

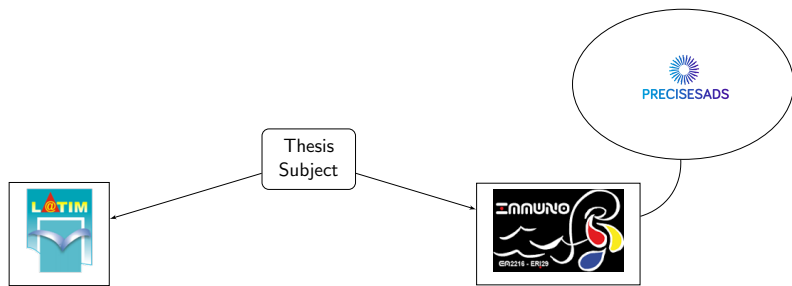
# Systemic simulation for the extraction of a patient's biological signature

Nathan Foulquier

LaTIM  
CHU Morvan, UMR 1227

Seminar, May 2017

# Context



Pascal Redou

Codirection

LaTIM

Lipab

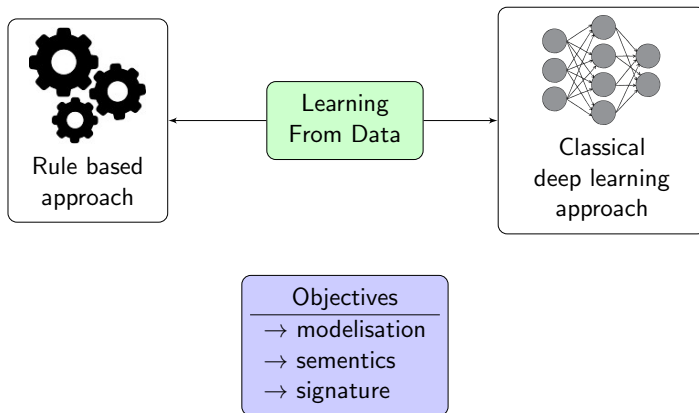
Alain Saraux

Christophe Le Gall

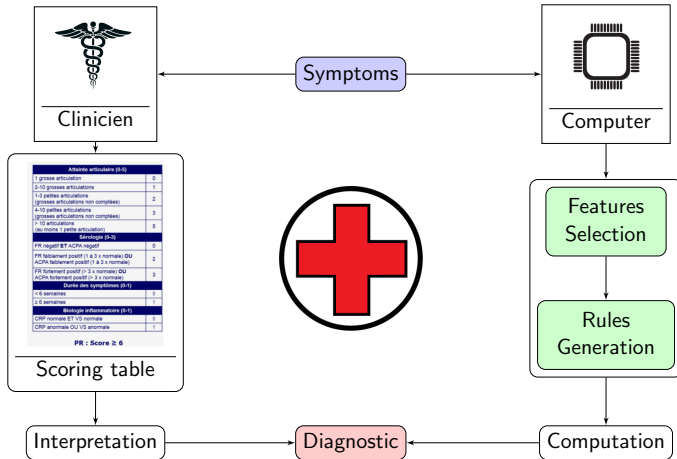
supervisors

Laurent Gobert

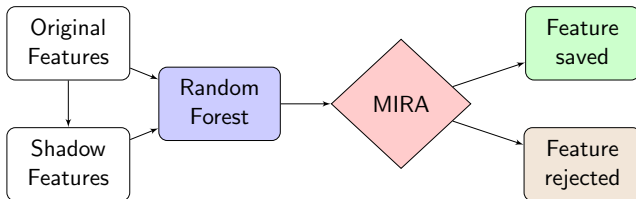
# The main idea



# The diagnostic



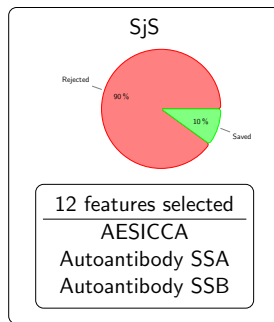
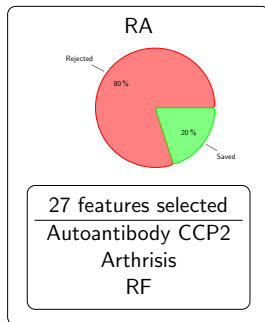
# The Boruta Algorithm



Shadow Feature: random permutation of an original feature  
MIRA: maximum importance of all shadow features

- Rank the original features
- Select the "important" features

# Results

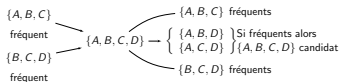


→ 80% of irrelevant features for RA  
→ 90% of irrelevant features for SjS

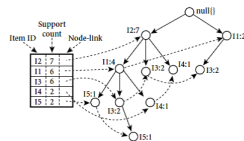
# Pattern Mining

$$\text{support}(E) = \frac{\text{card}(\{p \in P | E \subset p\})}{\text{card}(P)}$$

Pattern E is frequent if  $\text{support}(E) \geq \text{minsup}$



A priori Generation



FP-tree Mining

# Rules Extraction

$$r : (e - h) \rightarrow h$$

$$e \in \{items\}, card(e) \geq 2$$

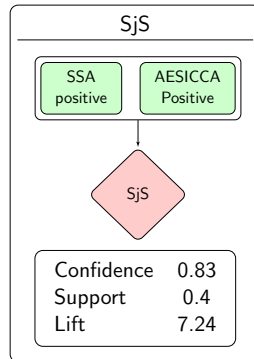
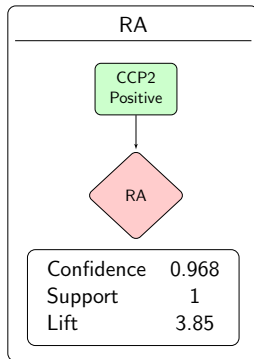
Generate  $h$  with  $h \neq \emptyset, h \neq e$

$$confidence(r) = \frac{support(e)}{support(e - h)}$$

a rule  $r$  is valid if  $confidence(r) \geq minconf$

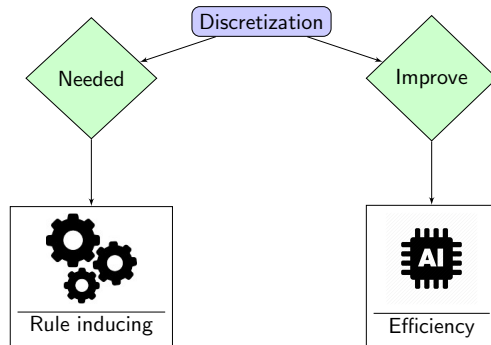


# Results



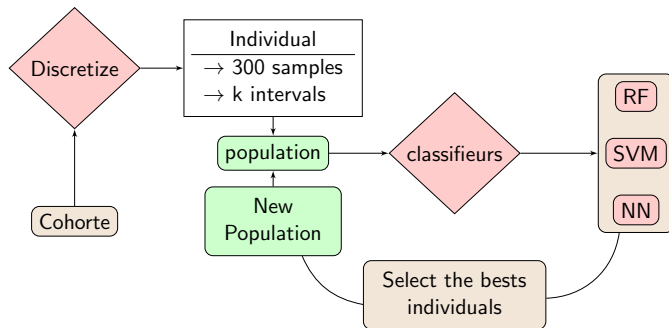
→ Easiest to identify RA than SjS

# Why?



- Some algorithms can handle discrete attributes only
- correct skewed distribution
- reduce the influence of outliers

# First Attempt



⇒ Learn the optimal value of  $k$   
Where  $k \in \{2, \dots, \max_k\}$

# The Ameva Algorithm

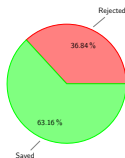
$$Ameva(k) = \frac{\chi^2(k)}{k(l-1)} \quad \chi^2(k) = N(-1 + \sum_{i=1}^l \sum_{j=1}^k \frac{n_{ij}^2}{n_i n_j})$$

$k$ : number of discrete intervals,  $l$ : number of classes

Maximize the dependency relationship  
between the class labels and the continuous-values attribute

Minimize the number of intervals  $k$

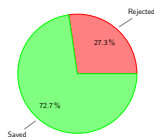
# Flow cytometry



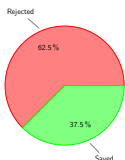
P1, 12 variables saved



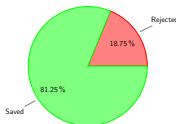
P2, 4 variables saved



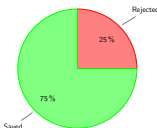
P3, 8 variables saved



P4, 3 variables saved

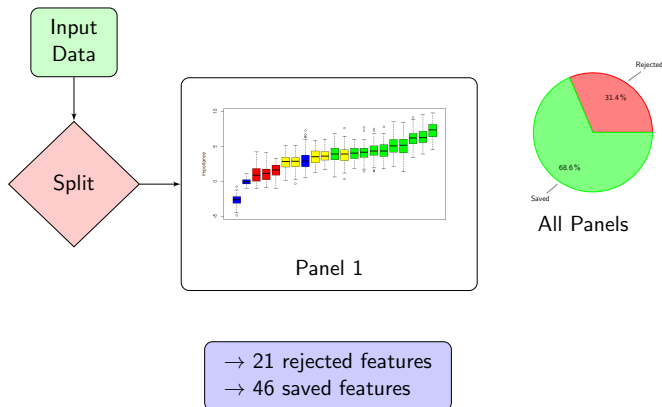


P5, 13 variables saved

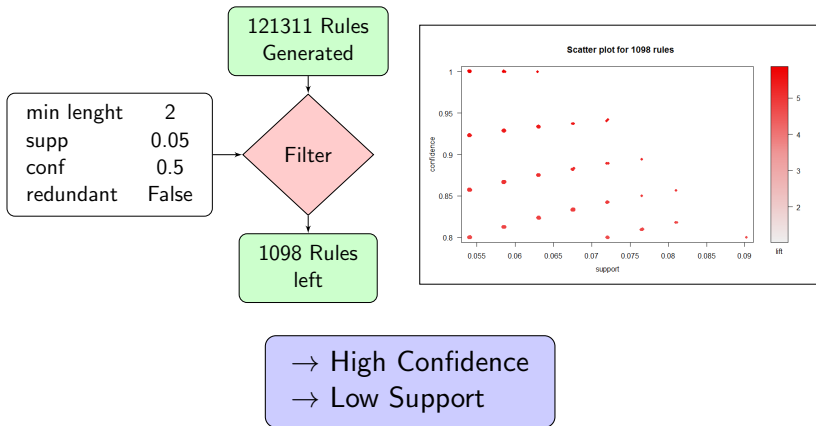


P6, 6 variables saved

# Flow cytometry



# Flow Cytometry

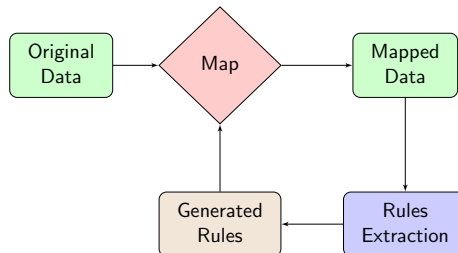


# Flow Cytometry





## Improvements



## Rule based discrimination

$$X_i \rightarrow [0, 1]$$

Probabilistic approach

# Perspectives

## Perspectives

- Refine the rules
- Implement the inference engine
- More data
- Include Inception

