# Probabilistic Rule Lists using the MDL Principle

#DS2018

**John Aoga**, Tias Guns, Siegfried Nijssen & Pierre Schaus

---

## ② Motivation

| Year | Month | Day | Part of Day | Minute | Door Opened |
|---|---|---|---|---|---|
| 2018 | October | Monday | Morning | 10 | Closed |
| | October | Monday | | | Opened |
| | October | Tuesday | | 30 | |
| | October | Tuesday | | 40 | |
| | October | Wednesday | | 50 | |
| | October | Wednesday | | 10 | |
| | October | Thursday | | 20 | |
| | October | Thursday | | 30 | |
| | October | Friday | | 40 | |
| | November | Monday | | 50 | |
| | November | Monday | | 10 | |
| | November | Tuesday | Afternoon | 20 | |
| | November | Tuesday | Morning | 30 | |
| | November | Wednesday | Afternoon | 40 | |
| | November | Wednesday | Morning | 50 | |
| | November | Thursday | Afternoon | 20 | |
| | November | Thursday | Morning | 20 | |
| | November | Friday | Afternoon | 30 | |
| | December | Monday | Afternoon | 40 | |
| | December | Monday | Afternoon | 50 | |
| | December | Tuesday | Morning | 10 | |
| | December | Tuesday | Afternoon | 20 | |
| | December | Wednesday | Morning | 30 | |
| | December | Wednesday | Afternoon | 40 | |
| | December | Thursday | Morning | 50 | |
| | December | Thursday | Afternoon | 10 | |
| | December | Friday | Morning | 20 | |

Can you summarize my life based on this data in interpretable way?

---

## ③ State of the Art

# always predict default label
# specific rules for exception

# Capture only the local behaviour in the data

# Capture only the local behaviour in the data

**Rule-based classification**
- set of rules that **predicts class of examples well**
- **CN2, RIPPER, AQ, C4.5,SBRL**

**Subgroup Discovery**
- set of rule that **describes subgroup of examples well**
- **Cortana, Vikamine**

**Probabilistic Rule List**
- set of rule that **describes all examples well, being small**
- **PRL**

Timeline: 1986 1989 1993 1995 2011 2012 2016 2018

AQ · CN2 · C4.5 · RIPPER · Cortana · Vikamine · SBRL · ? **PRL** — Rule-based methods

---

## ④ Simple Example

| | a | b | c | e | M/F |
|---|---|---|---|---|---|
| 1 | ✓ | ✓ | ✓ | ✓ | M |
| 2 | ✓ | ✓ | ✓ | ✓ | F |
| 3 | ✓ | ✓ | ✓ | ✓ | M |
| 4 | | ✓ | ✓ | ✓ | F |
| 5 | ✓ | | ✓ | | F |

How is the purchase of fruits depended on who buys them?

Summarize the target attribute based on what people buy in interpretable way

## 5 — Problem of PRL

▷ **Given** A database of instances (observations), with for each a *Boolean target attribute*

▷ **Find** A Probabilistic Rule List

▷ **Such that** this Rule List when applied to the given database *describe* it well being *small* and *interpretable*

---

## 6 — Contributions

CONTRIBUTION

**Goal:** Finding Rule List which learns rules with probabilities to characterize the class distribution over the entire data and favor smaller rule lists to ease interpretation

☑ New optimization criterion
   ▷ based on the MDL principle;
   ▷ aiming to find small-and-good rule lists

☑ New search algorithm
   ▷ based on branch-and-bound search;
   ▷ aiming to find the global optimum

---

## 7 — Methodology (1) Pattern-Based Approach

$$Pr = \frac{\text{Number of } \male \text{ in the rule cover}}{\text{The size of the rule cover}}$$

| | a | b | c | e | M/F |
|---|---|---|---|---|---|
| 1 | ✓ | ✓ | ✓ | ✓ | M |
| 2 | ✓ | ✓ | ✓ | | F |
| 3 | ✓ | ✓ | | ✓ | M |
| 4 | | ✓ | ✓ | ✓ | F |
| 5 | ✓ | | ✓ | | F |

IF { 🍎 🍏 }   $p = 2/3$
ELSE IF { 🍐 }   $p = 0/2$

**Standard Itemset Mining Task (pattern-sequence)**

---

## 8 — Methodology (2) Minimum Description Length

Rule List   Itemset List   Database

$$argmin_R \quad \text{score}(R, F, D)$$

• Score is (based on Shannon's Noiseless Channel Coding Theorem)

   • the number of bit to use each rule to encode the data (using log of probabilities)

   • the number of bit to encode the rule itself

## Slide 9

**Methodology (2)**
**Minimum Description Length**

$$\mathcal{R}^* = \operatorname*{argmin}_{\mathcal{R} \in \mathcal{L}(\mathcal{F}^* \cdot \emptyset)} L_{data}(\mathcal{D}|\mathcal{R}) + L_{model}(\mathcal{R})$$

$$L_{data}(\mathcal{D}|\mathcal{R}) = -\sum_{j=1}^{k} L_{local\ data}(\mathcal{D}|I^{(j)});$$

$$L_{model}(\mathcal{R}) = \log n + \sum_{j=1}^{} \Big( \log m + m_j \log m + \log n \Big)$$

$$|I^{(1)}| \qquad I_1^{(1)} \ldots I_{|I^{(1)}|}^{(1)} \qquad n_1^+$$

## Slide 10

**Methodology (2)**
**Why MDL is interressing**

| Rule-list | | Data | Model | Total |
|---|---|---|---|---|
| IF {} | $p = 2/5$ | 1243 | **6** | 1249 |

256x

## Slide 11

# Greedy solution

256x

|   | a | b | c | e | M/F |
|---|---|---|---|---|---|
| 1 | ✓ | ✓ | ✓ | ✓ | M |
| 2 | ✓ | ✓ | ✓ | ✓ | F … |
| 3 | ✓ | ✓ |  | ✓ | M |
| 4 |  | ✓ | ✓ | ✓ | F |
| 5 | ✓ |  | ✓ |  | F |

256x

**Select the best ⟨Itemset+default⟩ PRL**

until there is no more transaction to cover

⊙ { 🍎 🥝 }
$p = 2/3$          $p = 0/2$

**PRL with 722bits/1280 (56% of compression)**

## Slide 12

# Greedy solution

256x

|   | a | b | c | e | M/F |
|---|---|---|---|---|---|
| 1 | ✓ | ✓ | ✓ | ✓ | M |
| 2 | ✓ | ✓ | ✓ | ✓ | F |
| 3 | ✓ | ✓ |  | ✓ | M |
| 4 |  | ✓ | ✓ | ✓ | F |
| 5 | ✓ |  | ✓ |  | F |

Can We find Better?

$\langle (\phi, \frac{2}{5}) \rangle$
1243bits

Greedy Solution

$\langle (\{A\}, \frac{2}{4}), (\phi, \frac{0}{1}) \rangle$
1039 + 0bits

$\langle (\{C\}, \frac{1}{4}), (\phi, \frac{1}{1}) \rangle$
846 + 0bits

$\langle (\{A, B, C\}, \frac{1}{2}), (\phi, \frac{1}{3}) \rangle$
531 + 706bits

$\langle (\{B\}, \frac{2}{5}), (\phi, \frac{0}{1}) \rangle$
1039 + 0bits

$\langle (\{B, C\}, \frac{1}{3}), (\phi, \frac{1}{2}) \rangle$
722 + 512bits

$\langle (\{A, B\}, \frac{2}{3}), (\phi, \frac{0}{2}) \rangle$
722 + 0bits

$\langle (\{A, C\}, \frac{1}{3}), (\phi, \frac{1}{2}) \rangle$
722 + 512bits

Optimal Solution

$\langle (\{A, B, C\}, \frac{1}{2}), (\{A\}, \frac{1}{2}), (\phi, \frac{0}{1}) \rangle$
531 + 527 + 0bits

$\langle (\{A, B, C\}, \frac{1}{2}), (\{C\}, \frac{0}{2}), (\phi, \frac{1}{1}) \rangle$
531 + 15 + 0bits

$\langle (\{A, B, C\}, \frac{1}{2}), (\{E\}, \frac{1}{2}), (\phi, \frac{0}{1}) \rangle$
531 + 527 + 0bits

**Greedy Solution**

**Optimum Solution**

# Branch-and-Bound

**Algorithm 2:** *Branch-and-bound* $(\mathcal{F}, \mathcal{D})$

1   $PQ$ : PriorityQueue   ▷ *Partial rule lists ordered by code-length when adding default rule*
2   $best\mathcal{R} \leftarrow \langle\emptyset\rangle$, $best \leftarrow L(best\mathcal{R})$
3   $PQ$.enqueue-with-priority$(\langle\rangle, L(\langle\emptyset\rangle))$
4   **while** $\mathcal{R} \leftarrow PQ.dequeue()$ **do**
5     **for** *each* $I \in \mathcal{F} \setminus \mathcal{R}$ **do**
6       $\mathcal{R}' \leftarrow \langle\mathcal{R}, I\rangle$
7       **if** $L(\langle\mathcal{R}', \emptyset\rangle) < best$ **then**
8         $best\mathcal{R} = \langle\mathcal{R}', \emptyset\rangle$, $best \leftarrow L(best\mathcal{R})$
9       **if** *lower-bound*$(\mathcal{R}') < best$ **then**
10        $PQ$.enqueue-with-priority$(\mathcal{R}', L(\langle\mathcal{R}', \emptyset\rangle))$
11   **return** $best\mathcal{R}$

> Start by default rule list

> Add iteratively new rule in the rule List

> update the best if the new rule + default has minimum length than the curent best

> Compute the bound and store expandable rule-list in the PQ

*29/10/2018 — Aoga et al. — PRL — Discovery Science*
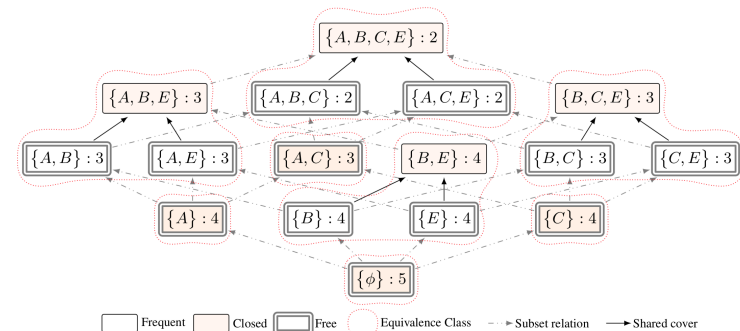
---

# Lower bound computation

- **When we have a partial rule list, can we remove some possibilities?**

- A good lower-bound is difficult to compute since there is an exponential number of rules that can be added to the list

- In the perfect case,

  - any expansion has to be greater than or equal in size to 1,

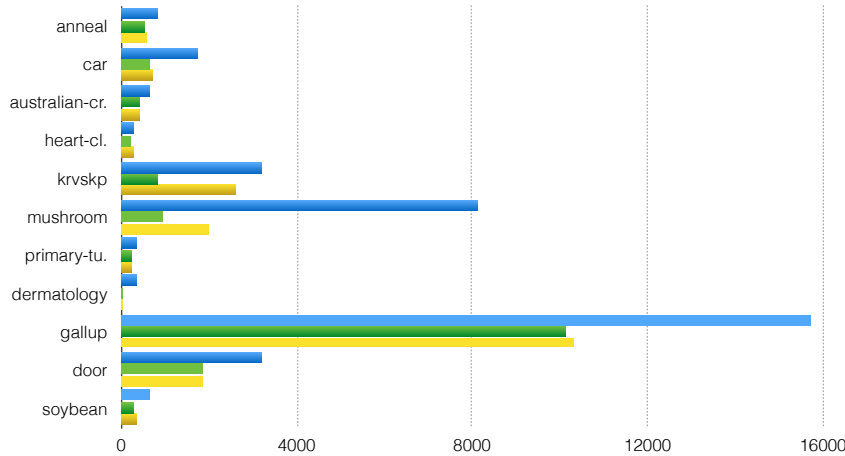  - and any expansion will achieve at best a data compression of 0

*29/10/2018 — Aoga et al. — PRL — Discovery Science*

---



**EXPERIMENTS**

I'm cleverer than Obelix

I'm stronger than Asterix

*29/10/2018 — Aoga et al. — PRL — Discovery Science*

---

# Implementation Details

- Set representation as a Bitvector + Bitwise operation
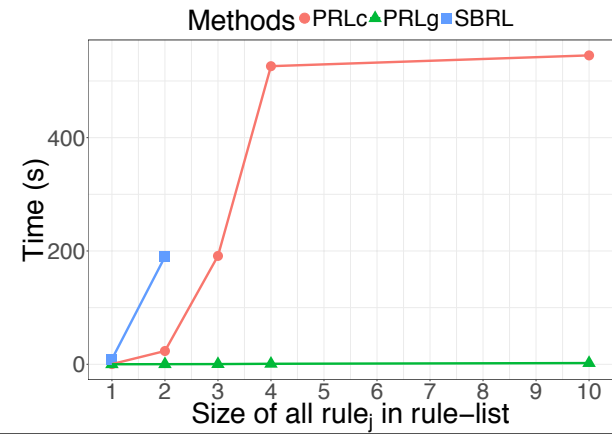
- Only find Rule in Free-sets (equivalent classes)



Frequent   Closed   Free   Equivalence Class   Subset relation   Shared cover

*29/10/2018 — Aoga et al. — PRL — Discovery Science*

**17 Compression Ratio**

Original Size ■ B&B Compression Size ■ Greedy Compression Size ■



**18 Impact of parameters**

Mushroom dataset (size = 8124x112) :: Varying rule list size (+ time limit=600)



**19 PRL vs Rule Learning algorithms**

Coding-length of: data ■ model ■

The best likelihood

The smallest size



**20 PRL vs Rule Learning algorithms**

Aoga et al. — PRL — Discovery Science

29/10/2018

# Prediction power of PRL

Gallup

# Conclusion

- We propose a New Descriptive method called Probabilistic Rule List

  - This Rule List is designed to be small and characterize well the target data

  - We found using a new optimization criterion based on MDL principle

  - We also designed a branch-and-bound method using Best First Search strategy