

Projet de conception et de programmation

Mickael LE DENMAT

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

5 février 2023

Table des matières

1 Introduction

2 *Rough Set Theory*

3 Application

4 Conclusion

5 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.

Rough Set Theory

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

Table – Dessin pour le calcul du IND_C

$$IND_C = \{\{o_1\}, \{o_2, o_3\}, \{o_4\}, \{o_5\}, \{o_6\}\}.$$

Rappel

$$X_1 = \{o_j \mid Flu(o_j) = \{Yes\}\} = \{o_1, o_2, o_5, o_6\}$$

$$X_2 = \{o_j \mid Flu(o_j) = \{No\}\} = \{o_3, o_4\}$$

$$IND_C = \{\{o_1\}, \{o_2, o_3\}, \{o_4\}, \{o_5\}, \{o_6\}\}.$$

Rappel

$$X_1 = \{o_j \mid Flu(o_j) = \{Yes\}\} = \{o_1, o_2, o_5, o_6\}$$

$$X_2 = \{o_j \mid Flu(o_j) = \{No\}\} = \{o_3, o_4\}$$

$$IND_C = \{\{o_1\}, \{o_2, o_3\}, \{o_4\}, \{o_5\}, \{o_6\}\}.$$

B-Lower

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

$$\underline{C}X_2 = \{o_4\}$$

B-Lower & B-Upper Approximation

Rappel

$$X_1 = \{o_j | Flu(o_j) = \{Yes\}\} = \{o_1, o_2, o_5, o_6\}$$

$$X_2 = \{o_j | Flu(o_j) = \{No\}\} = \{o_3, o_4\}$$

$$IND_C = \{\{o_1\}, \{o_2, o_3\}, \{o_4\}, \{o_5\}, \{o_6\}\}.$$

B-Lower

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

$$\underline{C}X_2 = \{o_4\}$$

B-Lower

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

$$\bar{C}X_2 = \{o_2, o_3, o_4\}$$

Rappel

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

$$\underline{C}X_2 = \{o_4\}$$

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

$$\bar{\bar{C}}X_2 = \{o_2, o_3, o_4\}$$

Positive & Negative Region

Rappel

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

$$\underline{C}X_2 = \{o_4\}$$

$$\bar{\underline{C}}X_1 = \{o_1, o_2, o_3, o_5, o_6\}$$

$$\bar{\underline{C}}X_2 = \{o_2, o_3, o_4\}$$

Positive Region

$$POS_C = \underline{C}X_1 \cup \underline{C}X_2$$

$$= \{o_1, o_5, o_6\} \cup \{o_4\}$$

$$= \{o_1, o_5, o_6, o_4\}$$

Positive & Negative Region

Rappel

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

$$\underline{C}X_2 = \{o_4\}$$

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

$$\bar{C}X_2 = \{o_2, o_3, o_4\}$$

Positive Region

$$POS_C = \underline{C}X_1 \cup \underline{C}X_2$$

$$= \{o_1, o_5, o_6\} \cup \{o_4\}$$

$$= \{o_1, o_5, o_6, o_4\}$$

Negative Region

$$NEG_C = U - (\bar{C}X_1 \cup \bar{C}X_2)$$

$$= \{o_1, o_2, o_3, o_4, o_5, o_6\} - (\{o_1, o_2, o_3, o_5, o_6\} \cup \{o_2, o_3, o_4\})$$

$$= \emptyset$$

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.

Exemple

[["Headache", "Temperature"], ["Muscle-pain", "Temperature"]]

Définition

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

Définition

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

Exemple

["Temperature"]

Rappel

$$POS_C = \{o_1, o_5, o_6, o_4\}$$

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

Exemple

$$\lambda = POS_C / U$$

$$= |\{o_1, o_5, o_6, o_4\}| / |\{o_1, o_2, o_3, o_4, o_5, o_6\}|$$

$$= 4/6$$

- Nous calculons toutes les dépendances pour un attributs.

- Nous calculons toutes les dépendances pour un attributs.
- Nous prenons l'attribut possédant la dépendance la plus haute et nous l'ajoutons dans la réduction.

- Nous calculons toutes les dépendances pour un attributs.
- Nous prenons l'attribut possèdent la dépendance la plus haute et nous l'ajoutons dans la réduction.
- Nous calculons la dépendance de la réduction en ajoutons à un à chaque attributs.

- Nous calculons toutes les dépendances pour un attributs.
- Nous prenons l'attribut possédant la dépendance la plus haute et nous l'ajoutons dans la réduction.
- Nous calculons la dépendance de la réduction en ajoutons à un à chaque attributs.
- Nous ajoutons l'attribut ayant la meilleur dépendance dans la réduction

- Nous calculons toutes les dépendances pour un attributs.
- Nous prenons l'attribut possédant la dépendance la plus haute et nous l'ajoutons dans la réduction.
- Nous calculons la dépendance de la réduction en ajoutons à un à chaque attributs.
- Nous ajoutons l'attribut ayant la meilleur dépendance dans la réduction
- Nous nous arrêtons quand nous obtenons la même dépendance qu'avec la totalité des attributs.

Application

Dataset Statlog Heart

	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 × 14 columns

Figure – Extrait du dataset

Preprocessing

- Détection des valeurs manquantes/ nulles.

- Détection des valeurs manquantes/ nulles.
- Discretisation.

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
Maximum heart rate	< 60	Low
	[60 - 100]	Normal
	> 100	Hight

Figure – Tableau de discrétisation

- Détection des valeurs manquantes/ nulles.
- Discretisation.
- Encodage.

- Détection des valeurs manquantes/ nulles.
- Discrétisation.
- Encodage.
- Équilibrer le nombre d'objet par classe.

Répartition des objets par classes

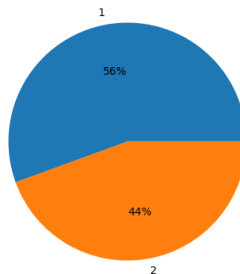


Figure – Répartition des objets par classe

- Détection des valeurs manquantes/ nulles.
- Discrétisation.
- Encodage.
- Équilibrer le nombre d'objet par classe.
- Normalisation (attributs/ classes).

Data visualization

- Matrice de corrélation.

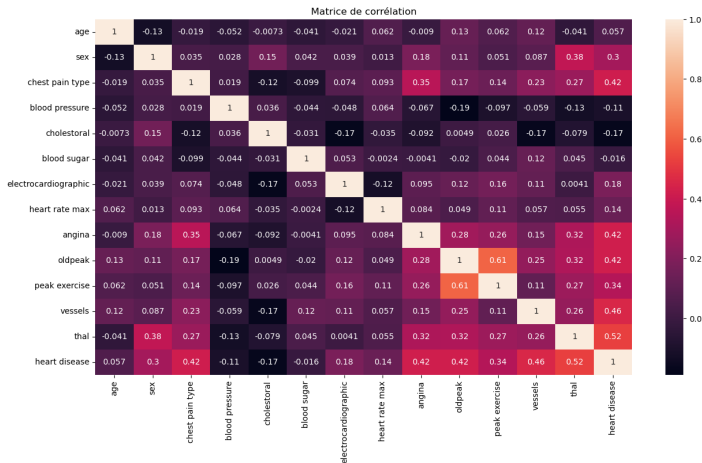


Figure – Matrice de corrélation

- Matrice de corrélation.
- Boxplot => suppression "outliers".

- Matrice de corrélation.
- Boxplot => suppression "outliers".
- Histogrammes.

- Matrice de corrélation.
- Boxplot => suppression "outliers".
- Histogrammes.
- Pairplot.

Réductions trouvées

Liste des réductions trouvées

Méthode	age	sex	chest pain type	blood pressure	cholesterol	blood sugar	electrocardiographic	heart rate max	angina	oldpeak	peak exercise	vessels	thal
Reduct	✓	✓	✓	✓	X	X	✓	X	X	✓	X	✓	✓
QuickReduct	✓	✓	✓	✓	X	X	✓	X	X	✓	X	✓	✓
Variance	X	✓	X	✓	✓	X	✓	X	✓	X	X	X	✓
Fuzzy	✓	✓	✓	✓	✓	✓	✓	X	X	✓	X	X	X

Table – Liste des réductions

Liste des réductions trouvées

Méthode	age	sex	chest pain type	blood pressure	cholesterol	blood sugar	electrocardiographic	heart rate max	angina	oldpeak	peak exercise	vessels	thal
Reduct	✓	✓	✓	✓	X	X	✓	X	X	✓	X	✓	✓
QuickReduct	✓	✓	✓	✓	X	X	✓	X	X	✓	X	✓	✓
Variance	X	✓	X	✓	✓	X	✓	X	✓	X	X	X	✓
Fuzzy	✓	✓	✓	✓	✓	✓	✓	X	X	✓	X	X	X

Table – Liste des réductions

Liste des réductions trouvées

Méthode	age	sex	chest pain type	blood pressure	cholesterol	blood sugar	electrocardiographic	heart rate max	angina	oldpeak	peak exercise	vessels	thal
Reduct	✓	✓	✓	✓	X	X	✓	X	X	✓	X	✓	✓
QuickReduct	✓	✓	✓	✓	X	X	✓	X	X	✓	X	✓	✓
Variance	X	✓	X	✓	✓	X	✓	X	✓	X	X	X	✓
Fuzzy	✓	✓	✓	✓	✓	✓	✓	X	X	✓	X	X	X

Table – Liste des réductions

Liste des réductions trouvées

Méthode	age	sex	chest pain type	blood pressure	cholesterol	blood sugar	electrocardiographic	heart rate max	angina	oldpeak	peak exercise	vessels	thal
Reduct	✓	✓	✓	✓	X	X	✓	X	X	✓	X	✓	✓
QuickReduct	✓	✓	✓	✓	X	X	✓	X	X	✓	X	✓	✓
Variance	X	✓	X	✓	✓	X	✓	X	✓	X	X	X	✓
Fuzzy	✓	✓	✓	✓	✓	✓	✓	X	X	✓	X	X	X

Table – Liste des réductions

Liste des réductions trouvées

Méthode	age	sex	chest pain type	blood pressure	cholesterol	blood sugar	electrocardiographic	heart rate max	angina	oldpeak	peak exercise	vessels	thal
Reduct	✓	✓	✓	✓	X	X	✓	X	X	✓	X	✓	✓
QuickReduct	✓	✓	✓	✓	X	X	✓	X	X	✓	X	✓	✓
Variance	X	✓	X	✓	✓	X	✓	X	✓	X	X	X	✓
Fuzzy	✓	✓	✓	✓	✓	✓	✓	X	X	✓	X	X	X

Table – Liste des réductions

Liste des réductions trouvées

Méthode	age	sex	chest pain type	blood pressure	cholesterol	blood sugar	electrocardiographic	heart rate max	angina	oldpeak	peak exercise	vessels	thal
Reduct	✓	✓	✓	✓	X	X	✓	X	X	✓	X	✓	✓
QuickReduct	✓	✓	✓	✓	X	X	✓	X	X	✓	X	✓	✓
Variance	X	✓	X	✓	✓	X	✓	X	✓	X	X	X	✓
Fuzzy	✓	✓	✓	✓	✓	✓	✓	X	X	✓	X	X	X

Table – Liste des réductions

Classification

Liste des modèles appliqués :

- Gaussian Process

Liste des modèles appliqués :

- Gaussian Process
- Random Forest

Liste des modèles appliqués :

- Gaussian Process
- Random Forest
- Nearest Neighbors

Liste des modèles appliqués :

- Gaussian Process
- Random Forest
- Nearest Neighbors
- Autres (Linear SVM, RBF SVM, Neural Net, Naive Bayes, Decision Tree, AdaBoost, QuadraticDiscriminantAnalysis)

Comparaison

```
sans_reduc  
  
[array([[23,  4],  
        [ 5, 15]], dtype=int64),  
 array([[22,  5],  
        [ 5, 15]], dtype=int64),  
 array([[24,  3],  
        [ 8, 12]], dtype=int64)]
```

Figure – Matrice de confusion de sans *feature selection*

Matrice de corrélation

```
reduct
Gaussian Process
[[-1  1]
 [-5  5]]
Random Forest
[[ 2 -2]
 [-4  4]]
Nearest Neighbors
[[-2  2]
 [-3  3]]
```

Figure – Différence entre avec *reduct* et sans

Matrice de corrélation

```
quickreduct
Gaussian Process
[[-1  1]
 [-5  5]]
Random Forest
[[ 2 -2]
 [-3  3]]
Nearest Neighbors
[[-2  2]
 [-3  3]]
```

Figure – Différence entre avec *quickreduct* et sans

```
variance
Gaussian Process
[[-4  4]
 [-4  4]]
Random Forest
[[-3  3]
 [-3  3]]
Nearest Neighbors
[[-3  3]
 [-2  2]]
```

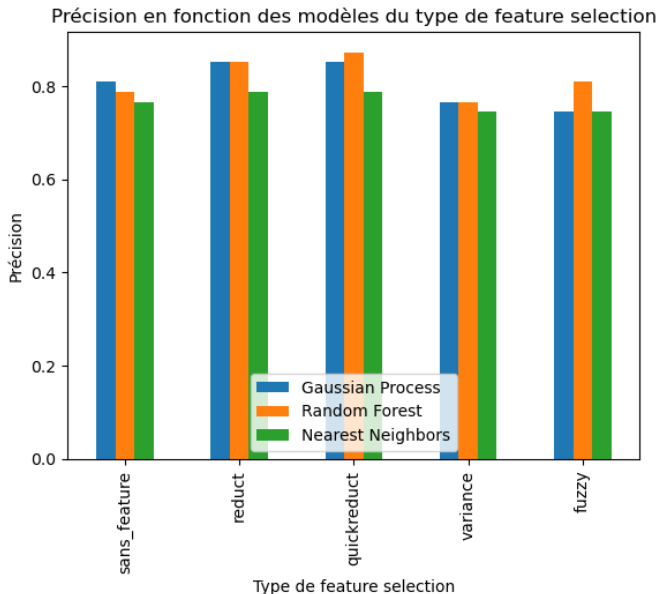
Figure – Différence entre avec *variance* et sans

Matrice de corrélation

```
fuzzy
Gaussian Process
[[-2  2]
 [-1  1]]
Random Forest
[[-0  0]
 [-1  1]]
Nearest Neighbors
[[-2  2]
 [-1  1]]
```

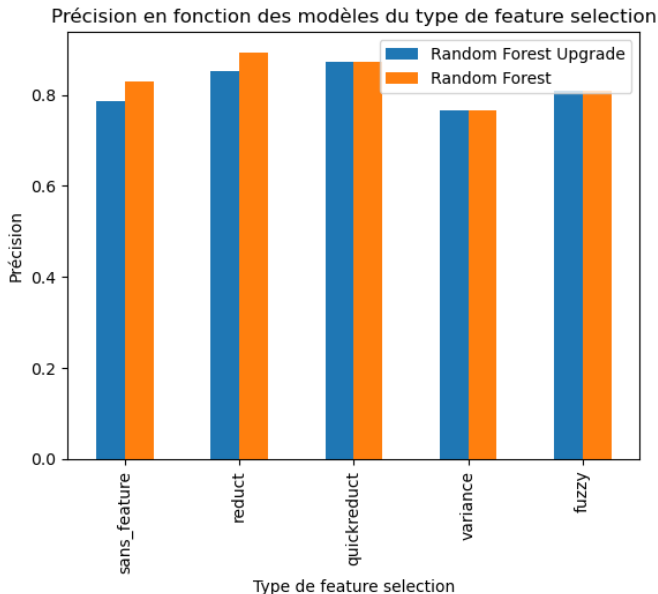
Figure – Différence entre avec *fuzzy* et sans

Comparaison de la précision



Amélioration des modèles.

Amélioration du *Random Forest*



Limite du Quickreduct

- Générer tous les sous ensembles possibles d'attributs.

- Générer tous les sous ensembles possibles d'attributs.
- Re-calcul tout à chaque fois.

Amélioration du Quickreduct

- Stocker les résultats et faire une version incrémentale.

- Stocker les résultats et faire une version incrémentale.
- Trouver un moyen de partir de l'ensemble des attributs et de supprimer les moyens utiles si c'est plus rapide (comment le savoir?).



Oracle.

Qu'est-ce que le data mining ?

Merci pour votre attention !
Avez vous des questions ?