for classification learnings.** In IJCAI 1993.

[11] Y. Yang, G.I. Webb, X., Wu: **Descretization methods.** In DMKD 2010.

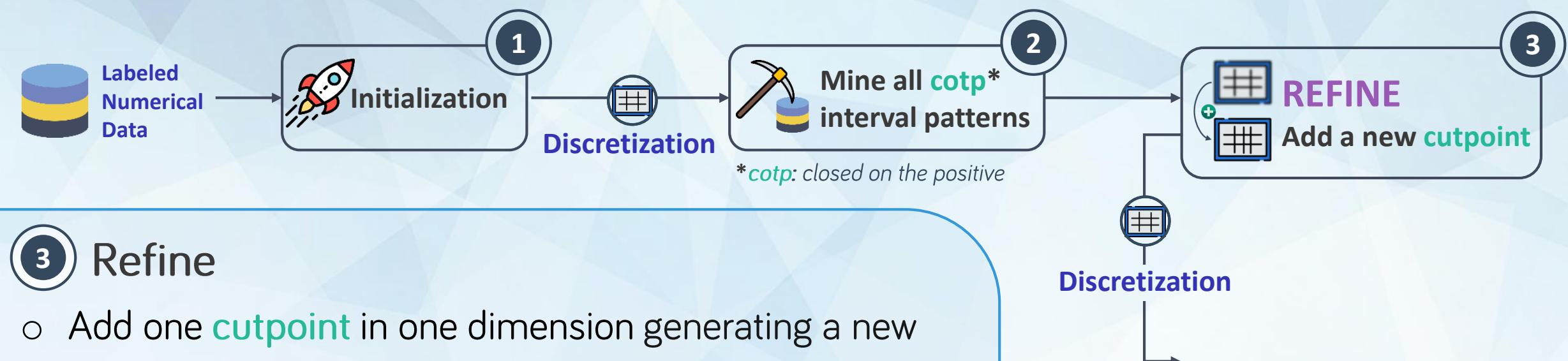


2 Mine

- Compute all closed-on-the-positive (cotp [12]) patterns using [MinIntChange](#) algorithm [6].
- Compute the first guarantees.

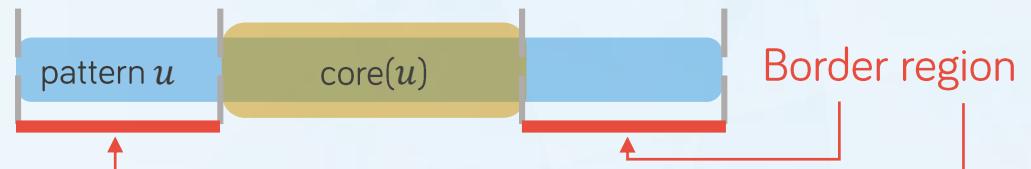
[6] M. Kaytoue, Sergei O. Kuznetsov, Amedeo Napoli. [Revisiting Numerical Pattern Mining with Formal Concept Analysis](#). In IJCAI 2011.

[12] C.G.. Garriga, P. Kralj, N. Lavrac. [Closed sets for labeled data](#). In JMLR 2008.



3 Refine

- Add one **cutpoint** in one dimension generating a new discretization, this can be done by:
 - Adding a **random cut-point** on a **random dimension** or,
 - adding a cut-point belonging to the **border region** of the pattern yielding the **highest accuracy bound**.

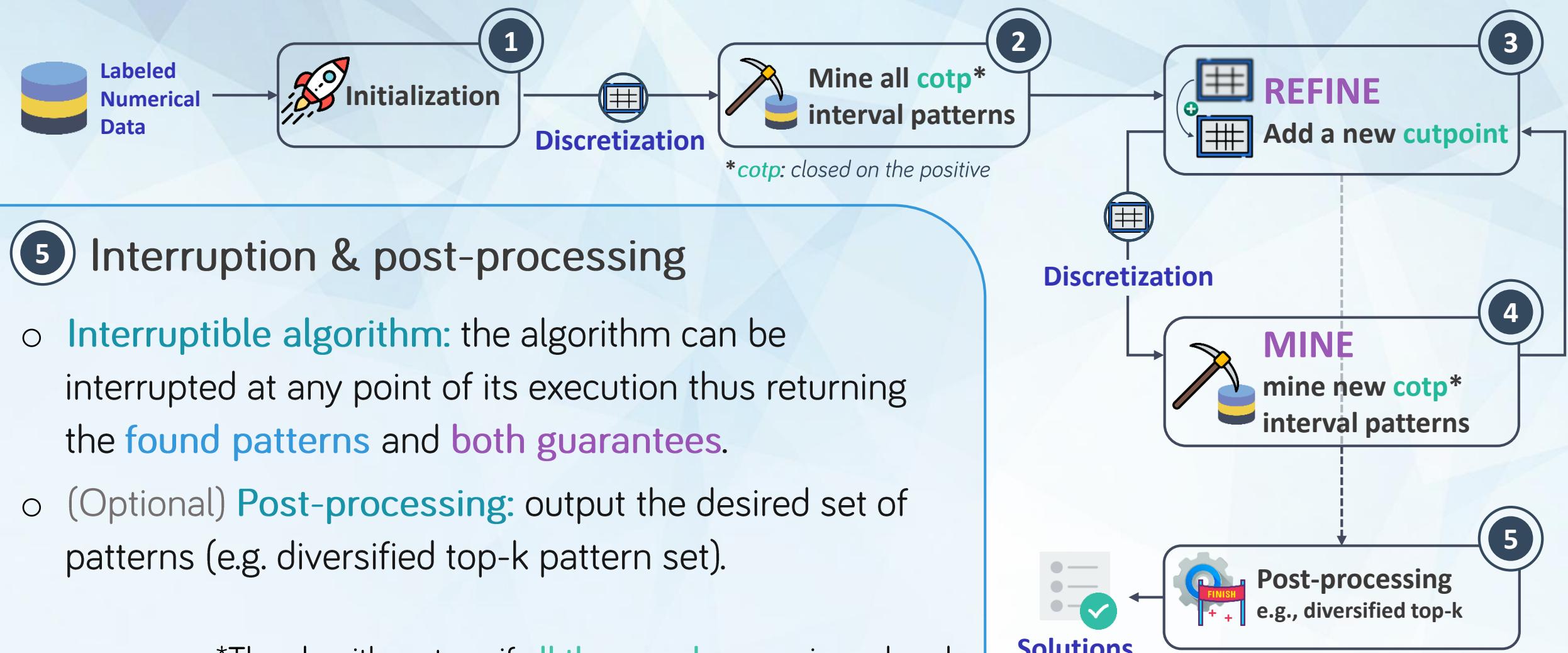




4 Mine

- Mine **new closed-on-the-positive patterns using [6]**
by only visiting those whose left or right bound on the split dimension is the new selected cut-point.
- Update both guarantees
An index structure is used to update quickly the guarantees

[6] M. Kaytoue, Sergei O. Kuznetsov, Amedeo Napoli. **Revisiting Numerical Pattern Mining with Formal Concept Analysis**. In IJCAI 2011.



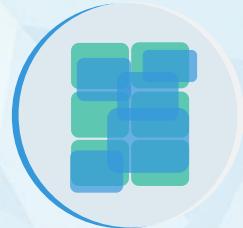


EMPIRICAL EVALUATION



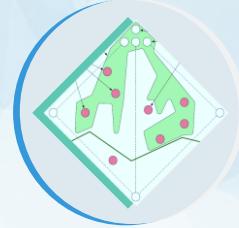
Accuracy Guarantee

A **bound** on the difference of quality between the **best quality pattern in D** and the **best found pattern in S** .



Specificity Guarantee

A **bound** on the distance between the **set of all patterns D** and the **set of already found patterns S** .



Comparison with MCTS4DM

MCTS4DM is the closest approach to **Refine&Mine** in terms of the underlying paradigm, i.e. **Anytime Algorithms**.



Accuracy Guarantee

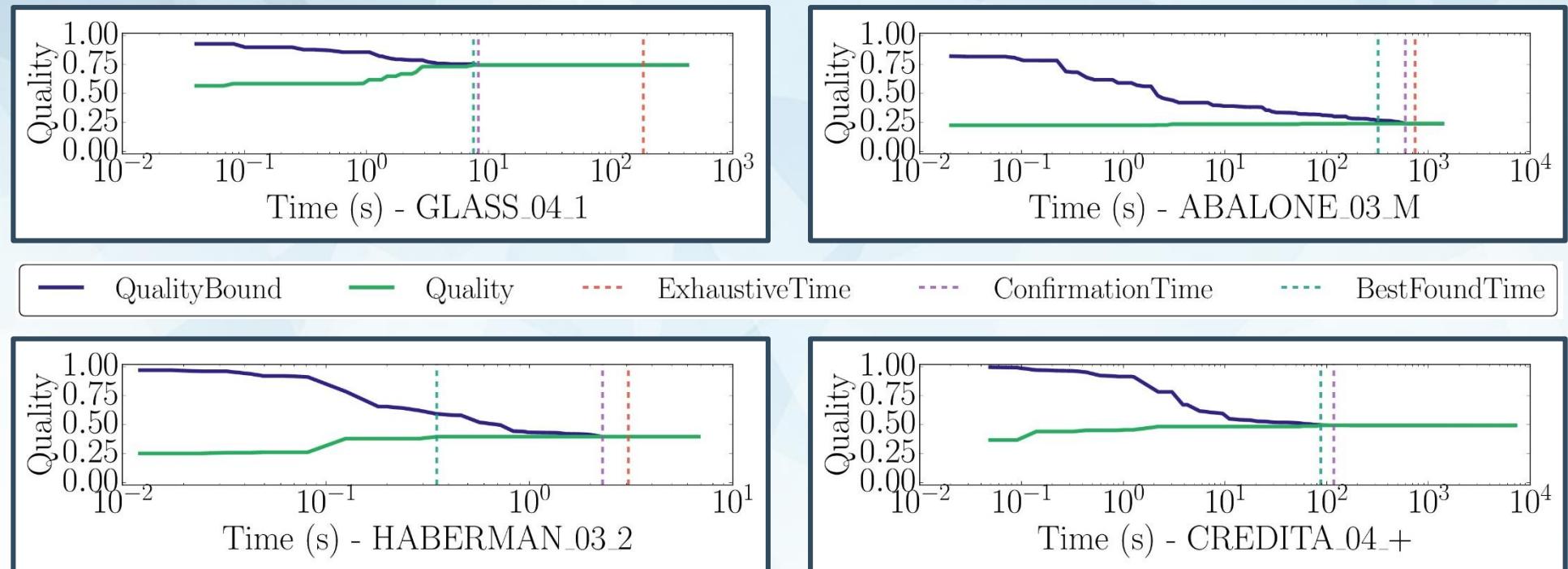
A **bound** on the difference of quality between the **best quality pattern in D** and the **best found pattern in S** .

Bottom line

Experiments were performed over 25 UCI Datasets

- 5 out of 25 Datasets: Exhaustive Algorithm achieved a better time in ensuring the discovery of the best pattern.
- 6 out of 25 Datasets: Exhaustive Algorithm was not able to finish within 2 Hours
- Refine&Mine required in average 40% of the time to guarantee the discovery of the best pattern

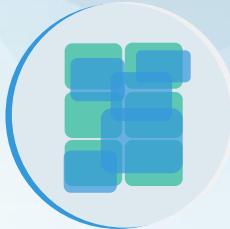
Experiment Plan: Accuracy Guarantee



Accuracy
Guarantee

Bottom line Experiments were performed over 25 UCI Datasets

- 5 out of 25 Datasets: Exhaustive Algorithm achieved a better time in ensuring the discovery of the best pattern.
- 6 out of 25 Datasets: Exhaustive Algorithm was not able to finish within 2 Hours
- Refine&Mine required in average 40% of the time to guarantee the discovery of the best pattern



Specificity Guarantee

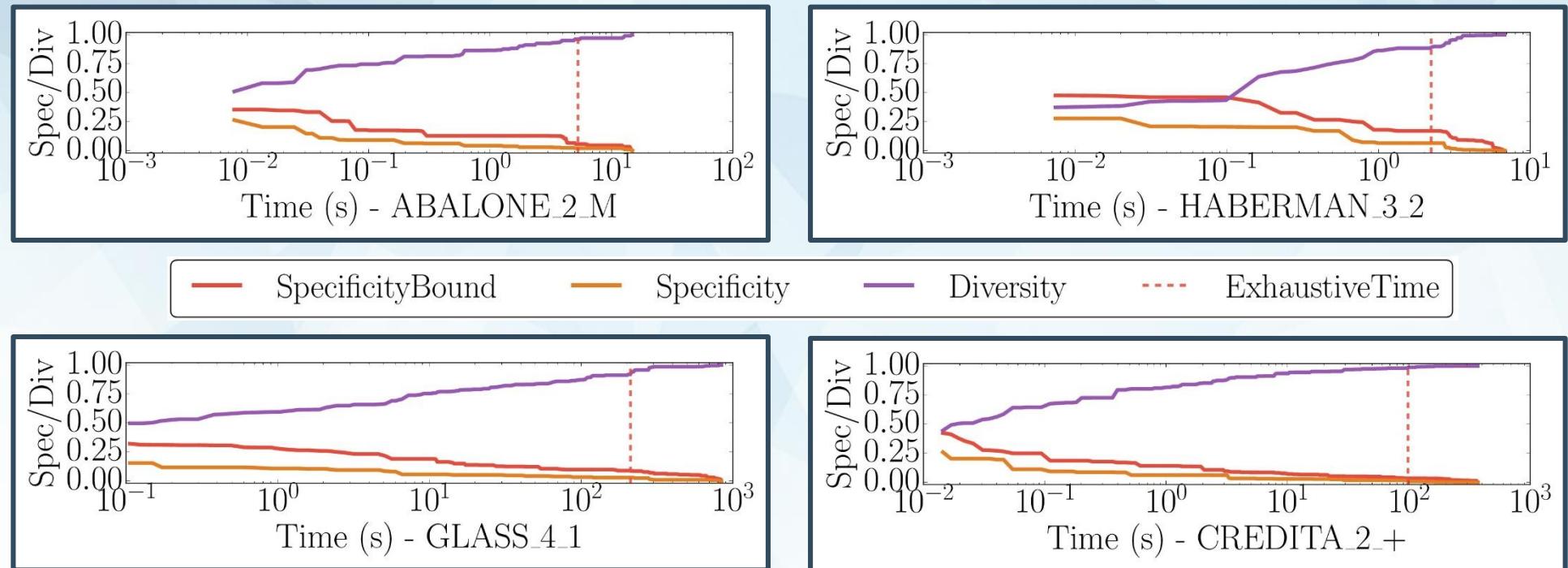
A **bound** on the distance between the **set of all patterns D** and the **set of already found patterns S** .

Bottom line

Experiments were performed over 15 UCI Datasets

Refine&Mine achieves to find a **80%** approximation of the top-k diversified patterns set found by the exhaustive search algorithm in less than **20%** of the required time of the exhaustive search algorithm

Experiment Plan: Specificity Guarantee



Bottom line Experiments were performed over 15 UCI Datasets

- Refine&Mine achieves to find a 80% approximation of the top-k diversified patterns set found by the exhaustive search algorithm in less than 20% of the required time of the exhaustive search algorithm



Specificity
Guarantee



Comparison with MCTS4DM

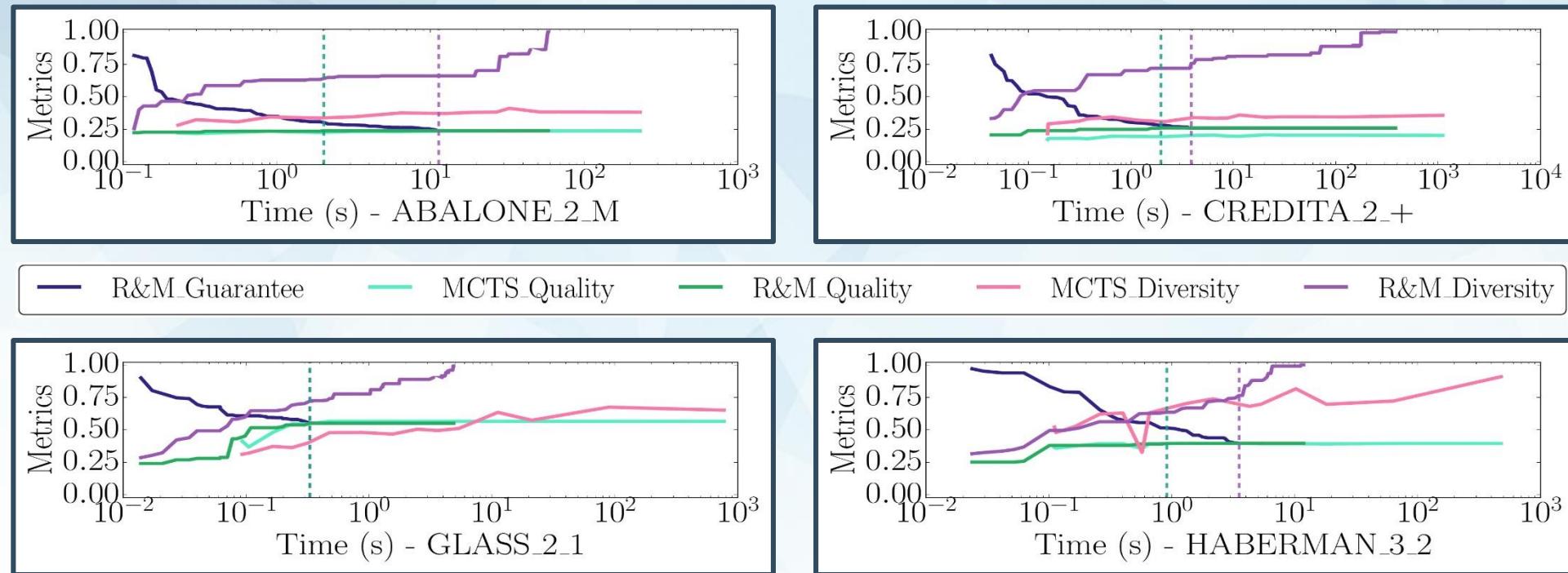
MCTS4DM is the closest approach to Refine&Mine in terms of the underlying paradigm, i.e. Anytime Algorithms.

Bottom line

Experiments were performed over 6 UCI Datasets

- In all configurations: Refine&Mine ensured a full traversal of the search space before MCTS4DM
- In all configurations: In average, Refine&Mine achieved 25% more diversity (after 10s of execution time) than MCTS4DM

Experiment Plan: Comparison with MCTS4DM

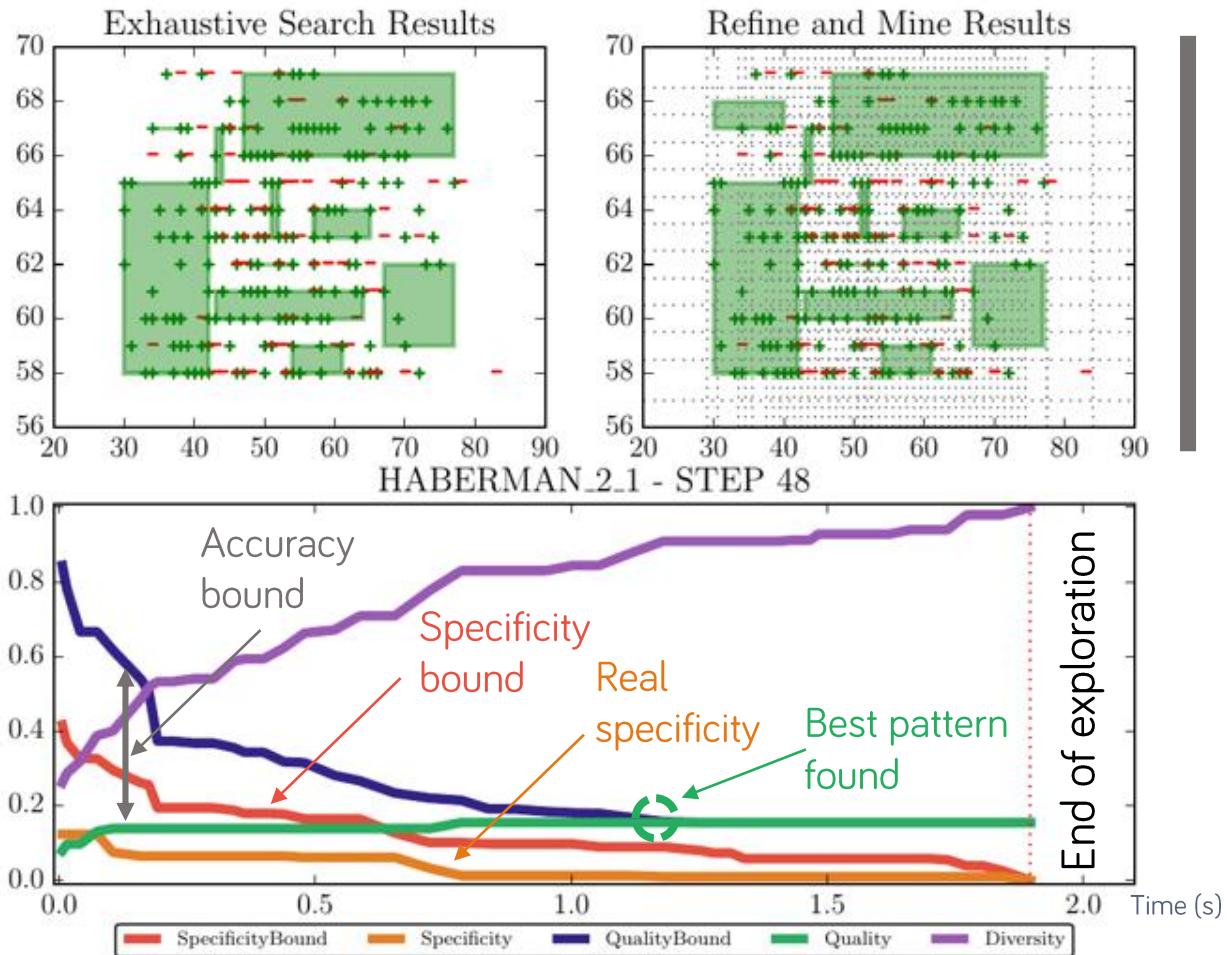


Bottom line Experiments were performed over 6 UCI Datasets



- In all configurations: Refine&Mine ensured a full traversal of the search space before MCTS4DM
- In all configurations: In average, Refine&Mine achieved 25% more diversity (after 10s of execution time) than MCTS4DM

R&M vs.
MCTS4DM



Empirical Evaluation

Experiments were performed over 25 UCI Datasets

- Refine&Mine achieves to find a 80% ground truth in less than 20% of the time required for the exhaustive search algorithm to finish
- Refine&Mine achieved better results both in finding the best pattern and a high diversified results set compared to MCTS4DM [∗]

[∗] G. Bosc, J.F. Boulicaut, C. Raïssi, M. Kaytoue. *Anytime discovery of a diverse set of patterns with MCTS*. In DMKD 2018.