



**UFR des Sciences**  
**CAMPUS DE VERSAILLES**

UNIVERSITÉ VERSAILLES SAINT-QUENTIN EN  
YVELINES

INVESTIGATING FEATURE SELECTION TECHNIQUES TO  
IMPROVE DATA MINING TASKS  
RAPPORT

---

## Projet de conception et de programmation

---

*Élève :*

Mickael LE DENMAT

*Enseignant :*

DR Zained GARCIA

27 janvier 2023

# Table des matières

<b>1</b>	<b>Introduction</b>	<b>3</b>
1.1	<i>Data Mining</i>	3
1.2	<i>Reduced Data</i>	4
1.3	Contexte	4
<b>2</b>	<b><i>Rough Set Theory - RST</i></b>	<b>5</b>
2.1	Un système d'information	5
2.2	Un système de décision	5
2.3	Classe d'équivalence	6
2.4	<i>Indiscernability Relation</i>	6
2.5	<i>rough set</i>	6
2.6	<i>B-lower approximation</i>	7
2.7	<i>B-upper approximation</i>	7
2.8	<i>B-boundary region</i>	7
2.9	<i>Positive Region</i>	7
2.10	<i>Negative Region</i>	7
2.11	<i>Reduct</i>	7
2.12	<i>Core</i>	8
<b>3</b>	<b>Algorithmes</b>	<b>9</b>
3.1	Pseudocode	9
3.1.1	<i>Indiscernability Relation</i>	9
3.1.2	<i>B-upper approximation</i> , <i>B-lower approximation</i> et <i>B-boundary region</i>	10
3.1.3	<i>Positive Region</i> et <i>Negative Region</i>	12
3.1.4	<i>Reduct</i> et <i>Core</i>	13
3.1.5	<i>quickReduct</i>	14
3.2	Implémentation	15
3.2.1	<i>Indiscernability Relation</i>	15
3.2.2	<i>B-lower approximation</i> , <i>B-upper approximation</i> et <i>B-boundary region</i>	17
3.2.3	<i>Positive Region</i> et <i>Negative Region</i>	19
3.2.4	<i>Reduct</i> et <i>Core</i>	21
3.2.5	<i>QuickReduct</i>	22
3.3	Test	23
3.3.1	<i>Indiscernability Relation</i>	23
3.3.2	<i>B-lower approximation</i> , <i>B-upper approximation</i> et <i>B-boundary region</i>	24
3.3.3	<i>Positive Region</i> et <i>Negative Region</i>	25
3.3.4	<i>Reduct</i> et <i>Core</i>	25
3.3.5	<i>QuickReduct</i>	26
<b>4</b>	<b>Applications</b>	<b>28</b>
4.1	Résumé des données	28
4.2	Analyse et visualisation des données	29
4.3	<i>Preprocessing</i>	33

4.4	Visualisation des données . . . . .	33
4.5	Application de modèles de classification . . . . .	33
4.5.1	Modèle A . . . . .	33
4.5.2	Modèle B . . . . .	33
4.5.3	Modèle C . . . . .	33
4.5.4	Modèle D . . . . .	33
4.6	Résultat . . . . .	33
4.6.1	Évaluation . . . . .	33
4.6.2	Comparaison . . . . .	33
<b>5</b>	<b>Conclusion</b>	<b>34</b>
5.1	Rappel des résultats attendues et obtenues . . . . .	34
5.2	Limites de l'algorithme . . . . .	34
5.3	Suggestions d'améliorations . . . . .	34
	<b>Bibliographie</b>	<b>35</b>

# 1 Introduction

## 1.1 *Data Mining*

La création des systèmes de stockages, notamment des bases de données a facilité la sauvegarde, l'organisation, et la recherche des données. Ils ont évolués avec les objets que nous voulions stocker, qui sont de plus en plus nombreux, complexes (types, structures, échelles, ...), et par conséquent plus lourd. Le stockage massif de ces données est motivé par de multiples raisons, par exemple, c'est simplement un stockage d'éléments personnels (photos, vidéos, textes, ...), ou alors des raisons comme le suivi d'évolutions (ventes de produit, ..) ou des changements (stockage de produits, ...) ou d'autres encore. Dans les raisons, les données peuvent être utilisées pour faire de la prise de décision. Cependant, dans ce cas précis, deux problèmes surviennent. Premièrement les informations représentées le sont sous formes de valeurs numériques, qui ne portent pas de contexte, Il faut donc l'extraire à partir des données brutes. Deuxièmement avec l'augmentation croissante des données stockées, l'analyse de patterns ou des contextes deviennent longue et complexe. Les conclusions peuvent être fausses car induites en erreurs du à une donnée inutile, ou elles peuvent être inexploitable dans la réalité. Il faut réduire l'ensemble des données. La question suivante s'est alors naturellement posée, comment comprendre ces données et les transformer en informations c'est-à-dire en savoir exploitable et utile. C'est dans ce contexte qu'est né le *Data Mining*. Le *Data Mining* est "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. [...]. Il est aussi connu sous le nom de découvert de connaissances dans la données." [Ora]. Pour découvrir ce savoir, il existe un processus, appelé *Knowledge Discovery in Database* décrivant les étapes successives à faire (cf figure ci-dessous).

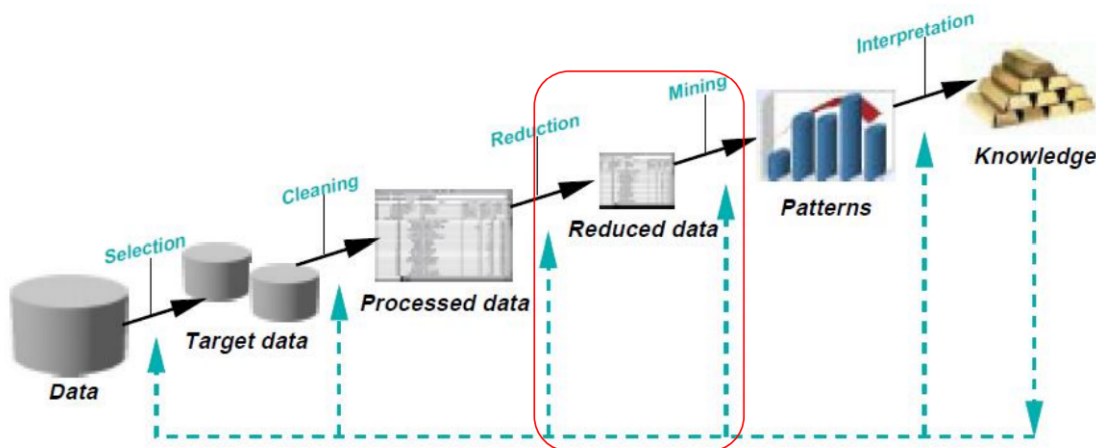


FIGURE 1 – KDD process

Dans ce papier nous nous concentrerons sur la partie *Reduced Data*

## 1.2 *Reduced Data*

La réduction de dimensionnalité (ou la sélection de variable) est une phase très importantes du KDD. Elle a plusieurs objectifs, réduire le nombre d'attributs afin d'accélérer l'apprentissage des modèles de machine learning sur les données, sélectionner les données les plus importantes mais aussi apporter plus de sémantique aux données. Pour faire de la réduction de donnée nous avons deux cas possibles, soit nous souhaitons garder toutes les variables mais nous les "simplifions" en les combinant. Cette méthode implique la perte de la sémantique des données. Soit, à l'inverse, nous prenons uniquement les variables les plus importantes et nous gardons la sémantique des données. La première s'appelle la *transformation based*. Nous pouvons citer comme exemple le *Principal Component Analysis - PCA*. Il va projeter les données sur des axes (les Composantes Principales) en minimisant la distance entre les points et leurs projections. Ainsi nous pouvons réduire notre problème de  $x$  variables à  $y$  principales composantes. Nous devons aussi assurer que les données garde une certaine sémantique, nous ne voulons pas transformer nos données n'importe comment. Dans le cas du *PCA*, nous calculons la variance cumulée et nous vérifions quelle soit entre 95% et 99%. Dans le cas de la prise de décision ne nous pouvons pas appliquer une telle méthode car le sens des données nous importe. La deuxième, la *Selection based*, utilise des algorithmes afin de connaître l'importance de chaque variable pour conserver les plus importantes. Ici nous parlerons de l'une des méthodes majeures : le *Rough Set theory - RST*.

## 1.3 Contexte

Ce travail est un projet réalisé par un étudiant en deuxième année de master à l'université de Versailles Saint-Quentin en Yvelines. Il est effectué dans le cadre d'une matière, "Projet conception et programmation". Il est supervisé et sera évalué par .... Les objectifs de ce travail sont multiples. Premièrement nous devons comprendre comment fonctionne le *QuickReduc*, un algorithme basé sur le RST. Deuxièmement nous devons l'implémenter et le tester. Troisièmement nous devons le comparer à d'autres algorithmes. Pour cela nous commencerons par expliquer ce qu'est le RST. Ensuite nous détaillerons son fonctionnement, nous l'implémenterons et nous le testerons sur des exemples simples. Puis nous appliquerons plusieurs algorithmes similaires sur plusieurs *datasets* et nous les comparerons. Enfin nous conclurons.

Commençons par expliquer ce qu'est le RST.

## 2 Rough Set Theory - RST

*Rough Set Theory - RST* ou la theorie des ensembles approximatifs est basé sur un ensemble de définition que nous allons détaillées dans ce chapitre.

### 2.1 Un système d'information

Commençons par la base : un système d'information. Un système d'information (*Information System - IS* en anglais) est une table ou une matrice avec une ligne par objets renseignés (l'ensemble des objets est appelé l'univers ( $U$ )) et  $x$  colonnes, appelé des attributs, décrivant les propriétés de l'objet (l'ensemble des attributs est appelé  $A$ ).

C'est données sont utilisées pour faire de la prise de dcision, il nous faut pour cela un attribut spécifique, que nous utiliserons comme décision. C'est pour cela que nous avons créés les systèmes d'informations.

### 2.2 Un système de décision

Un système de décision (*Decision System - DS* en anglais) est un système d'information avec une, comme dit précédemment, colonne supplémentaire appelée décision (l'ensemble des décision est appelé  $\theta$ ). Elle correspond à la classe de l'objet (exemple un patient peut avoir comme décision/ classe "malade" ou "non malade")

$x \in U$	a	b	c	d
$x_1$	1	0	2	2
$x_2$	0	1	1	1
$x_3$	2	0	0	1
$x_4$	1	1	0	2
$x_5$	1	0	2	0
$x_6$	2	2	0	1
$x_7$	2	1	1	1
$x_8$	0	1	1	0

TABLE 1 – Exemple d'un système d'information

$x \in U$	a	b	c	d	e
$x_1$	1	0	2	2	0
$x_2$	0	1	1	1	2
$x_3$	2	0	0	1	1
$x_4$	1	1	0	2	2
$x_5$	1	0	2	0	1
$x_6$	2	2	0	1	1
$x_7$	2	1	1	1	2
$x_8$	0	1	1	0	1

TABLE 2 – Exemple d'un système de décision

Dans un IS ou un DS, il est possible d'avoir de la redondance dans les données. Ces données peuvent avoir un, deux voir même tous leurs attributs égales. Nous pouvons créer des sous ensembles de données grâce à ce critère, c'est ce qu'on appelle : les classes d'équivalence.

## 2.3 Classe d'équivalence

Mathématiquement la classe d'équivalence est défini comme suit : Soit un ensemble E muni d'une relation d'équivalence, noté  $\sim$ . Ici deux objets sont équivalent si et seulement si ils possèdent les mêmes valeurs pour chaque attributs. On définit une classe d'équivalence, noté  $[x]$ , d'un élément  $x$  de E comme l'ensemble des  $y$  de E tels que  $x \sim y$ . [Wik]

$$y \in [x] \iff x \sim y \quad (1)$$

La redondance des données n'est pas forcément problématique dans le machine learning, sauf dans un cas particulier : si deux données possédant les mêmes attributs ont une valeur de décision différentes. Lorsque deux objets respectent cette règle, nous avons une *Indiscernability Relation* entre les deux.

## 2.4 Indiscernability Relation

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  $IND(P)$ , défini comme :

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

La relation  $IND(P)$  est la *Indiscernability Relation*.

Autrement dit deux objets sont indiscernable si pour toutes propriétés  $a \in P, a(x) = a(y)$ . Dans l'exemple TABLE 2, nous nous pouvons définir trois relations indiscernables sur les attributs  $\{a\}$ ,  $\{b, c\}$  et  $\{a, b, c\}$  qui définissent les trois portions de l'univers :

$$IND(a) = \{\{x_1, x_4, x_5\}, \{x_2, x_8\}, \{x_3, x_6, x_7\}\} \quad (3)$$

$$IND(b, c) = \{\{x_3\}, \{x_1, x_5\}, \{x_2, x_7, x_8\}, \{x_8\}\} \quad (4)$$

$$IND(a, b, c) = \{\{x_1, x_5\}, \{x_2, x_8\}, \{x_3\}, \{x_4\}, \{x_6\}, \{x_7\}\} \quad (5)$$

Si nous calculons la relation indiscernable sur l'ensemble des attributs du système, nous notons cela  $IND_C$  ou  $U/C$ .

## 2.5 rough set

Soit un ensemble X d'objets cibles tel que  $X \subseteq U$ . Nous souhaitons représenté X en utilisant un sous-ensemble d'attributs P. X comprend une seule classe et nous voulons décrire cette classe en à partir des classes d'équivalence des attributs dans P. En général, X ne peut pas être décrit précisément car il peut inclure ou exclure des objets qui ne sont pas distingués sur la base des attributs de P. Pour résoudre ce problème nous approchons X en le bornant avec la *B-lower approximation* notée  $\underline{P}(X)$ , et la *B-upper approximation* notée  $\bar{P}(X)$ .

## 2.6 *B-lower approximation*

La *B-lower approximation* est défini comme suit :

$$\underline{B}(X) = \{o_j | [o_j]_B \subseteq X\} \quad (6)$$

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

## 2.7 *B-upper approximation*

La *B-upper approximation* est défini comme suit :

$$\bar{B}(X) = \{o_j | [o_j]_B \cap X \neq \emptyset\} \quad (7)$$

C'est l'ensemble des objets dans  $U/P$  qui peuvent possiblement être classés dans  $X$ . Ce sont des objets possèdent un objets indiscernables.

Grâce à ces deux approximations nous pouvons calculer *B-boundary region* notée  $BN_B(X)$ .

## 2.8 *B-boundary region*

La *B-boundary region* est défini comme :

$$BN_B(X) = \bar{B}(X) - \underline{B}(X) \quad (8)$$

Une fois que nous avons calculer les différentes approximations pour chaque valeurs de décisions, nous pouvons calculer les régions.

## 2.9 *Positive Region*

La *Positive Region* est défini comme ceci :

$$POS_C\{(d)\} = \bigcup_{X \in U/\{(d)\}} \bar{C}(X) \quad (9)$$

## 2.10 *Negative Region*

La *Negative Region* est défini comme ceci :

$$NEG_C\{(d)\} = U - \bigcup_{X \in U/\{(d)\}} \underline{C}(X) \quad (10)$$

## 2.11 *Reduct*

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



## 2.12 *Core*

*Core* est un ensemble d'attributs indépendant incluant tous les *Reduct* .

Ceci conclut l'explication du *Rough Set Theory - RST* . A présent nous allons détailler le fonctionnement des algorithmes du *Rough Set Theory - RST* et nous parlerons d'une des heuristique, le *QuickReduct* .

## 3 Algorithmes

Dans ce chapitre nous expliquerons le fonctionnement des algorithmes du RST en commençant par leur pseudocode puis par leur implémentation. Nous parlerons par la suite d'une heuristique, le *QuickReduct*

### 3.1 Pseudocode

Concernant les algorithmes du RST, nous les avons nommé ainsi :

- IND et IND\_C
- b\_lower et b\_upper
- POS\_C et NEG\_C

Pour décrire les algorithmes nous utiliserons la notation suivante :  
(numéro de ligne) : explication.

Commençons par IND et IND\_C.

#### 3.1.1 Indiscernability Relation

Les deux algorithmes sont similaires, IND regroupe les objets possédants les mêmes valeurs pour les C attributs, alors que IND\_C regroupe les objets possédants les mêmes valeurs pour tous les attributs.

---

#### Algorithm 1: Algorithme IND

---

**Input:** DS : Le système de décision

C : La liste des attributs

d : le nom de la colonne de décision.

**Result:** IND : La liste des identifiants identiques

```

1  $ind \leftarrow \{\}$  ;
2  $IS \leftarrow DS.drop(d)$  ;
3  $groupe \leftarrow IS.groupby(C)$  ;
4 for  $groupe$  in  $groupes$  do
5    $ind.append(g[index])$  ;
6 end
```

---



---

#### Algorithm 2: Algorithme IND\_C

---

**Input:** DS : Le système de décision

d : le nom de la colonne de décision.

**Result:** IND\_C : La liste des identifiants identiques

```

1  $ind\_c \leftarrow \{\}$  ;
2  $IS \leftarrow DS.drop(d)$  ;
3  $groupe \leftarrow IS.groupby(IS.columns)$  ;
4 for  $groupe$  in  $groupes$  do
5    $ind\_c.append(g[index])$  ;
6 end
```

---

Le fonctionnement des algorithmes est le suivant :

- (1) : Nous déclarons un ensemble vide. Cet ensemble contiendra des sous-ensembles d'index d'objets possédant les mêmes caractéristiques.
- (2) : Nous gardons uniquement le système d'information en supprimant la colonne décision.
- (3) : Nous regroupons les objets en fonction d'une sous-ensemble d'attributs C ou en fonction de tous les attributs IS.colonne. Cela nous donne une liste de tuple composé de deux éléments. Le premier est un ensemble d'attributs, et le deuxième est un ensemble d'objets possédant les mêmes valeurs pour cet ensemble d'attributs.
- (4) : Pour chaque groupe, nous ajoutons dans l'ensemble initial l'ensemble des des objets en gardant uniquement leurs index.

Passons maintenant à la fonction calculant la *B-upper approximation* .

### 3.1.2 *B-upper approximation* , *B-lower approximation* et *B-boundary region*

---

#### Algorithm 3: Algorithme *b\_lower*

---

**Input:** DS : Le système de décision.

ind\_c : La liste des objets indiscernables.

d : le nom de la colonne de décision.

d\_value : la valeur de décision.

**Result:** Cxi : la *B-lower approximation*

```

1  X ← DS.groupby(d);
2  Xi ← X[d_value];
3  for index ← Xi.index do
4      idc_obj ← groupe_ind_c_obj(index) ;
5      possede_indiscernable ← Faux ;
6      if len(idc_obj) ≠ 0 then
7          for index_obj2 in idc_obj do
8              if DS[index][d] ≠ DS[index_obj2][d] then
9                  possede_indiscernable ← Vrai
10             end
11         end
12     end
13     if possede_indiscernable == Faux then
14         CXi.append(index)
15     end
16 end

```

---

Pour cet algorithme nous commençons par grouper les objets en fonction de leur valeurs de décision d. Nous aurons donc pour chaque valeur de d un groupe avec x indices d'objets. Puis nous gardons le groupe correspondant à la valeur de *d\_value*. Pour chaque index d'objet o nous l'ajoutons dans CXi sauf s'il est dans un groupe indiscernable avec un autre objet o' tel que  $d(o) \neq d(o')$ .

Expliquons le fonctionnement de l'algorithme :

- (1) : Nous regroupons les objets du système de décision en fonction de leurs valeurs de décision.
- (2) : Nous gardons uniquement le groupe dont la valeur est égal à  $d\_value$ .
- (3) : Nous parcourons les indexes des objets de chaque groupes.
- (4) : `groupe_ind_c_obj` est une fonction qui renvoie le groupe indiscernable où se situe l'objet `obj` en le supprimant du groupe. Ainsi si il est seul dans le groupe, la liste renvoyé sera vide, sinon nous aurons la liste des objets indiscernables.
- (6) : Si l'objet n'est pas seul dans son groupe.
- (7) : Nous parcourons les objets indiscernables.
- (8) : Si un des objets indiscernable possède une valeur de décision différente, nous ne devons pas les ajouter dans CXi.

Nous avons un algorithme similaire pour calculer la *B-lower approximation* .

---

**Algorithm 4:** Algorithme *b\_lower*

---

**Input:** DS : Le système de décision.

ind\_c : La liste des objets indiscernables.

d : le nom de la colonne de décision.

d\_value : la valeur de décision.

**Result:** CXi : la *B-lower approximation*

```

1  $X \leftarrow DS.groupby(d);$ 
2  $Xi \leftarrow X[d\_value];$ 
3 for  $index \leftarrow Xi.index$  do
4    $idc\_obj \leftarrow groupe\_ind\_c\_obj(index) ;$ 
5    $possede\_indiscernable \leftarrow Faux ;$ 
6    $sous\_liste \leftarrow \{ \} ;$ 
7   if  $len(idc\_obj) \neq 0$  then
8     for  $index\_obj2$  in  $idc\_obj$  do
9       if  $DS[index][d] \neq DS[index\_obj2][d]$  then
10          $sous\_liste.append(index\_obj2)$ 
11       end
12     end
13   end
14    $CXi.concat(sous\_liste)$ 
15 end
```

---

La différence entre ces deux algorithmes se trouve lignes 6, 10, 14 où nous créons une liste contenant un objet ainsi que tous ces objets indiscernables que nous ajoutons à CXi.

Avec ces deux algorithmes, nous pouvons calculer la *B-boundary region* .

---

**Algorithm 5:** Algorithme *b\_boundary*

---

**Input:** DS : Le système de décision.  
           bupper : La *B-upper approximation* .  
           blower : La *B-lower approximation* .

**Result:**  $BX_C$  : La *B-boundary region*

1  $BX_C \leftarrow bupper - blower$  ;

---

Nous passons aux fonctions *Positive Region* et *Negative Region* .

### 3.1.3 *Positive Region* et *Negative Region*

---

**Algorithm 6:** Algorithme *POS\_C*

---

**Input:** DS : Le système de décision.  
           d : le nom de la colonne de décision.

**Result:**  $POS_C\{(d)\}$  : La *Positive Region*

1  $d\_values \leftarrow DS[d]$  ;  
 2  $POS_C \leftarrow \{\}$  ;  
 3  $ind\_c \leftarrow IND\_C(DS, d)$  ;  
 4 **for**  $d\_value$  **in**  $d\_values$  **do**  
 5    $POS_C.concat(b\_lower(DS, ind\_c, d, d\_value))$   
 6 **end**

---

Le fonctionnement de l'algorithme est le suivant :

(1) : Nous commençons par prendre toutes les décisions (sans doublons) du système de décision.

(4) : Pour chaque valeur de décision nous calculons la *B-lower approximation* et nous l'ajoutons dans  $POS_C$ .

Nous avons aussi la version calculant la *Positive Region* pour C attributs.

---

**Algorithm 7:** Algorithme POS

---

**Input:** DS : Le système de décision.

          C : le nom des attributs. d : le nom de la colonne de décision.

**Result:**  $POS\{(d)\}$  : La *Positive Region*

1  $ds \leftarrow DS[C + d]$  ;  
 2  $ind \leftarrow IND(ds, d, C)$  ;  
 3  $d\_values \leftarrow ds[d]$  ;  
 4  $POS \leftarrow \{\}$  ;  
 5 **for**  $d\_value$  **in**  $d\_values$  **do**  
 6    $POS.concat(b\_lower(ds, ind, d, d\_value))$   
 7 **end**

---

L'algorithme pour la *Negative Region* est presque identique.

---

**Algorithm 8:** Algorithme  $NEG\_C$

---

**Input:** DS : Le système de décision.

ind\_c : La liste des objets indiscernables.

d : le nom de la colonne de décision.

**Result:**  $NEG_C\{d\}$  : La *Negative Region*

```

1 d_values ← DS[d] ;
2 NEG_C ← {} ;
3 for d_value in d_values do
4   | NEG_C.concat(b_lower(DS, ind_c, d, d_value))
5 end
6 NEG_C ← diff_list(DS.index, NEG_C)
```

---

La différence étant ligne (6), nous faisons la différences entre l'univers et la *Negative Region*.

A l'aide de la *Positive Region* nous pouvons calculer le *Reduct* et le *Core*.

### 3.1.4 Reduct et Core

L'algorithme du *Reduct* est le suivant :

---

**Algorithm 9:** Algorithme reduct

---

**Input:** DS : Le système de décision.

d : le nom de la colonne de décision.

**Result:** reducts : Le *Reduct*

```

1 reducts ← {} ;
2 POS_C ← POS_C(DS, d) ;
3 C ← DS.columns ≠ d ;
4 for combi in combinaisons(DS, D) do
5   | if combi ≠ C then
6     | pos ← POS(DS, d, combi) ;
7     | if pos == POS_C then
8       | reducts.append(combi)
9     | end
10  | end
11 end
```

---

Le détail de cet algorithme est le suivant :

- (1) : Nous créons un ensemble qui contiendra des sous-ensembles de *Reduct*.
- (2) : Nous calculons la *Positive Region* de l'ensemble des attributs.
- (3) : Nous prenons l'ensemble des attributs sans la colonne de décision.
- (4) : La fonction combinaisons renvoie l'ensemble des combinaisons possibles d'attributs du DS en paramètres. Nous parcourons donc chaque une des combinaisons.

- (5) : Si la combinaison est différente de l'ensemble d'attributs.
- (6) : Nous calculons la *Positive Region* de la combinaison d'attributs.
- (7) : Si la *Positive Region* est égal à la *Positive Region* de tous les attributs alors la combinaison d'attributs est un *Reduct*.
- (8) : Nous l'ajoutons dans l'ensemble des *Reduct*.

A l'aide des *Reduct* nous pouvons calculer le *Core*. Nous renvoyons la liste des attributs en communs dans tous les *Reduct*.

Nous avons terminé la partie *Rough Set Theory - RST*, nous allons passer à une heuristique le *QuickReduct*.

### 3.1.5 *quickReduct*

---

**Algorithm 10:** Algorithme quickReduct

---

**Input:** C : the set of all conditional features

D : the set of decision features

**Result:** R : the reduced set of conditional features

```

1  $R \leftarrow \{\}$ ;
2 do
3    $T \leftarrow R$ ;
4    $\forall x \in (C - R)$ ;
5   if  $\lambda_{RU\{x\}}(D) > \lambda_T(D)$  then
6      $T \leftarrow RU\{x\}$ ;
7   end
8    $R \leftarrow T$ ;
9 while  $\lambda_R(D) == \lambda_C(D)$ ;
```

---

Explication de l'algorithme :

- (1) Initialisation d'un ensemble, noté R, vide. Nous stockerons dedans l'ensemble réduit des attributs.
- (2) Tant que la dépendance de R n'est pas égal à la dépendance de C.
- (3) Nous créons un ensemble temporaire, noté T, égal à l'ensemble R.
- (4) Pour chaque attributs, noté x, présent dans C mais pas dans R.
- (5) Si la dépendance de R en ajoutant x est supérieur à la dépendance de T.
- (6) On ajoute x dans l'ensemble T.

Ceci terminant les pseudocodes, nous allons pouvoir les implémenter. Avant de parler à proprement parlé de l'implémentation, nous allons définir l'environnement de travail. Pour implémenter les algorithmes, nous utiliserons :

- Un ordinateur portable Dell G5 5587 avec :
  - Windows 10 Famille.
  - 16Go de RAM.
  - Intel Core i7-8750 avec une fréquence de 2.20GHz.
- Le logiciel Visual Studio Code pour coder.

- L’extension Jupyter v2022.11.1003412109 pour écrire notre notebook.
- Nous coderons avec le langage Python.
- Anaconda pour gérer l’environnement virtuel, avec un ensemble de bibliothèques comme par exemple **Panda** et **Numpy** pour la manipulation d’objets complexes. Mais aussi *sklearn* pour les modèles de machine learning.

Nous pouvons à présent parler de l’implémentation des différentes algorithmes.

## 3.2 Implémentation

Dans cette sous partie nous montrerons l’implémentation des algorithmes. Comme nous utilisons le Python, nous avons du convertir les objets mathématiques en objet implémenté par le langage Python. Par exemple les ensembles sont devenus des listes.

### 3.2.1 Indiscernibility Relation

```

1 def IND(DS, d, C):
2     """
3     Calcule la Indiscernibility Relation pour un attribut.
4
5     @param DS le système de décision.
6     @param d la décision.
7     @param C La liste des attributs.
8
9     @return ind la Indiscernibility Relation.
10    """
11    ind = []
12    IS = DS.drop(d, axis=1)
13    group = IS.groupby(C)
14
15    for g in group:
16        gDS = pd.DataFrame(g[1])
17        ind.append(list(gDS.index))
18    ind.sort()
19    return ind

```

Les paramètres sont DS un **DataFrame** issue de la bibliothèque **Panda** . Le paramètre d, une chaîne de caractères , est le nom de la colonne de décision. Le dernier paramètre C, une liste de chaîne de caractères , contenant le nom des attributs.

La fonctionnement de cet algorithme est le suivant :

(12) : La fonction **drop** est une fonction issue de la bibliothèque **Panda** . Nous l’utilisons pour supprimer la colonne d.

(13) : Nous utilisons la fonction **groupby** pour regrouper les objets.

(15) : Pour chaque groupe d’objets, sous forme de tuple.

(16) : Nous convertissons les objets pour utiliser la fonction **index**.

(17) : Nous ajoutons les tableaux convertis des indexes des objets en liste.



```

1 def IND_C(DS, d):
2     """
3     Calcule la Indiscernability Relation pour tous les
4     attributs du système de décision.
5
6     @param DS le système de décision.
7     @param d la décision.
8
9     @return ind_c la Indiscernability Relation.
10    """
11    ind = []
12    IS = DS.drop(d, axis=1)
13    group = IS.groupby(list(IS.columns))
14
15    for g in group:
16        gDS = pd.DataFrame(g[1])
17        ind.append(list(gDS.index))
18    ind.sort()
19    return ind

```

L'algorithme fonctionne de la même manière, sauf que nous groupons selon **IS.columns**, qui est une fonction de **Panda** nous renvoyant la liste des attributs de DS.

Passons maintenant à une fonction intermédiaires dont nous avons parlé durant la partie pseudocode.

```

1 def groupe_ind_c_obj(ind_c, obj):
2     """
3     Renvoie le groupe indiscernable de l'objet obj sans ce dernier.
4
5     @param ind_c la liste des objets indiscernable.
6     @param obj l'index d'un objet.
7
8     @return le groupe indiscernable de l'objet.
9     """
10    for g in ind_c:
11        for obj2 in g:
12            if obj == obj2:
13                copy_g = g.copy()
14                copy_g.remove(obj)
15                return copy_g
16    return []

```

L'idée de l'algorithme est la suivante :

Pour chaque sous liste dans *ind\_c*, nous parcourons la liste des indexes. Si nous trouvons l'objet que nous cherchons, nous copions son groupe et nous le supprimons de son groupe et nous renvoyons son groupe.

### 3.2.2 *B-lower approximation* , *B-upper approximation* et *B-boundary region*

```

1 def b_lower(DS, ind_c, d, d_value):
2     """
3     Renvoie la B-lower approximation.
4
5     @param DS le système de décision.
6     @param ind_c la liste des objets indiscernables.
7     @param d le nom de la colonne de décision.
8     @param d_value la valeur de décision.
9
10    @return la B-lower approximation.
11    """
12    X = DS.groupby(d)
13    Xi = None
14    for name, groupe in X:
15        if name == d_value:
16            Xi = pd.DataFrame(groupe)
17
18    if Xi is not None:
19        CXi = []
20
21        for index in Xi.index:
22            # Nous regardons le groupe indiscernable de cet objet.
23            idc_obj = groupe_ind_c_obj(ind_c, index)
24
25            # Si aucun des objets indiscernables de cet objet n'a une
26            # décision différente.
27            if not any(DS.at[index_obj2, d] != DS.at[index, d] for
28                    index_obj2 in idc_obj):
29                CXi.append(index)
30            CXi.sort()
31            return CXi
32    else:
33        return "error"

```

Le fonctionnement de cet algorithme a été évoqué durant la partie pseudocode, mais il y a quelques différences. La première étant entre les lignes (12) et (16) qui permettent de trouver les groupes où se trouve la *d\_value*. La deuxième étant les lignes (30) et (31) qui permettent de sélectionner la valeur de décision en fonction de leur indexes.

```

1 def b_upper(DS, ind_c, d, d_value):
2     """
3     Renvoie la B-upper approximation.
4
5     @param DS le système de décision.
6     @param ind_c la liste des objets indiscernables.
7     @param d le nom de la colonne de décision.
8     @param d_value la valeur de décision.
9
10    @return la B-upper approximation.
11    """
12    X = DS.groupby(d)
13    Xi = None
14    for name, groupe in X:
15        if name == d_value:
16            Xi = pd.DataFrame(groupe)
17
18    if Xi is not None:
19        CXi = list(Xi.index)
20        for index in Xi.index:
21            # Nous regardons le groupe indiscernable de cet objet.
22            idc_obj = groupe_ind_c_obj(ind_c, index)
23
24            # Si il possède des objets indiscernables.
25            list_add = [index_obj2 for index_obj2 in idc_obj
26                        if DS.at[index_obj2, d] != DS.at[index, d]]
27            CXi += list_add
28        CXi.sort()
29        return CXi
30    else:
31        return "error"

```

```

1 def b_boundary(DS, bupper, blower):
2     """
3     Calcule la B-boundary region.
4
5     @param DS le système de décision.
6     @param bupper la B-upper approximation.
7     @param blower la B-lower approximation.
8
9     @return la différence entre les deux approximations.
10    """
11    return diff_list(bupper, blower)

```

Cette fonction utilise une fonction que nous avons conçu, la fonction *diff\_list*.

```

1 def diff_list(list1, list2):
2     """
3     Renvoie la différence entre deux listes.
4
5     @param list1 une première liste.
6     @param list2 une deuxième liste.
7
8     @return les éléments présents dans list1 mais pas dans list2.
9     """
10    set_list2 = set(list2)
11    diff = [x for x in list1 if x not in set_list2]
12    return diff

```

Fonction issue du cite [Len22].

### 3.2.3 Positive Region et Negative Region

```

1 def POS(DS, d, C):
2     """
3     Calcul la Positive Region avec les C attributs.
4
5     @param DS Le système de décision.
6     @param d La colonne de décision.
7     @param C La liste des attributs.
8
9     @return La Positive region des C attributs.
10    """
11    attr = C.copy()
12    attr.append(d)
13    ds = DS[attr]
14    ind = IND(ds, d, C)
15    d_values = list(ds[d])
16    d_values = [*set(d_values)]
17    POS = []
18    for d_value in d_values:
19        POS += b_lower(ds, ind, d, d_value)
20    POS.sort()
21    return POS

```

(11) à (13) : Création d'un **DataFrame** avec les C attributs qui nous intéresses.

(16) : Permet de supprimer les doublons.

(17) à (19) : Nous ajoutons toutes les *B-lower approximation* en fonction de la valeur de décision *d\_value*.

```

1 def POS_C(DS, d):
2     """
3     Calcule la Positive Region avec tous les attributs.
4
5     @param DS Le système de décision.
6     @param d La colonne de décision.
7
8     @return La Positive region de tous les attributs.
9     """
10    d_values = list(DS[d])
11    d_values = [*set(d_values)]
12    POS_C = []
13    ind_c = IND_C(DS, d)
14    for d_value in d_values:
15        POS_C += b_lower(DS, ind_c, d, d_value)
16    POS_C.sort()
17    return POS_C

1 def NEG_C(DS, d):
2     """
3     Calcule la Negative Region avec tous les attributs.
4
5     @param DS Le système de décision.
6     @param d La colonne de décision.
7
8     @return La Negative region de tous les attributs.
9     """
10    d_values = list(DS[d])
11    d_values = [*set(d_values)]
12    NEG_C = []
13    ind_c = IND_C(DS, d)
14    for d_value in d_values:
15        NEG_C += b_upper(DS, ind_c, d, d_value)
16    NEG_C = [*set(NEG_C)]
17    neg = diff_list(list(DS.index), NEG_C)
18    neg.sort()
19    return neg

```

(7) : Nous supprimons les doublons.

(8) : Nous appliquons la formule donnée dans le chapitre 1.

### 3.2.4 Reduct et Core

```

1 def reduct(DS, d):
2     """
3     Renvoie toutes les réductions d'attributs possibles.
4
5     @param DS Le système de décision.
6     @param d La colonne de décision.
7
8     @return reducts : La liste des réductions.
9     """
10    reducts = []
11    pos_c = POS_C(DS, d)
12    C = list(DS.columns)
13    C.remove(d)
14    for combi in combinaisons(DS, d):
15        liste_combi = list(combi)
16        if liste_combi != C:
17            pos = POS(DS, d, liste_combi)
18            if pos == pos_c:
19                reducts.append(liste_combi)
20    reducts.sort()
21
22    # minimal set
23    min_len = min([len(x) for x in reducts])
24    reducts = [x for x in reducts if len(x) == min_len]
25
26    return reducts

```

```

1 def core(reducts):
2     """
3     Calcul l'intersection de toutes les réductions.
4
5     @param reducts la liste des réductions.
6
7     @return le core du système de décision.
8     """
9     return list(reduce(lambda i, j: i & j, (set(x) for x in reducts)))

```

Fonction issue du cite [man19].

### 3.2.5 QuickReduct

```

1 def dependance_attributs(DS, C, d):
2     """
3     Calcul la dépendance des attributs.
4
5     @param DS Le système de décision.
6     @param C la liste des attributs.
7     @param d la colonne de décision.
8
9     @return la dépendance des attributs.
10    """
11    ds = DS[C]
12    if len(list(ds.columns)) == 1 and list(ds.columns)[0] == d:
13        return 0
14    ind = IND_C(ds, d)
15    pos_c = POS_C(ds, ind, d)
16    dep = Fraction(len(pos_c), len(ds.index))
17    return dep

1 def quickReduct(DS, d):
2     """
3     L'algorithme de quickReduct.
4
5     @param DS Le système de décision.
6     @param d la colonne de décision.
7
8     @return une sous liste d'attributs.
9    """
10    C = list(DS.columns)
11    C.remove(d)
12    print("C :", C)
13    # calcul de lambda C.
14    dep_C = my_lambda(DS, DS.columns, d)
15    print("dep_C :", dep_C)
16    R = []
17    while True:
18        T = R
19        # calcul de lambda T.
20        attr = T.copy() + [d]
21        dep_T = my_lambda(DS, attr, d)
22        print("lambda({}) = {}".format(attr[:-1], dep_T))
23
24        # calcul de C - R.
25        Rset = set(R)
26        C_R = [x for x in C if x not in Rset]
27
28        # parcours des attributs.
29        for x in C_R:
30            # calcul de lambda RU{x}.
31            attr2 = R.copy() + [x, d]
32            dep_RUx = my_lambda(DS, attr2, d)
33            print("lambda({}) = {}".format(attr2[:-1], dep_RUx))
34
35            # changement meilleur dépendance.
36            if dep_RUx > dep_T:
37                T = R.copy() + [x]

```

```

38     print("changement R =", T)
39     attr3 = T.copy() + [d]
40     dep_T = my_lambda(DS, attr3, d)
41
42     R = T.copy()
43     # calcul de lambda R.
44     attr4 = R.copy() + [d]
45     dep_R = my_lambda(DS, attr4, d)
46     print("lambda({}) = {}".format(attr4[:-1], dep_R))
47     if dep_R == dep_C:
48         break
49     return R

```

### 3.3 Test

Passons maintenant aux tests. Pour cela nous devons créer des jeux de tests. Voici un exemple que nous utiliserons pour les algorithmes issues du *Rough Set Theory - RST* :

	Headache	Muscle-pain	Temperature	Flu
Patients				
o1	Yes	Yes	very high	Yes
o2	Yes	No	high	Yes
o3	Yes	No	high	No
o4	No	Yes	normal	No
o5	No	Yes	high	Yes
o6	No	Yes	very high	Yes

FIGURE 2 – Dataframe "df"

#### 3.3.1 Indiscernability Relation

Commençons par calculer IND.

Si nous prenons uniquement la colonne **Headache**, nous voyons les objets  $o_1$ ,  $o_2$  et  $o_3$  ont la même valeur, "Yes", et les  $o_4$ ,  $o_5$  et  $o_6$  ont la même valeur, "No". Nous avons donc le résultat suivant :

$$IND_{\{Headache\}} = \{\{o_1, o_2, o_3\}, \{o_4, o_5, o_6\}\}$$

Mais si nous prenons les colonnes **Headache** et **Muscle-pain**, nous voyons que l'objet est le seul à posséder le couple de valeur ("Yes", "Yes"), les objets  $o_2$  et  $o_3$  possèdent le couple de valeur ("Yes", "No"), enfin les objets  $o_4$ ,  $o_5$  et  $o_6$  possèdent le couple de valeur ("No", "Yes"). Nous avons donc le résultat suivant :

$$IND_{\{Headache, Muscle-pain\}} = \{\{o_1\}\{o_2, o_3\}, \{o_4, o_5, o_6\}\}$$

```

1 ind = IND(df, "Flu", ["Headache"])
2 ind
3 Python> [['o1', 'o2', 'o3'], ['o4', 'o5', 'o6']]
4 ind2 = IND(df, "Flu", ["Headache", "Muscle-pain"])
5 ind2
6 Python> [['o1'], ['o2', 'o3'], ['o4', 'o5', 'o6']]

```



Nous notons la sortie des cellules Jupiter à la suite de "Python>".

Calculons maintenant  $IND_C$ .

Dans le premier exemple nous voyons que  $o_1$ ,  $o_4$ ,  $o_5$ , et  $o_6$  ne sont identiques à aucun objets et que  $o_2$  et  $o_3$  sont identiques entre eux. Nous devons obtenir l'ensemble suivant :  $IND_C = \{\{o_1\}, \{o_2, o_3\}, \{o_4\}, \{o_5\}, \{o_6\}\}$ .

Dans le deuxième exemple nous voyons que  $x_1$ ,  $x_4$ ,  $x_5$  et  $x_7$  ne sont identiques à aucun autres objets,  $x_2$  et  $x_8$  sont identiques entre eux et  $x_3$  et  $x_6$  de même. Nous devons donc obtenir l'ensemble suivant :

$IND_C = \{\{x_1\}, \{x_2, x_8\}, \{x_3, x_6\}, \{x_4\}, \{x_5\}, \{x_7\}\}$ .

Avec notre algorithme, nous obtenons les résultats suivants :

```
1 ind_c1 = IND_C(df, "Flu")
2 ind_c1
3 Python> [['o1'], ['o2', 'o3'], ['o4'], ['o5'], ['o6']]
4 ind_c2 = IND_C(df2, "e")
5 ind_c2
6 Python> [['x1'], ['x2', 'x8'], ['x3', 'x6'], ['x4'],
7 ['x5'], ['x7']]
```

### 3.3.2 $B$ -lower approximation , $B$ -upper approximation et $B$ -boundary region

Nous avons deux valeurs pour la décisions, "Yes" ou "No". Cela implique que nous avons deux classes d'équivalences :

$X_1 = \{o_j | Flu(o_j) = \{Yes\}\} = \{o_1, o_2, o_5, o_6\}$  et  $X_2 = \{o_j | Flu(o_j) = \{No\}\} = \{o_3, o_4\}$

Si nous calculons les  $B$ -lower approximation , notées  $\underline{C}X_1$  et  $\underline{C}X_2$ , nous obtenons les résultats suivants :

$\underline{C}X_1 = \{o_1, o_5, o_6\}$  et  $\underline{C}X_2 = \{o_4\}$

Nous enlevons  $o_2$  dans  $\underline{C}X_1$  car  $o_2$  est dans le même groupe  $IND_C$  que  $o_3$ . Pour la même raison nous supprimons  $o_3$  dans l'équation  $\underline{C}X_2$ .

Si nous calculons les  $B$ -upper approximation , notées  $\bar{C}X_1$  et  $\bar{C}X_2$ , nous obtenons les résultats suivants :

$\bar{C}X_1 = \{o_1, o_2, o_3, o_5, o_6\}$  et  $\bar{C}X_2 = \{o_2, o_3, o_4\}$

Enfin nous pouvons calculer les  $B$ -boundary region .

$BN_C(X_1) = \{o_2, o_3\}$  et  $BN_C(X_2) = \{o_2, o_3\}$

Avons nos algorithmes nous avons les résultats suivants :

```
1 blower_yes = b_lower(df, ind_c, "Flu", "Yes")
2 Python> ['o1', 'o5', 'o6']
3 blower_no = b_lower(df, ind_c, "Flu", "No")
4 Python> ['o4']
5
6 bupper_yes = b_upper(df, ind_c, "Flu", "Yes")
7 Python> ['o1', 'o2', 'o3', 'o5', 'o6']
8 bupper_no = b_upper(df, ind_c, "Flu", "No")
9 Python> ['o2', 'o3', 'o4']
10
11 b_boundary(df, bupper_yes, blower_yes)
12 Python> ['o2', 'o3']
13 b_boundary(df, bupper_no, blower_no)
14 Python> ['o2', 'o3']
```

$blower\_yes$  correspond à  $\underline{C}X_1$ ,  $bupper\_yes$  correspond à  $\bar{C}X_1$ ,  $blower\_no$  correspond à  $\underline{C}X_2$ , et  $bupper\_no$  correspond à  $\bar{C}X_2$ .

### 3.3.3 Positive Region et Negative Region

Nous avons la *Positive Region* pour un ou plusieurs attributs :

$$\begin{aligned} POS_{\{Headache, Muscle-pain\}}(\{d\}) &= \{o_1\} \\ POS_{\{Headache, Temperature\}}(\{d\}) &= \{o_1, o_4, o_5, o_6\} \\ POS_{\{Muscle-pain, Temperature\}}(\{d\}) &= \{o_1, o_4, o_5, o_6\} \\ POS_{\{Headache\}}(\{d\}) &= \emptyset \\ POS_{\{Muscle-pain\}}(\{d\}) &= \emptyset \\ POS_{\{Temperature\}}(\{d\}) &= \{o_1, o_4, o_6\} \end{aligned}$$

Avec notre algorithmes nous avons le résultat suivants :

```
1 for combi in combinaisons(df, "Flu"):
2 liste_combi = list(combi)
3 print("POS{{{}}} = {}".format(liste_combi, POS(df, "Flu", liste_combi)))
4 Python> POS[['Headache']] = []
5 POS[['Muscle-pain']] = []
6 POS[['Temperature']] = ['o1', 'o4', 'o6']
7 POS[['Headache', 'Muscle-pain']] = ['o1']
8 POS[['Headache', 'Temperature']] = ['o1', 'o4', 'o5', 'o6']
9 POS[['Muscle-pain', 'Temperature']] = ['o1', 'o4', 'o5', 'o6']
10 POS[['Headache', 'Muscle-pain', 'Temperature']] = ['o1', 'o4', 'o5', 'o6']
```

Dans le code ci-dessus, nous calculons la *Positive Region* pour toutes les combinaisons d'attributs possibles. Ensuite nous avons la *Positive Region* et *Negative Region* pour tous les attributs :

$$\begin{aligned} POS_C &= \underline{C}X_1 \cup \underline{C}X_2 = \{o_1, o_5, o_6\} \cup \{o_4\} = \{o_1, o_5, o_6, o_4\} \\ NEG_C &= U - (\bar{C}X_1 \cup \bar{C}X_2) = \{o_1, o_2, o_3, o_4, o_5, o_6\} - (\{o_1, o_2, o_3, o_5, o_6\} \cup \{o_2, o_3, o_4\}) = \emptyset \end{aligned}$$

Avec nos algorithmes nous avons les résultats suivants :

```
1 POS_C(df, "Flu")
2 Python> ['o1', 'o4', 'o5', 'o6']
3 NEG_C(df, "Flu")
4 Python> []
```

### 3.3.4 Reduct et Core

Si nous reprenons les calculs des *Positive Region* pour chaque attributs, nous voyons que les *Reduct* sont :

$$\{Headache, Temperature\} \text{ et } \{Muscle - pain, Temperature\}$$

Avec notre algorithmes nous avons le résultat suivants :

```
1 reducts = reduct(df, "Flu")
2 Python> [['Headache', 'Temperature'], ['Muscle-pain', 'Temperature']]
```

Nous utilisons les *Reduct* pour trouver le *Core* , c'est-à-dire l'intersection des *Reduct* .

$$\{Headache, Temperature\} \cap \{Muscle - pain, Temperature\} = \{Temperature\}$$

Avec notre algorithmes nous avons le résultat suivants :

```
1 core(reducts)
2 Python> ['Temperature']
```

### 3.3.5 QuickReduct

Pour tester cet algorithme nous utiliserons un nouvel exemple.  
Appliquons étapes par étapes l'algorithme *QuickReduct*.

	a	b	c	d	e	X
x						
x1	1	2	4	0	1	1
x2	0	3	3	2	1	2
x3	2	3	1	3	3	2
x4	1	1	2	1	2	1
x5	0	2	0	1	2	1
x6	1	1	2	4	3	2
x7	2	2	1	3	2	2
x8	1	2	0	2	2	1

FIGURE 3 – Dataframe "df3"

Nous commençons avec  $R = \emptyset$ .

Puis nous calculons la dépendance pour chaque attributs :

$$\gamma_{\{a\}}(X) = 2/8$$

$$\gamma_{\{b\}}(X) = 2/8$$

$$\gamma_{\{c\}}(X) = 6/8$$

$$\gamma_{\{d\}}(X) = 6/8$$

$$\gamma_{\{e\}}(X) = 2/8$$

Nous ajoutons l'attribut possédant la plus haute dépendance, ici  $\{c\}$ .

Nous changeons la valeur de R,  $R = \{c\}$ .

Nous re-calculons la dépendance en faisant des combinaisons avec  $\{c\}$ .

$$\gamma_{\{c,a\}}(X) = 6/8$$

$$\gamma_{\{c,b\}}(X) = 6/8$$

$$\gamma_{\{c,d\}}(X) = 8/8$$

Nous obtenons la même dépendance qu'avec tous les attributs, nous pouvons arrêter de calculer, nous avons trouvé la réduction.

$$R = \{c, d\}$$

Avec notre algorithme nous avons :

```

1 R3 = quickReduct(df3, "X")
2 Python> C : ['a', 'b', 'c', 'd', 'e']
3 dep_C : 1
4 lambda([]) = 0
5 lambda(['a']) = 1/4
6 changement R = ['a']
7 lambda(['b']) = 1/4
8 lambda(['c']) = 3/4
9 changement R = ['c']
10 lambda(['d']) = 3/4
11 lambda(['e']) = 1/4
12 lambda(['c']) = 3/4
13 lambda(['c']) = 3/4
14 lambda(['c', 'a']) = 3/4
15 lambda(['c', 'b']) = 3/4

```

```
16 lambda(['c', 'd']) = 1
17 changement R = ['c', 'd']
18 lambda(['c', 'e']) = 1
19 lambda(['c', 'd']) = 1
20 Réduction = ['c', 'd']
```

## 4 Applications

Pour appliquer les différents algorithmes sur un *dataset*, nous devons trouver en trouvant permettant de faire de la classification. A l'aide du site UCI qui met à disposition des *dataset* pour faire du *Machine Learning*.

Nous utiliserons le *dataset* Statlog heart.

### 4.1 Résumé des données

Le *Statlog heart* est un contient des données sur les maladies cardiaque. Il contient 13 attributs et une classe. Les 13 attributs sont :

- *age* = 0 ou 1, mais nous ne savons pas si 1 correspond au sexe masculin et 0 féminin ou l'inverse.
- *sex*
- *chest pain type* (4 values) = type de douleur thoracique
- *resting blood pressure* = tension artérielle au repos
- *serum cholestoral in mg/dl*
- *fasting blood sugar > 120 mg/dl* = glycémie à jeun
- *resting electrocardiographic results* (values 0,1,2) = résultats électrocardiographiques au repos
- *maximum heart rate achieved* = fréquence cardiaque maximale atteinte
- *exercise induced angina* = angine induite par l'effort
- *oldpeak* = *ST depression induced by exercise relative to rest* = Dépression ST induite par l'exercice par rapport au repos
- *the slope of the peak exercise ST segment* = la pente du segment ST d'effort maximal
- *number of major vessels (0-3) colored by flourosop* = nombre de vaisseaux principaux (0-3) colorés par le fluorosop
- *thal* : 3 = normal; 6 = fixed defect; 7 = reversable defect

Nous avons aussi une classe, indiquant s'il y a présent, ou non, de maladie cardiaque.

Si nous affichons le *dataset* une fois chargé, nous avons ceci :

	age	sex	chest pain type	blood pressure	cholesterol	blood sugar	electrocardiographic	heart rate max	angina	oldpeak	peak exercise	vessels	thal	heart disease
0	70.0	1.0	4.0	130.0	322.0	0.0	2.0	109.0	0.0	2.4	2.0	3.0	3.0	2
1	67.0	0.0	3.0	115.0	564.0	0.0	2.0	160.0	0.0	1.6	2.0	0.0	7.0	1
2	57.0	1.0	2.0	124.0	261.0	0.0	0.0	141.0	0.0	0.3	1.0	0.0	7.0	2
3	64.0	1.0	4.0	128.0	263.0	0.0	0.0	105.0	1.0	0.2	2.0	1.0	7.0	1
4	74.0	0.0	2.0	120.0	269.0	0.0	2.0	121.0	1.0	0.2	1.0	1.0	3.0	1
...	...	...	...	...	...	...	...	...	...	...	...	...	...	...
265	52.0	1.0	3.0	172.0	199.0	1.0	0.0	162.0	0.0	0.5	1.0	0.0	7.0	1
266	44.0	1.0	2.0	120.0	263.0	0.0	0.0	173.0	0.0	0.0	1.0	0.0	7.0	1
267	56.0	0.0	2.0	140.0	294.0	0.0	2.0	153.0	0.0	1.3	2.0	0.0	3.0	1
268	57.0	1.0	4.0	140.0	192.0	0.0	0.0	148.0	0.0	0.4	2.0	0.0	6.0	1
269	67.0	1.0	4.0	160.0	286.0	0.0	2.0	108.0	1.0	1.5	2.0	3.0	3.0	2

270 rows x 14 columns

FIGURE 4 – Une partie du dataset

## 4.2 Analyse et visualisation des données

Commençons par regarder les statistiques et la répartition des données.

	sex	chest pain type	blood sugar	electrocardiographic	angina	oldpeak	peak exercise	vessels	thal
count	270.00	270.00	270.00	270.00	270.00	270.00	270.00	270.00	270.00
mean	0.68	3.17	0.15	1.02	0.33	1.05	1.59	0.67	4.70
std	0.47	0.95	0.36	1.00	0.47	1.15	0.61	0.94	1.94
min	0.00	1.00	0.00	0.00	0.00	0.00	1.00	0.00	3.00
25%	0.00	3.00	0.00	0.00	0.00	0.00	1.00	0.00	3.00
50%	1.00	3.00	0.00	2.00	0.00	0.80	2.00	0.00	3.00
75%	1.00	4.00	0.00	2.00	1.00	1.60	2.00	1.00	7.00
max	1.00	4.00	1.00	2.00	1.00	6.20	3.00	3.00	7.00

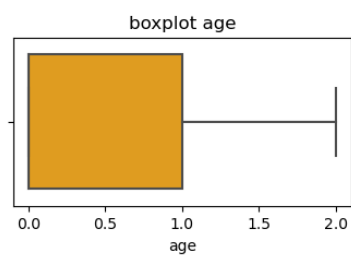
FIGURE 5 – Statistiques sur les objets

Pour pouvoir visualiser nos données et pour pouvoir appliquer les algorithmes du *Rough Set* par la suite, nous devons discrétiser nos données. Voici comment nous allons procéder :

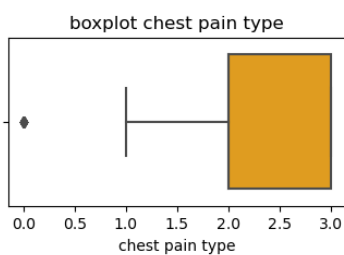
Attribut	Interval	Discretisation
age	[20 - 40)	Young
	[40 - 60)	Mid
	[60 - 80)	Old
Blood pressure	[90 - 139]	Normal
	> 140	Abnormal
Cholestoral	[150 - 250]	Normal
	< 150 or > 250	Abnormal
Mximun heart rate	< 60	Low
	[60 - 100]	Normal
	> 100	Hight

TABLE 3 – Tableau de discretisation

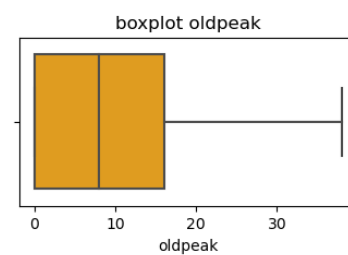
Regardons ensuite les *boxplot* de chaque variable.



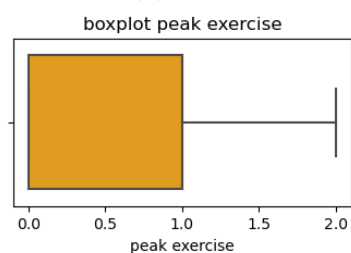
(a) Age



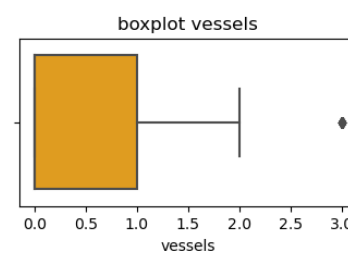
(b) Chest pain type



(c) Oldpeak



(d) The slope of the peak  
exercise ST segment



(e) Number of major vessels

FIGURE 6 – Résumé des données en fonction de l'attribut

Nous constatons qu'il y a une valeur aberrante en *Chest pain type* et en *Number of major vessels*. Nous devons traiter ça dans la partie *preprocessing*.

Après cette transformation, affichons la répartition des données par attributs.

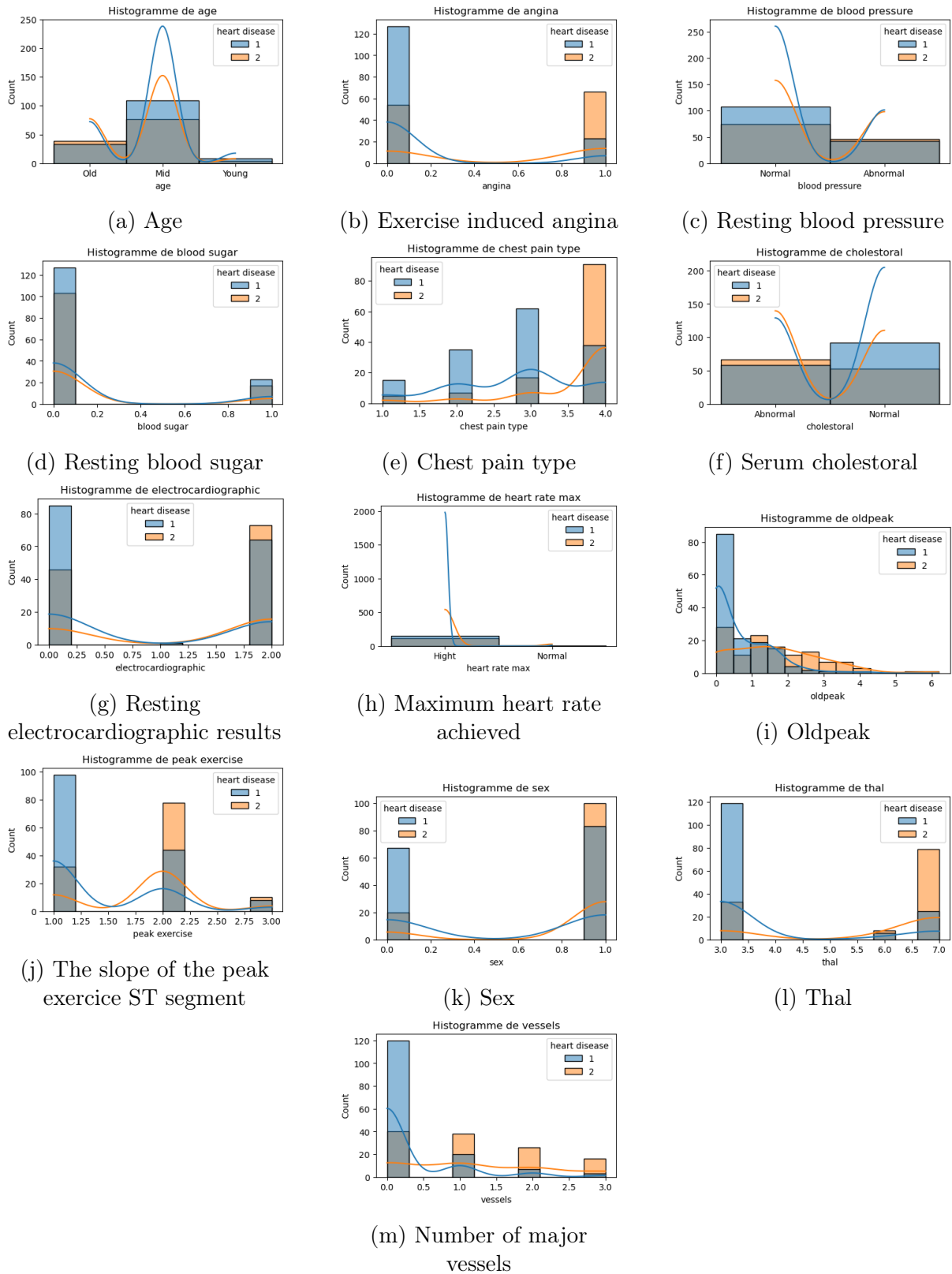


FIGURE 7 – Résumé des données en fonction de l'attribut



Pour l'âge, notre *dataset* est largement composé de **Mid** personnes (entre 40 et 60 ans) et il y a une prévalence de présence de maladie cardiaque.  
A FINIR...

Enfin, regardons la répartition des objets en fonction des classes. Nous avons une

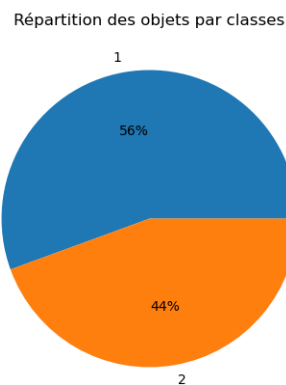


FIGURE 8 – Répartition des objets par classe

répartition quasiment 50/50, nous ne serons pas obligé d'équilibrer le *dataset* dans la partie *preprocessing*.

### 4.3 *Preprocessing*

Après avoir analysé les données nous allons pouvoir les traiter. Comme nous avons pu le remarquer dans les 6 nous devons supprimer les valeurs aberrantes. Ensuite

Pour la phase de *preprocessing* nous commençons par **encoder** les attributs du *dataset*. Pour cela nous utilisons un *OrdinalEncoder* pour convertir les attributs de type ordinal en numérique.

### 4.4 Visualisation des données

### 4.5 Application de modèles de classification

#### 4.5.1 Modèle A

#### 4.5.2 Modèle B

#### 4.5.3 Modèle C

#### 4.5.4 Modèle D

### 4.6 Résultat

#### 4.6.1 Évaluation

#### 4.6.2 Comparaison

## 5 Conclusion

### 5.1 Rappel des résultats attendues et obtenues

### 5.2 Limites de l'algorithme

### 5.3 Suggestions d'améliorations

## Bibliographie

- [Len22] Chinmoy LENKA. *Python / Difference between two lists*. 2022. URL : <https://www.geeksforgeeks.org/python-difference-two-lists/>.
- [man19] MANJEET\_04. *Python / Find common elements in list of lists*. 2019. URL : <https://www.geeksforgeeks.org/python-find-common-elements-in-list-of-lists/>.
- [Ora] ORACLE. *Qu'est-ce que le Data Mining?* URL : <https://www.oracle.com/fr/database/data-mining-definition.html>.
- [Wik] WIKIPEDIA. *Relation d'équivalence*. URL : [https://fr.wikipedia.org/wiki/Relation\\_d'equivalence#](https://fr.wikipedia.org/wiki/Relation_d'equivalence#).