

UNIVERSITÉ DES SCIENCES ET DE LA TECHNOLOGIE HOUARI
BOUMEDIENE



DATA MINING

Rapport Projet Partie 02
Technique de Data Mining

HEBBACHE IMAD EDDINE,

KHOUCHA MADANI,

MEZIANI MOHAMED IMAD,

Mme. Drias

Mr. Khennak

1^{er} janvier 2022

Introduction

Le Data Mining est un domaine très large qui consiste à appliquer un ensemble de processus ou de techniques sur les données brutes dont le but d'explorer et analyser ces données pour les rendre utiles.

Dans cette deuxième partie de projet on va implémenter ces techniques qui permettent en premier lieu de normaliser et de discrétiser les données, ensuite d'extraire des motifs fréquents de ces données et à la fin d'appliquer une classification à partir des données initiales.

Toutes ces techniques seront appliquées à partir d'une IHM conçue avec l'utilisation du langage Java.

1 Objectifs

Tout cela nous a mené à fixer les objectifs suivants :

- Implémenter les deux algorithmes de normalisation :
 - La normalisation Min-Max.
 - La normalisation Z-score
- Développer et appliquer les deux types de discrétisation suivantes :
 - La discrétisation en classes d'effectifs égaux.
 - La discrétisation en classes d'amplitudes égales.
- Appliquer les deux algorithmes (Apriori et Eclat) d'extraction de motifs fréquents sur les données discrétisées.
- Classification des instances du dataset par l'implémentation de l'algorithme de la classification Naïve Bayésienne et KNN.

2 Pré-traitement des données

Le pré-traitement des données est une étape importante du processus d'exploration et d'analyse des données qui prend les données dans leur format initial et les transforme en un format pouvant être compris et analysé par les ordinateurs.

On dit souvent que les données pré-traitées sont plus importantes que les algorithmes les plus puissants.

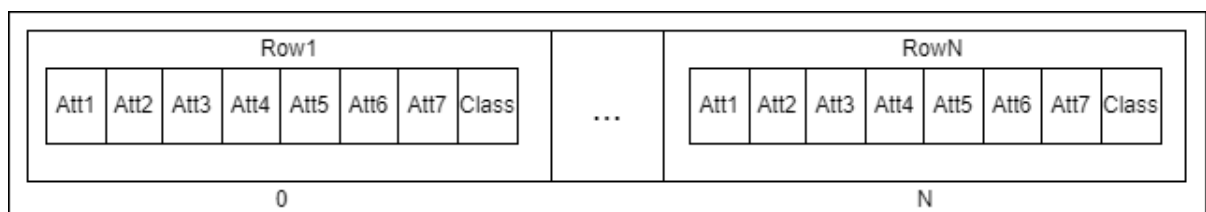
2.1 La normalisation

1. Description :

La normalisation permet de mettre les données à l'échelle dans un intervalle régularisé spécifique à fin de donner une précision pour la comparaison des données.

2. Structures utilisées :

Le résultat est retourné sous forme de tableau de lignes (Class Row).



Normalisation avec Min-Max

C'est une technique de normalisation de données, une transformation linéaire est effectuée sur les données. La valeur minimale et maximale sont extraites des données et chaque valeur est remplacée selon la formule suivante :

$$VALEUR_{(i,nouvelle)} = \frac{VALEUR_{(i,courante)} - VALEUR_{(min,courante)}}{VALEUR_{(max,courante)} - VALEUR_{(min,courante)}} (VALEUR_{(max,nouvelle)} - VALEUR_{(min,nouvelle)}) + VALEUR_{(min,nouvelle)}$$

La normalisation Min-Max préserve les relations entre les valeurs de données d'origine.

Algorithm

input : Un tableau d'attributs *Attr*, Valeur Minimale *Min*, Valeur Maximale *Max*
output : Un tableau de lignes (Row) *Normalisedb*

```
1 while Le tableau d'attributs non vide do
2   for Chaque ligne do
3     Valeur(i,nouvelle) ←  $\frac{Valeur(i,courante) - Valeur(min,courante)}{Valeur(max,courante) - Valeur(min,courante)} \times (Valeur(max,nouvelle) -$   

4      $Valeur(min,nouvelle)) + Valeur(min,nouvelle)$ 
5   end
6 while Attribut non vide do
7   Créer une nouvelle ligne avec les valeurs de chaque attribut avec le même indice Ajouter la ligne
8   dans le tableau de lignes Normalised
9 end
```

Normalisation avec Z-score

Un score Z est une mesure numérique qui décrit la relation d'une valeur avec la moyenne d'un ensemble de valeurs.

La normalisation avec Z-score exprime l'écart par rapport à la valeur de la moyenne selon la formule suivante :

$$VALEUR_{(i,nouvelle)} = \frac{VALEUR_{(i,courante)} - VALEUR_{(moyenne,courante)}}{S}$$
$$S = \frac{1}{N} \sum_{i=1}^N |VALEUR_{(i,courante)} - VALEUR_{(moyenne,courante)}|$$

S représente formule de calcul de l'écart-type.

```

output : Un tableau de lignes (Row) Normalised
9 while Le tableau d'attributs non vide do
10   moyenne ← Moyenne(Attribut)  $S = \frac{1}{N} \sum_{i=0}^N \text{Valeur}(i, \text{courante}) - \text{moyenne}$ ;
11   for Chaque ligne do
12      $\text{Valeur}(i, \text{nouvelle}) \leftarrow \frac{\text{Valeur}(i, \text{courante}) - \text{moyenne}}{S}$ 
13   end
14 end
15 while Attribut non vide do
16   Créer une nouvelle ligne avec les valeurs de chaque attribut avec le même indice   Ajouter la ligne
    dans le tableau de lignes Normalised
17 end

```

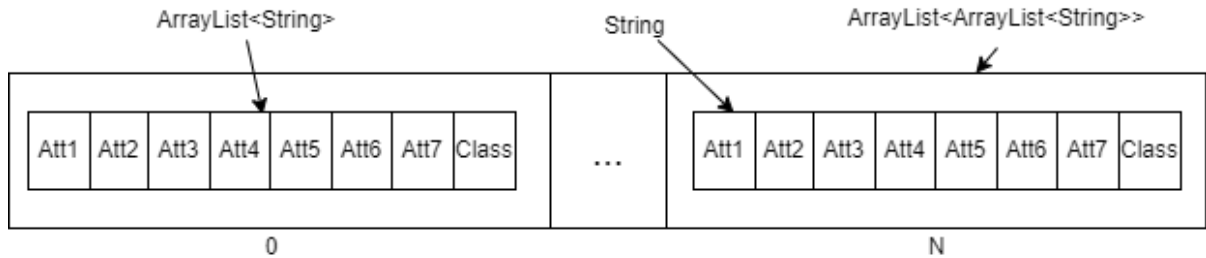
2.2 La discrétisation

1. Description :

La discrétisation regroupe les données en intervalles plus petits, c'est un peu similaire au groupement des données par classe mais cela se fait après le nettoyage des données.

2. Structures utilisées :

La structure utilisée est un tableau de tableaux de chaîne de caractères.



Discrétisation en classes d'effectifs égaux.

Cette méthode permet de créer des intervalles avec le même nombre d'attributs.

```

input : Un tableau d'attributs Attr, NombredequantilesQ
output : Un tableau de tableau de chaîne de caractère Discretised
18 while Le tableau d'attributs non vide do
19    $\text{NombreIntervals} \leftarrow \frac{N}{Q}$  Créer des intervalles
20   for Chaque ligne do
21     Attribuer à chaque valeur la classe de l'intervalle auquel elle appartient
22   end
23 end
24 while Attribut non vide do
25   Créer un tableau de chaîne de caractère avec la classe de chaque attribut avec le même indice   Ajouter
    le tableau dans le tableau Discretised
26 end

```

Discrétisation en classes d'amplitudes égales

Cette méthode permet de créer des intervalles avec la même étendue.

```

input  : Un tableau d'attributs Attr, Nombred'intervallesQ
output : Un tableau de tableau de chaine de caractere Dicretised
27 while Le tableau d'attributs non vide do
28   |  $LongeurInterval \leftarrow \frac{Valeur(max)-Valeur(min)}{Q}$    Créer des intervalles
29   for Chaque ligne do
30   |   Attribuer a chaque valeur la classe de l'intervalle auquel elle appartient
31   end
32 end
33 while Attribut non vide do
34   | Cree un tableau de chaîne de caractère avec la classe de chaque attribut avec le même indice   Ajouter
35   | le tableau dans le tableau Dicretised
36 end

```

3 Extraction des motifs fréquents

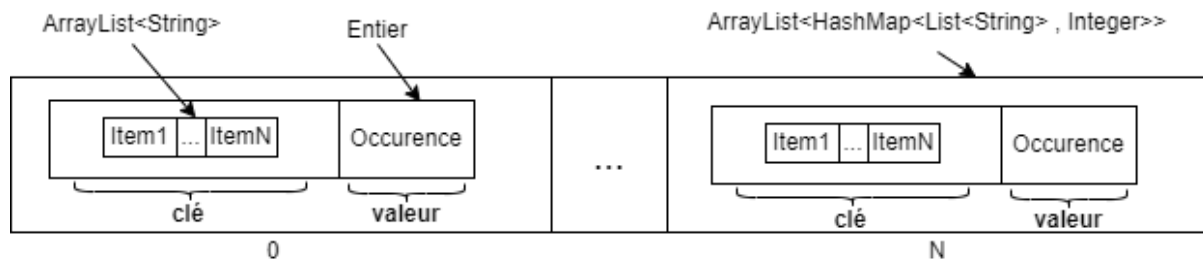
3.1 Les itemsets fréquents

1. Description :

Les itemsets fréquents représentent les items du dataset avec une fréquence supérieur au support minimal qui sera définit comme paramètre dans l'algorithme.

2. Structures utilisées :

La structure utilisée est un tableau (ArrayList) de HashMap qui est une structure de données utilisée pour stocker les objets paire clé-valeur qui a comme clé un tableau de chaînes de caractères qui contient les items et comme valeur un entier qui représente le support de la clé.



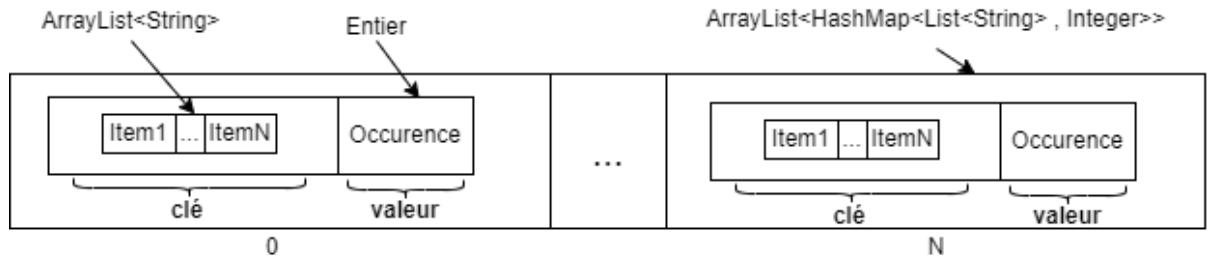
3.2 Apriori

(a) Description :

L'algorithme Apriori est le premier algorithme propose pour l'extraction des motifs fréquents, une approche itérative pour découvrir les itemsets les plus fréquents. Pour chaque itération il réduit l'espace de recherche en détectant les motifs les plus fréquents par rapport une valeur dite Support minimum.

(b) Structures utilisées :

La structure utilisée est un tableau (ArrayList) de HashMap qui est une structure de données utilisée pour stocker les objets paire clé-valeur qui a comme clé un tableau de chaînes de caractères qui contient les items et comme valeur un entier qui représente le support de la clé.



(c) Pseudo-Code :

```

input   : Un tableau d'attributs discretisee, Support Minimal
output  : La liste des itemsets frequents
36  $L_1 \leftarrow (\text{Itemsets frequents})$  while  $L_{new}$  nonvide do
37    $C_{new} \leftarrow \text{generer toute les combinaison possible}$ 
38   for Chaque ligne de  $C_{new}$  do
39     Calculer le nombre de transactions dans les quelles figure l'intersection des items de la ligne
40     if la fréquence de ligne > support minimal then
41        $L_{new} \leftarrow L_{new} + \text{les itemsets de cette ligne}$ 
42     end
43   end
44    $L \leftarrow L + L_{new}$ 
45 end

```

3.3 Eclat

(a) Description :

L'algorithme ECLAT signifie Equivalence Class Clustering and bottom-up Lattice Traversal, il est considéré comme une version plus améliorée de Apriori. Il utilise l'approche de recherche en profondeur d'abord alors que Apriori utilise l'approche en largeur d'abord ce qui rend Eclat plus rapide parce qu'il utilise moins de mémoire.

(b) Pseudo-code :

```

input   : Support Minimale  $supMin$ , Nombre de transactions
output  : Items-fréquents
2.  $table\ verticale \leftarrow \text{inverse la BD}$ 
3. Éliminé les linges avec support inférieur a  $supMin$ 
4. .3 for Toutes les combinaisons possible des items do
   $. .a créer une nouvelle ligne contenant
   | i.l'union du contenu des premières colonnes de  $table\ verticale$ 
   | ii.l'intersection des contenus des deuxièmes colonnes de  $table\ verticale$ 
   end
6. .4 répétez les étapes 2 et 3 jusqu'à ce qu'aucun nouvel ensemble d'éléments ne puisse être créé.

```

4 Classification supervisée des instances du dataset

4.1 Classification Naïve Bayésienne.

La classification Naïve Bayésienne est une classification qui se base sur les probabilités, elle s'applique sur des données discrétisées. Cette classification est peu fiable car elle traite chaque caractéristique indépendamment mais requiert peu de données d'entraînement.

```
input  : Un tableau d'attributs discrétisés, Données à classifier
output : La classe
46 for Chaque Attribut k do
47   for Chaque Classe i do
48      $P(X_k|Classe_i) \leftarrow \text{lenombre d'apparition de } X_k \text{ avec la classe } i$  end
49   end
50   for Chaque Classe i do
51      $P(X|Classe_i) \leftarrow \prod P(X_j|Classe_i)$  end
52   Choisir La classe avec la plus haute probabilité
```

4.2 Classification K-nearest neighbors(KNN).

La classification Knn est une classification qui se base sur les k plus proches voisins, elle peut être appliquée sur des données quantitatives en calculant la distance avec les mesures euclidienne, Manhattan ou plusieurs autres mesures, comme il peut s'appliquer sur des données discrétisées avec comme exemple la mesure de hamming.

```
input  : Un tableau d'attributs discrétisés, Données à classifier
output : La classe
53 for Chaque ligne du dataset do
54   Calculer la distance entre nos données et la ligne
55 end
56 Trier les résultats selon la distance par ordre DESC ;
   Récupérer les K premiers instances ;
   Choisir La classe avec la plus fréquente entre les K instances ;
```

5 Interfaces

Pour visualiser nos résultats nous avons conçu un IHM simple à utiliser qui facilite cette opération en utilisant la bibliothèque **javaFX** de JAVA, ci-dessous les différentes fonctionnalités de notre interface.

5.1 Interfaces de normalisation

L'onglet ci-dessous permet la normalisation avec les méthodes Min-Min ou Z-Score.

Interface de normalisation Min-Max

Pour la normalisation avec Min-Max il faut introduire les deux valeurs Min et MAX.

Informations sur la base	Informations sur les attributs	Histogrammes et graphes	Normalisation	Discrétisation	Extraction de motifs	Classification
--------------------------	--------------------------------	-------------------------	---------------	----------------	----------------------	----------------

Nombre d'intervalles

Discrétiser

4

Instance	surface	périmètre	compacité	longueur du grain	largeur du grain	coefficient d'asymétrie	longueur du sillon du noyau	Classe
1	I12	I23	I33	I42	I52	I61	I72	Kama
2	I12	I22	I33	I42	I53	I61	I71	Kama
3	I12	I22	I34	I41	I53	I62	I71	Kama
4	I12	I22	I34	I41	I53	I61	I71	Kama
5	I13	I23	I34	I42	I53	I61	I72	Kama
6	I12	I22	I34	I42	I52	I61	I71	Kama
7	I12	I22	I33	I42	I52	I62	I72	Kama
8	I12	I22	I34	I42	I52	I62	I71	Kama
9	I13	I23	I33	I43	I53	I61	I73	Kama
10	I13	I23	I33	I43	I53	I61	I72	Kama
11	I12	I23	I33	I42	I52	I62	I72	Kama
12	I12	I22	I33	I42	I52	I61	I71	Kama
13	I12	I22	I33	I42	I52	I62	I71	Kama
14	I12	I22	I33	I42	I52	I62	I71	Kama
15	I12	I22	I33	I42	I52	I62	I71	Kama
16	I12	I22	I34	I42	I53	I62	I71	Kama
17	I12	I22	I34	I41	I53	I63	I71	Kama
18	I12	I22	I34	I42	I53	I61	I72	Kama
19	I12	I22	I34	I41	I53	I61	I71	Kama

Discrétisation avec Quantiles

Pour la discrétisation avec Quantiles il faut introduire le nombre de quantiles.

Informations sur la base	Informations sur les attributs	Histogrammes et graphes	Normalisation	Discrétisation	Extraction de motifs	Classification
--------------------------	--------------------------------	-------------------------	---------------	----------------	----------------------	----------------

Nombre de Quantiles

Discrétiser

70

Instance	surface	périmètre	compacité	longueur du grain	largeur du grain	coefficient d'asymétrie	longueur du sillon du noyau	Classe
1	I1-2	I2-2	I3-2	I4-2	I5-2	I6-1	I7-2	Kama
2	I1-2	I2-2	I3-2	I4-2	I5-2	I6-1	I7-1	Kama
3	I1-2	I2-2	I3-3	I4-1	I5-2	I6-1	I7-1	Kama
4	I1-2	I2-2	I3-3	I4-1	I5-2	I6-1	I7-1	Kama
5	I1-3	I2-2	I3-3	I4-2	I5-3	I6-1	I7-2	Kama
6	I1-2	I2-2	I3-3	I4-2	I5-2	I6-1	I7-1	Kama
7	I1-2	I2-2	I3-2	I4-2	I5-2	I6-2	I7-2	Kama
8	I1-2	I2-2	I3-3	I4-2	I5-2	I6-1	I7-1	Kama
9	I1-3	I2-3	I3-2	I4-3	I5-3	I6-1	I7-3	Kama
10	I1-3	I2-3	I3-3	I4-3	I5-3	I6-1	I7-3	Kama
11	I1-2	I2-2	I3-2	I4-2	I5-2	I6-3	I7-2	Kama
12	I1-2	I2-2	I3-2	I4-2	I5-2	I6-1	I7-1	Kama
13	I1-2	I2-2	I3-3	I4-2	I5-2	I6-2	I7-1	Kama
14	I1-2	I2-2	I3-2	I4-2	I5-2	I6-2	I7-1	Kama
15	I1-2	I2-2	I3-2	I4-2	I5-2	I6-2	I7-1	Kama
16	I1-2	I2-2	I3-3	I4-2	I5-2	I6-2	I7-1	Kama
17	I1-2	I2-2	I3-3	I4-1	I5-2	I6-3	I7-1	Kama
18	I1-2	I2-2	I3-3	I4-2	I5-3	I6-1	I7-1	Kama
19	I1-2	I2-2	I3-3	I4-1	I5-3	I6-1	I7-1	Kama

5.3 Interfaces d'extraction de motifs fréquents

l'onglet ci-dessous permet l'extraction des motifs fréquents avec les différentes algorithmes Apriori et Eclat ainsi l'extraction des règles d'association et de corrélation

Informations sur la base

Informations sur les attributs

Histogrammes et graphes

Normalisation

Discrétisation

Extraction de motifs

Classification

Apriori

Support minimum

Nombre d'intervalle de discrétisation

Appliquer

Eclat

Support minimum

Nombre d'intervalle de discrétisation

Appliquer

Confiance minimum

Génération des règles d'association

Règle d'association

Confiance

Aucun contenu dans la table

Génération des règles de corrélation

Règle de corrélation

Lift

Aucun contenu dans la table

Interface de l’application des algorithmes Apriori et Eclat

Tous les deux demande en entrée le support minimum ainsi que la valeur de l’intervalle de discrétisation.

Apriori

20

4

[I72] => 92

[I71] => 49

[I12, I22] => 47

[I51, I72] => 42

[I41, I51] => 47

[I52] => 57

[I51] => 62

[I32] => 50

[I53] => 52

[I12] => 57

[I34] => 43

[I11] => 81

[I33] => 98

[I11, I51, I72] => 42

[I11, I41, I51] => 47

[I11, I21, I51] => 55

[I21, I41] => 61

[I11, I41] => 62

[I42, I72] => 55

[I33, I62] => 49

[I611] => 55

Appliquer

48 milliseconds

Eclat

20

4

[I72] => 92

[I71] => 49

[I12, I22] => 47

[I51, I72] => 42

[I41, I51] => 47

[I52] => 57

[I51] => 62

[I32] => 50

[I53] => 52

[I12] => 57

[I34] => 43

[I11] => 81

[I33] => 98

[I11, I51, I72] => 42

[I11, I41, I51] => 47

[I11, I21, I51] => 55

[I21, I41] => 61

[I11, I41] => 62

[I42, I72] => 55

[I33, I62] => 49

[I611] => 55

Appliquer

21 milliseconds

Extraction des règles d’association et de corrélation

Après l’extraction des motifs fréquents en utilisant l’un des deux algorithmes en peut générer les règles d’association et les règles de corrélation en introduisant la confiance minimum.

9

60

Génération des règles d'association

Génération des règles de corrélation

Règle d'association	Confiance
[I12] => [I22]	82.4561403508772
[I22] => [I12]	79.66101694915254
[I51] => [I72]	67.74193548387096
[I41] => [I51]	70.1492537313433
[I51] => [I41]	75.80645161290323
[I51] => [I11, I72]	67.74193548387096
[I11, I51] => [I72]	67.74193548387096
[I11, I72] => [I51]	87.50000000000001
[I51, I72] => [I11]	100.0
[I41] => [I11, I51]	70.1492537313433
[I11, I41] => [I51]	75.80645161290323
[I51] => [I11, I41]	75.80645161290323
[I11, I51] => [I41]	75.80645161290323
[I41, I51] => [I11]	100.0
[I11] => [I21, I51]	67.90123456790124
[I21] => [I11, I51]	79.71014492753625
[I11, I21] => [I51]	79.71014492753625

Règle de corrélation	Lift
[I12] => [I22]	2.934879571810883
[I22] => [I12]	2.934879571810883
[I51] => [I72]	1.546283309957924
[I41] => [I51]	2.376023110255176
[I51] => [I41]	2.376023110255176
[I51] => [I11, I72]	2.963709677419355
[I11, I51] => [I72]	1.546283309957924
[I11, I72] => [I51]	2.963709677419355
[I51, I72] => [I11]	2.592592592592592
[I41] => [I11, I51]	2.376023110255176
[I11, I41] => [I51]	2.567637877211238
[I51] => [I11, I41]	2.567637877211238
[I11, I51] => [I41]	2.376023110255176
[I41, I51] => [I11]	2.592592592592592
[I11] => [I21, I51]	2.592592592592592
[I21] => [I11, I51]	2.699859747545586
[I11, I21] => [I51]	2.699859747545586

5.4 Interfaces pour la classification

l'onglet ci-dessous permet la classification des instances avec les différents algorithmes de classification Naïve Bayésienne ou K-nearest neighbors(KNN) ainsi que la comparaison entre les deux algorithmes.

Informations sur la base

Informations sur les attributs

Histogrammes et graphes

Normalisation

Discretisation

Extraction de motifs

Classification

Classification

Sélectionner la classification

Zone

Surface

Compacité

Longueur du grain

Largeur du grain

Coefficient d'asymétrie

Longueur du sillon

Ensemble de test (N°)

Classifier

Classification naïve

Ensemble de test (N°)

Nombre d'intervalle

Tester

Accuracy

Sensitivity

Specificity

Precision

Rappel

F-Score

Classe

Prédite

TP

FN

FP

TN

1

2

3

Accuracy moyenne

Sensitivity moyenne

Specificity moyenne

Precision moyenne

Rappel moyenne

F-Score moyenne

Matrice de confusion

KNN

Ensemble de test (N°)

Le nombre de voisins "K"

La mesure

Tester

Accuracy

Sensitivity

Specificity

Precision

Rappel

F-Score

Classe

Prédite

TP

FN

FP

TN

1

2

3

Accuracy moyenne

Sensitivity moyenne

Specificity moyenne

Precision moyenne

Rappel moyenne

F-Score moyenne

Matrice de confusion

classification Naïve Bayésienne

Pour classier une instance on doit introduire les items de chaque attribut et le pourcentage de l'ensemble de test (T), et le nombre d'intervalles.

Classification

classification Naïve Bayésienne ▼

113

123

133

143

153

163

172

20

4

Classifier

La classe prédite est : Canadian

classification KNN

Pour classier une instance on doit introduire les valeurs des items de chaque attribut et le pourcentage de l'ensemble de test (T), sélectionnée une mesures de distance (euclidienne ou Manhattan) et introduire le nombre de voisins.

Classification

classification KNN ▼

15.26

14.84

0.871

5.763

3.312

2.221

5.22

20

Manhattan ▼

5

Classifier

La classe prédite est : Canadian
le temps d'exécution est : 27 miliseconds

5.5 Interface pour calculs des mesures d'évaluation d'un algorithme de classification

l'onglet ci-dessous permet les calculs des différentes mesures d'évaluation des classifieurs ainsi que la matrice de confusion pour chaque méthodes.

Classification naive

Ensemble de test (N°)

Nombre d'intervalle

Tester

1

Accuracy	
Sensitivity	
Specificity	
Precision	
Rappel	
F-Score	

Classe	Prédite	
Réelle	TP	FN
	FP	TN

2

Accuracy	
Sensitivity	
Specificity	
Precision	
Rappel	
F-Score	

Classe	Prédite	
Réelle	TP	FN
	FP	TN

3

Accuracy	
Sensitivity	
Specificity	
Precision	
Rappel	
F-Score	

Classe	Prédite	
Réelle	TP	FN
	FP	TN

Accuracy moyenne	
Sensitivity moyenne	
Specificity moyenne	
Precision moyenne	
Rappel moyenne	
F-Score moyenne	

Matrice de confusion

KNN

Ensemble de test (N°)

Le nombre de voisins "K"

La mesure

Tester

1

Accuracy	
Sensitivity	
Specificity	
Precision	
Rappel	
F-Score	

Classe	Prédite	
Réelle	TP	FN
	FP	TN

2

Accuracy	
Sensitivity	
Specificity	
Precision	
Rappel	
F-Score	

Classe	Prédite	
Réelle	TP	FN
	FP	TN

3

Accuracy	
Sensitivity	
Specificity	
Precision	
Rappel	
F-Score	

Classe	Prédite	
Réelle	TP	FN
	FP	TN

Accuracy moyenne	
Sensitivity moyenne	
Specificity moyenne	
Precision moyenne	
Rappel moyenne	
F-Score moyenne	

Matrice de confusion

Classification naive

20

4

Tester

1

Accuracy	0,967
Sensitivity	1,000
Specificity	0,967
Precision	1,000
Rappel	0,900
F-Score	0,947

Classe	Prédite	
Réelle	18	2
	0	40

2

Accuracy	0,983
Sensitivity	0,975
Specificity	0,983
Precision	0,952
Rappel	1,000
F-Score	0,976

Classe	Prédite	
Réelle	20	0
	1	39

3

Accuracy	0,983
Sensitivity	0,975
Specificity	0,983
Precision	0,952
Rappel	1,000
F-Score	0,976

Classe	Prédite	
Réelle	20	0
	1	39

Accuracy moyenne	0,978
Sensitivity moyenne	0,967
Specificity moyenne	0,983
Precision moyenne	0,968
Rappel moyenne	0,967
F-Score moyenne	0,966

Matrice de confusion

KNN

20

5

Euclidienne

Tester

1

Accuracy	0,933
Sensitivity	0,975
Specificity	0,933
Precision	0,944
Rappel	0,850
F-Score	0,895

Classe	Prédite	
Réelle	17	3
	1	39

2

Accuracy	0,983
Sensitivity	0,975
Specificity	0,983
Precision	0,952
Rappel	1,000
F-Score	0,976

Classe	Prédite	
Réelle	20	0
	1	39

3

Accuracy	0,950
Sensitivity	0,950
Specificity	0,950
Precision	0,905
Rappel	0,950
F-Score	0,927

Classe	Prédite	
Réelle	19	1
	2	38

Accuracy moyenne	0,956
Sensitivity moyenne	0,933
Specificity moyenne	0,967
Precision moyenne	0,934
Rappel moyenne	0,933
F-Score moyenne	0,932

Matrice de confusion

6 Résultats

D'après les différentes exécutions des algorithmes d'extraction des items fréquents on peut remarquer que l'algorithme Eclat est plus rapide que Apriori (La mesure du temps est en milliseconds). Le tableau ci-dessous montre les résultats.

Support min et intervalles	Apriori	Eclat
sup-min = 20 et int = 4	47	16
sup-min = 30 et int = 4	29	12
sup-min = 40 et int = 4	11	8
sup-min = 20 et int = 5	34	18
sup-min = 30 et int = 5	13	9
sup-min = 40 et int = 5	8	7

TABLE 1 – Tableau comparatif entre les algorithmes des items fréquents Apriori et Eclat

Par rapport aux résultats de la classification on remarque la classification Naïve Bayésienne et plus rapide que la classification KNN.

Conclusion

A partir de ce projet nous avons pu mieux comprendre les taches du data mining les plus importantes tel que l'extraction de règles d'associations et la classification. Ce projet nous a permis d'appliquer les algorithmes d'extraction des motifs fréquents A-Priori et Eclat, ainsi que la classification supervisé avec les deux algorithmes Naïve Bayésienne et KNN, voir l'utilité et l'importance de ces derniers et comprendre les avantages et inconvénients de chaque méthode.