

# Projet de conception et de prog...

Mickael LE DENMAT

Université Versailles Saint-Quentin en Yvelines  
Investigating feature selection techniques to improve data mining tasks

27 janvier 2023

# Table des matières

1 Introduction

2 *Rough Set Theory*

3 Références

- Le monde d'aujourd'hui : Beaucoup de données !

- Le monde d'aujourd'hui : Beaucoup de données !
- Problème de prise de décision.

- Le monde d'aujourd'hui : Beaucoup de données !
- Problème de prise de décision.

## Solution

*Data Mining* : la pratique consistant à rechercher automatiquement de grandes quantités de données afin de découvrir des tendances et des modèles qui vont au delà de la simple analyse. [Oracle, ]

- Complexité du monde réel.

- Complexité du monde réel.
- Accélérer la prise de décision.

- Complexité du monde réel.
- Accélérer la prise de décision.

## Solution

*Feature selection*  $\Rightarrow$  Rough Set Theory.



# *Rough Set Theory*

Patient	Headache	Muscle-Pain	Temperature
$o_1$	Yes	Yes	Very High
$o_2$	Yes	No	High
$o_3$	Yes	No	High
$o_4$	No	Yes	Normal
$o_5$	No	Yes	High
$o_6$	No	Yes	Very High

Table – Exemple d'un système d'information

# Un système de décision

Patient	Headache	Muscle-Pain	Temperature	Flu
$o_1$	Yes	Yes	Very High	Yes
$o_2$	Yes	No	High	Yes
$o_3$	Yes	No	High	No
$o_4$	No	Yes	Normal	No
$o_5$	No	Yes	High	Yes
$o_6$	No	Yes	Very High	Yes

Table – Exemple d'un système d'information

## Définition

Soit  $I = (U, A)$ , avec  $I$  un système d'information,  $U$  un ensemble d'objets et  $A$  un ensemble d'attributs. Avec n'importe quel sous ensemble  $P \subseteq A$ . Il existe une relation d'équivalence, noté  $IND(P)$ , définit comme :

$$IND(P) = \{(x, y) \in U^2 \mid \forall a \in P, a(x) = a(y)\} \quad (1)$$

## Exemple

Patient	Headache	Muscle-Pain	Temperature	Flu
$o_1$	Yes	Yes	Very High	Yes
$o_2$	Yes	No	High	Yes
$o_3$	Yes	No	High	No
$o_4$	No	Yes	Normal	No
$o_5$	No	Yes	High	Yes
$o_6$	No	Yes	Very High	Yes

Soit un ensemble  $X$  d'objets cibles tel que  $X \subseteq U$ .

- Représenter  $X$  avec  $P$  (un sous ensemble d'attributs).

Soit un ensemble  $X$  d'objets cibles tel que  $X \subseteq U$ .

- Représenter  $X$  avec  $P$  (un sous ensemble d'attributs).
- $X$  n'a qu'une seule classe  $\Rightarrow$  classe d'équivalence de  $P$ .

Soit un ensemble  $X$  d'objets cibles tel que  $X \subseteq U$ .

- Représenter  $X$  avec  $P$  (un sous ensemble d'attributs).
- $X$  n'a qu'une seule classe  $\Rightarrow$  classe d'équivalence de  $P$ .

## Problème

On ne peut pas calculer précisément  $X$ .

On doit l'approximer.



## *B-Lower Approximation*

$$\underline{B}(X) = \{o_j \mid [o_j]_B \subseteq X\} \quad (2)$$

C'est l'ensemble complet des objets dans  $U/P$  qui peuvent être classés dans  $X$  sans ambiguïté.

## *B-Lower Approximation*

$$\underline{B}(X) = \{o_j | [o_j]_B \subseteq X\} \quad (3)$$

C'est l'ensemble complet des objets dans  $U/P$  qui peuvent être classés dans  $X$  sans ambiguïté.

## *B-Upper Approximation*

$$\bar{B}(X) = \{o_j | [o_j]_B \cap X \neq \emptyset\} \quad (4)$$

C'est l'ensemble des objets dans  $U/P$  qui peuvent être classés dans  $X$ . Ce sont des objets possédant un objet indiscernable.

## *B-boundary Approximation*

$$BN_B(X) = \bar{B}(X) - \underline{B}(X) \quad (5)$$

## Définition

### Positive Region

$$POS_C\{(d)\} = \bigcup_{X \in U/\{(d)\}} \bar{C}(X) \quad (6)$$

# Positive & Negative Region

## Définition

### Positive Region

$$POS_C\{(d)\} = \bigcup_{X \in U/\{(d)\}} \bar{C}(X) \quad (7)$$

### Negative Region

$$NEG_C\{(d)\} = U - \bigcup_{X \in U/\{(d)\}} \underline{C}(X) \quad (8)$$

## Définition

Un *reduct* est un sous ensemble minimal d'attributs ayant la même *Positive Region* que l'ensemble des attributs.

## Définition

Un *reduct* est un sous ensemble minimal d'attributs ayant la même *Positive Region* que l'ensemble des attributs.

## Définition

Un *core* est un ensemble d'attributs indépendant incluant tous les *reduct*.



Oracle.

Qu'est-ce que le data mining ?



Merci pour votre attention !  
Des questions ?