



MS SIO

Introduction au Machine Learning

Compte rendu de Travaux Pratiques

Jérémy CAO

William AFONSO





Code source accessible à :
<https://github.com/will-afs/IML>

01

Analyse descriptive de notre ensemble de données

02

Entraînement d'un modèle via la méthode KNN

03

Entraînement d'un arbre de décisions

04

Conclusion

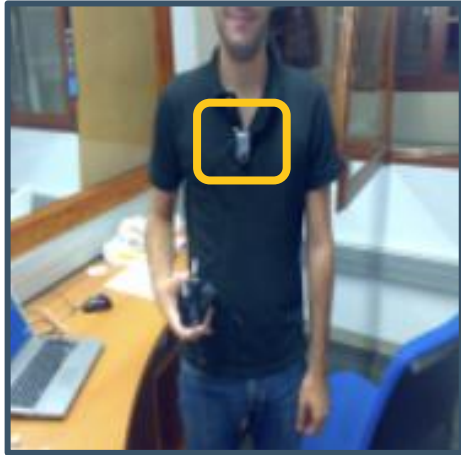


1. Analyse descriptive de notre ensemble de données

1. Analyse descriptive de notre ensemble de données

4

1. Présentation des données brutes



Accéléromètre (x, y, z)

- 1 - Working at Computer
- 2 - Standing up
- 3 - Standing
- 4 - Walking
- 5 - Going up/down stairs
- 6 - Walking + talking with someone
- 7 - Talking while standing

Activités

Mesure n°	Accélération x	Accélération y	Accélération z	N° activité
1	1800	2200	1899	3
Mesure n°	Accélération x	Accélération y	Accélération z	N° activité
1	1502	2215	2153	1
...	-	-	-	-
53440	1878	2277	1988	3
...	-	-	-	-
162501	1550	1600	1455	7

15 individus

Jeu de données brut

Objectif :

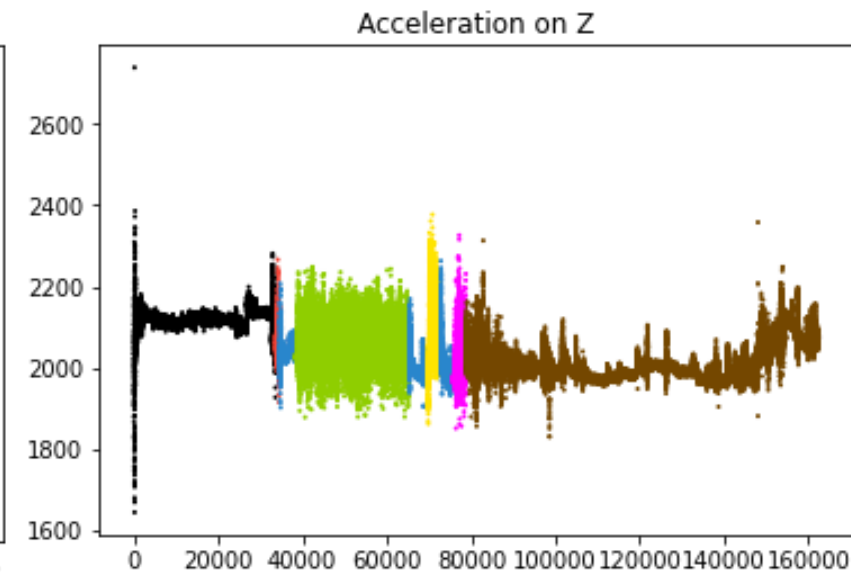
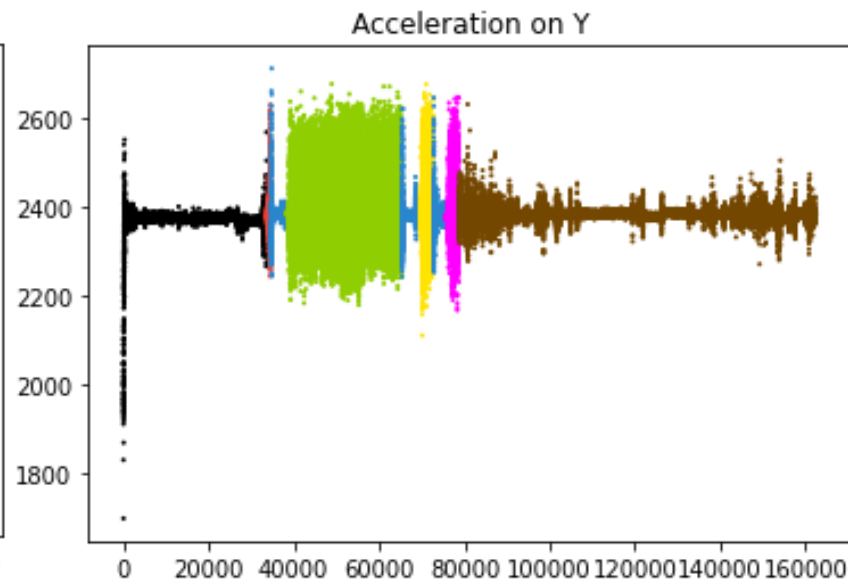
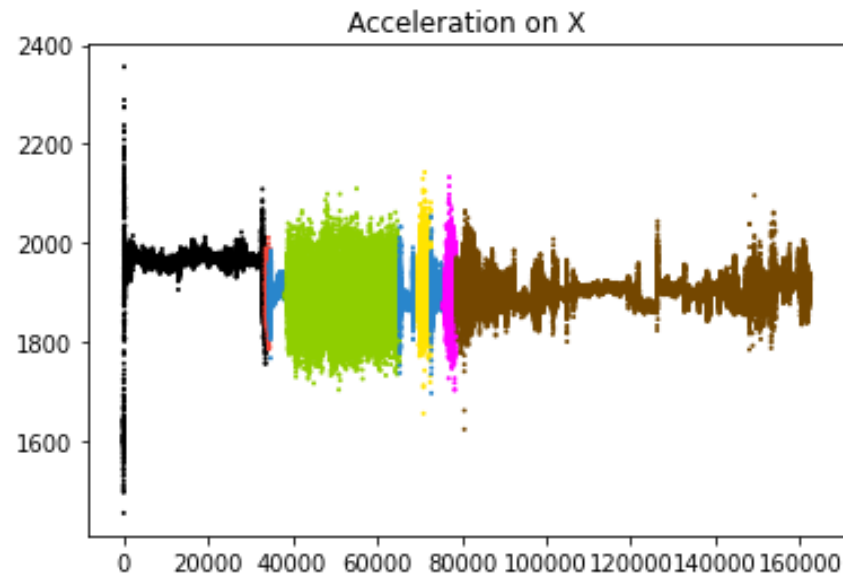
Raffiner nos données pour les rendre exploitables

Trouver l'algorithme le plus approprié pour déduire l'activité à partir d'accélération (x, y, z)

1. Analyse descriptive de notre ensemble de données

2. Analyse des données brutes et commentaires

5

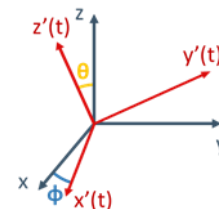


Activities

- Working at a computer
- Standing up, walking and going up\down stairs
- Standing
- Walking
- Going up\down stairs
- Walking and talking with someone
- Talking while standing



3,7 g au repos?!



Système désaxé



Régimes transitoires



Mauvaise répartition



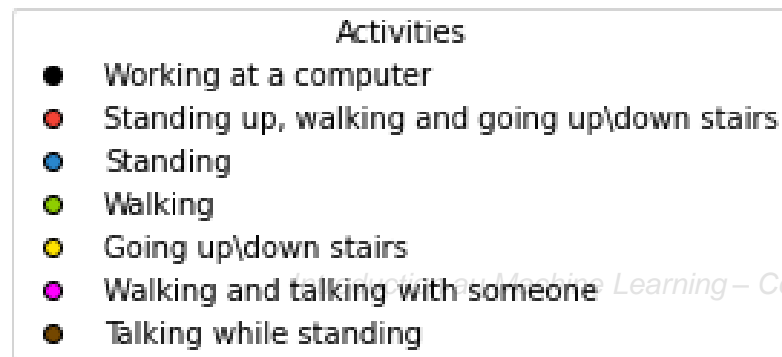
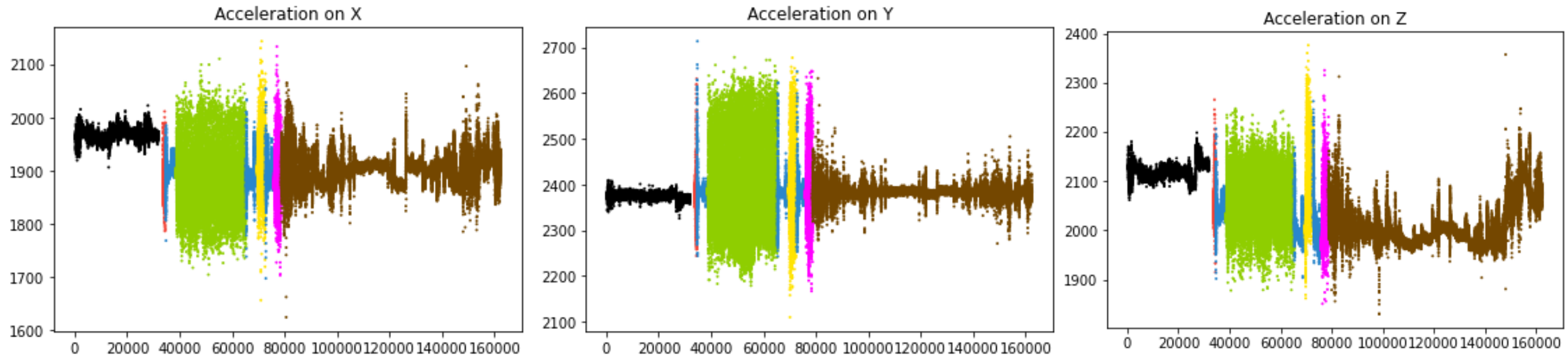
Mouvements chaotiques ?
Périodiques ?

1. Analyse descriptive de notre ensemble de données

3. Suppression des régimes transitoires sur l'activité 1

6

Suppression des transitoires sur l'activité 1



1. Analyse descriptive de notre ensemble de données

4. Augmentation de l'échantillon de données

Mesure n°	Accélération x	Accélération y	Accélération z	N° activité
1	1502	2215	2153	1
...	-	-	-	-
53440	1878	2277	1988	3
...	-	-	-	-
162501	1550	1600	1455	7
...				
Mesure n°	Accélération x	Accélération y	Accélération z	N° activité
1	1897	2198	2110	1
...	-	-	-	-
53440	1877	2455	2322	3
...	-	-	-	-
145099	1887	2110	2232	7

Individu 1

Individu 15

Datasets concaténés



2. Prédiction avec la méthode KNN

2. Prédiction avec la méthode KNN

9

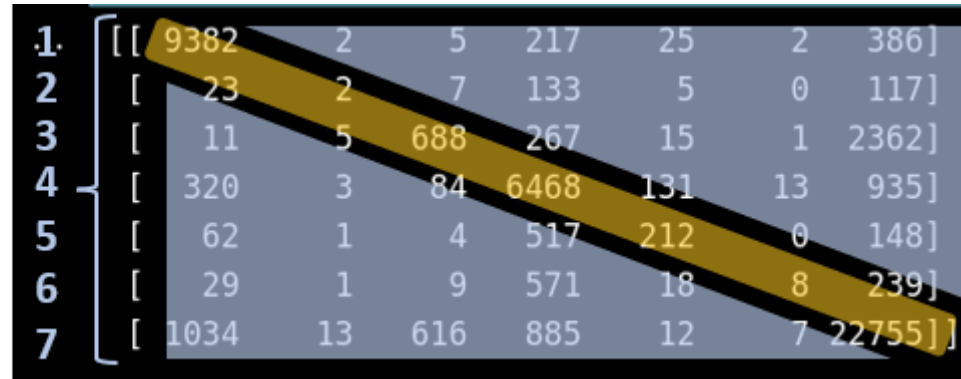
1. Comparaison de la précision du modèle avant et après data-processing

$K_{opt} = 9$

Suppression des transitoires sur l'activité 1

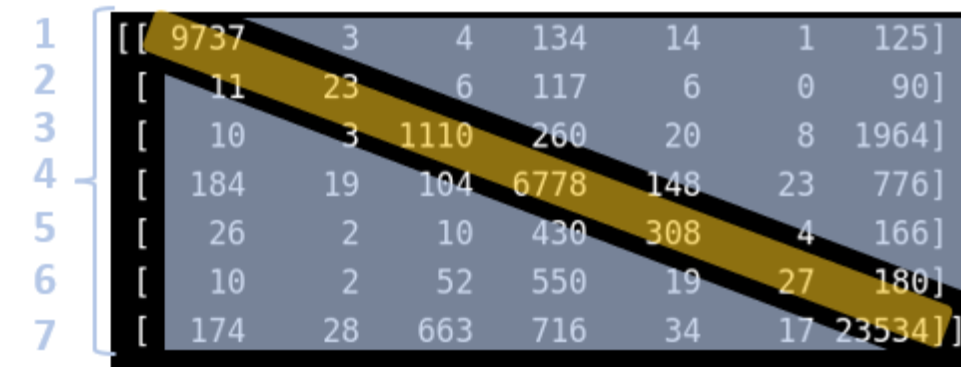
Avant : 85%

Activités
réelles



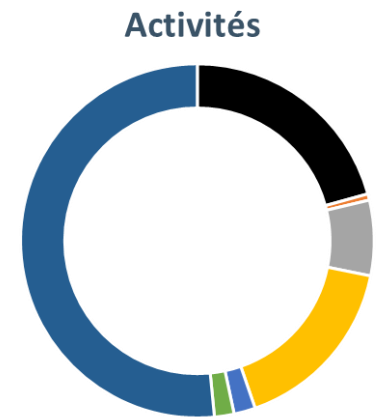
Après : 87%

Activités
réelles



1 2 3 4 5 6 7

Activités prédites par le modèle



- 1 - Working at Computer
- 2 - Standing up
- 3 - Standing
- 4 - Walking
- 5 - Going up/down stairs
- 6 - Walking + talking with someone
- 7 - Talking while standing

2. Prédiction avec la méthode KNN

1. Comparaison de la précision du modèle avant et après data-augmentation

$K_{opt} = 9$

Concaténation des données de chaque individu

Avant : 85 %

Activités réelles

1.	[9382	2	5	217	25	2	386]
2	[23	2	7	133	5	0	117]
3	[11	5	688	267	15	1	2362]
4	[320	3	84	6468	131	13	935]
5	[62	1	4	517	212	0	148]
6	[29	1	9	571	18	8	239]
7	[1034	13	616	885	12	7	22755]

Après : 77 %

1	[166142	1465	1973	5470	167	65	7342]
2	[5534	2885	857	2925	43	42	2103]
3	[3843	381	31033	13931	1361	738	13628]
4	[7823	623	8271	77083	672	389	12125]
5	[1368	50	3316	6084	1886	267	2417]
6	[670	57	1976	2845	475	3028	5184]
7	[7961	555	9267	12862	590	1883	145299]

1 2 3 4 5 6 7

Activités prédites par le modèle



3. Prédiction avec la méthode d'arbre de décisions

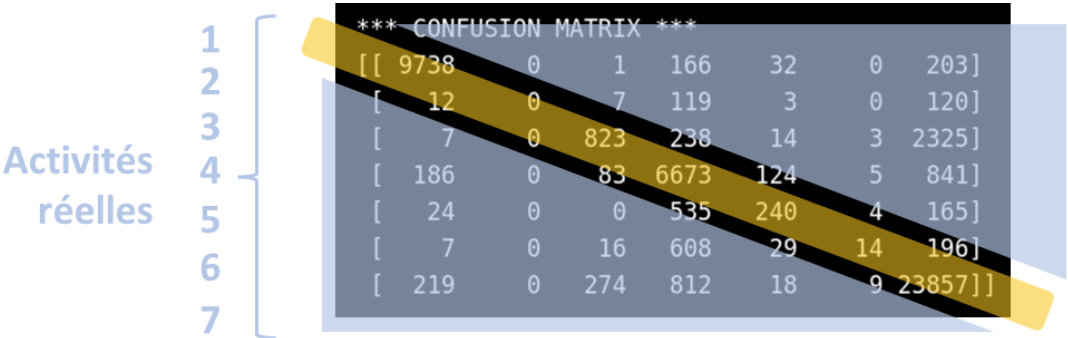
3. Construction d'un arbre de décisions

1. Comparaison de la précision du modèle avant et après data-processing

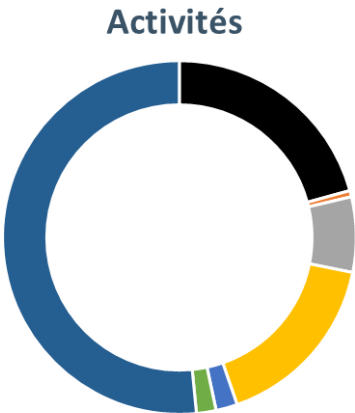
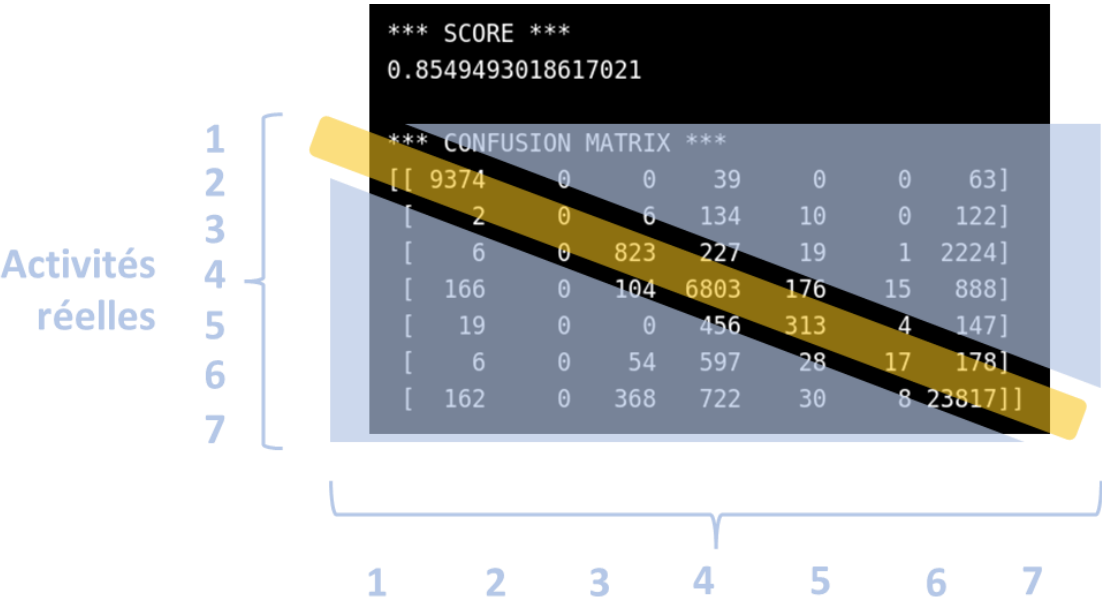
Kleaf = +- 209

Suppression des transitoires sur l'activité 1

Avant : 85-%



Après : 85+%



- 1 - Working at Computer
- 2 - Standing up
- 3 - Standing
- 4 - Walking
- 5 - Going up/down stairs
- 6 - Walking + talking with someone
- 7 - Talking while standing

3. Construction d'un arbre de décisions

13

1. Comparaison de la précision du modèle avant et après data-augmentation

Kleaf = +- 209

Concaténation des données de chaque individu

Avant : 85-%

Activités
réelles

*** CONFUSION MATRIX ***							
1	[9738	0	1	166	32	0	203]
2	[12	0	7	119	3	0	120]
3	[7	0	823	238	14	3	2325]
4	[186	0	83	6673	124	5	841]
5	[24	0	0	535	240	4	165]
6	[7	0	16	608	29	14	196]
7	[219	0	274	812	18	9	23857]]

Après : 71%

Activités
réelles

*** CONFUSION MATRIX ***							
1	[161536	0	2649	6358	7	47	12141]
2	[7268	0	743	3132	3	11	3328]
3	[4169	0	25912	13743	1225	470	19412]
4	[10404	0	7452	70884	23	59	17934]
5	[1916	0	2844	6357	1393	96	3014]
6	[812	0	1327	2967	762	972	7483]
7	[9369	0	7369	14128	260	748	146227]]

1 2 3 4 5 6 7

Activités prédites par le modèle

Conclusion

MERCI !

Des questions?

ANNEXES

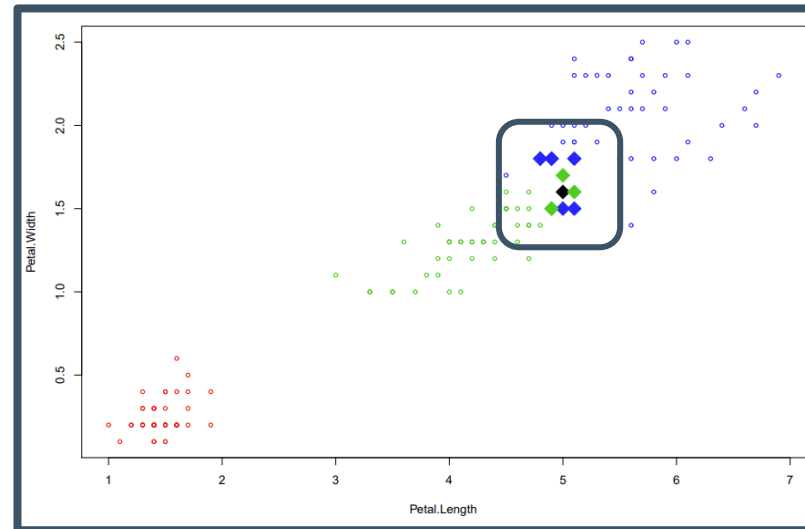
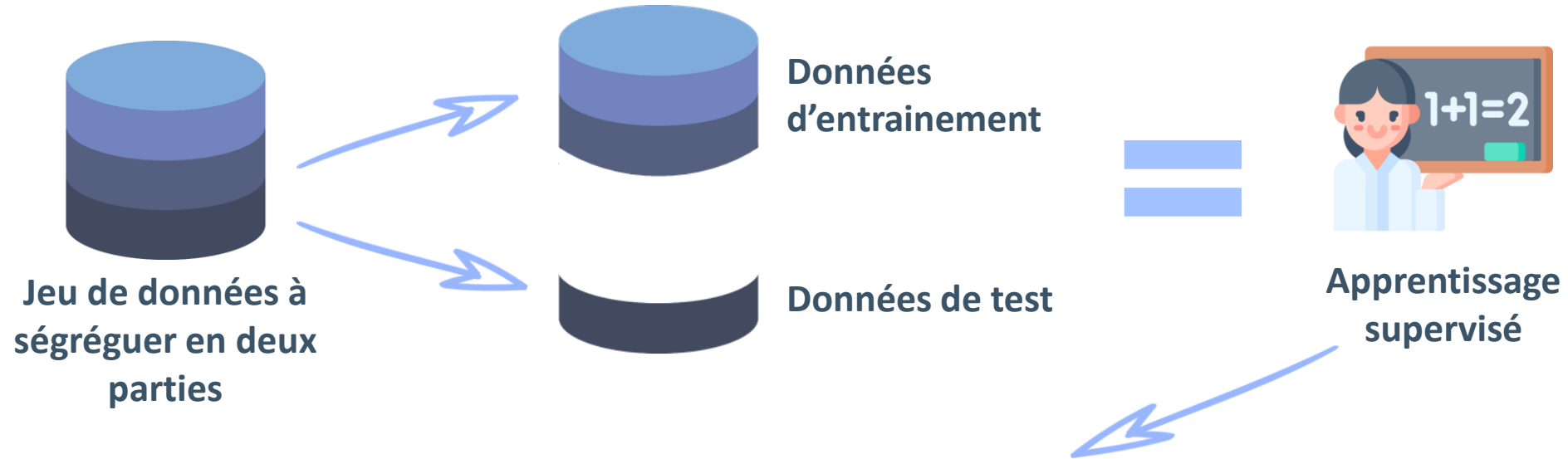


Synthèse du cours

Synthèse du cours

1. Méthode KNN

18



On prédit la caractéristique du point à partir de ses K voisins les plus proches

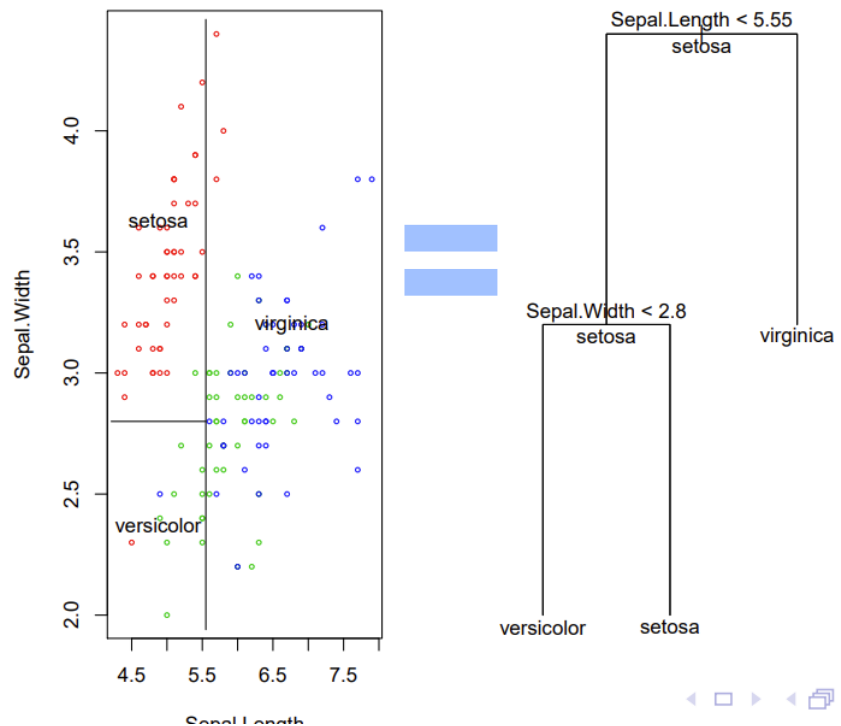


Attention aux données et au paramétrage !

2. Arbres de décision

Apprentissage supervisé (comme Knn)

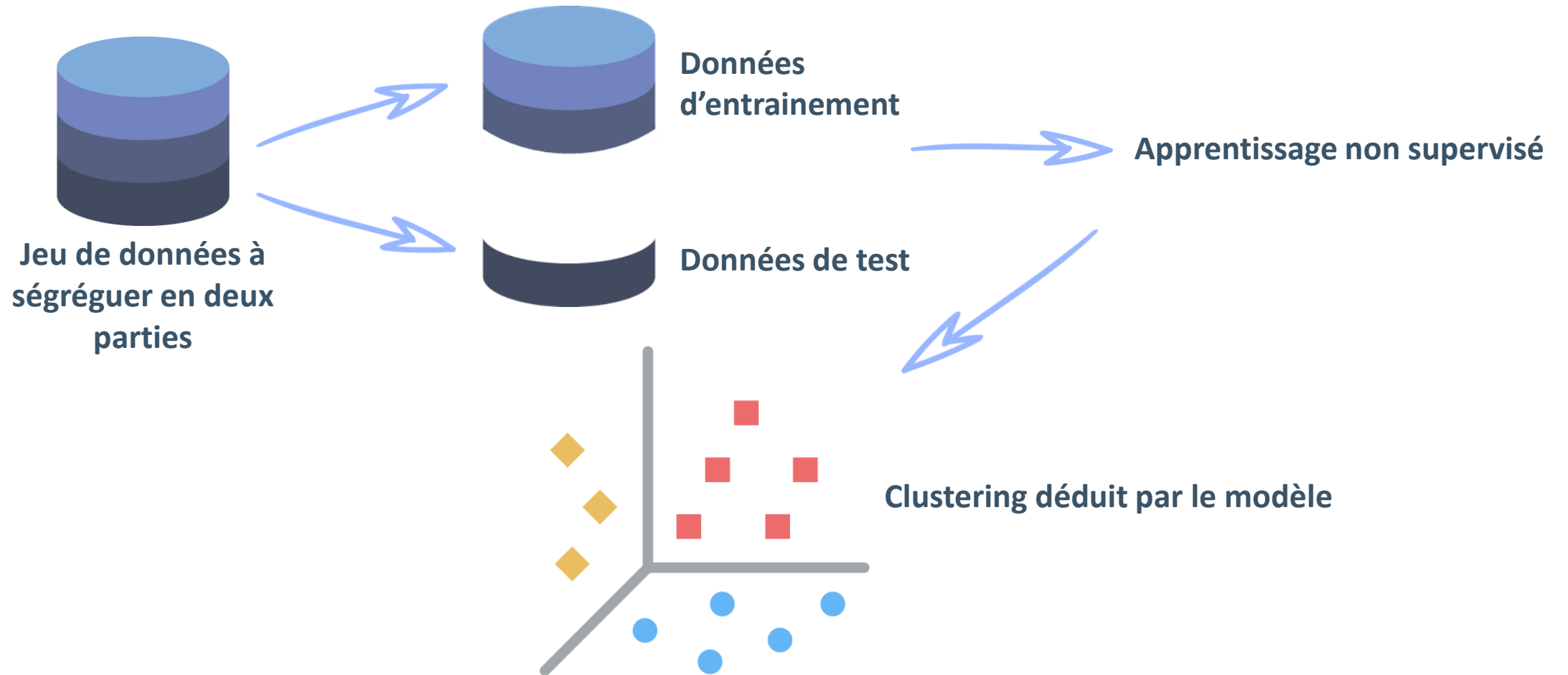
Segmentation de l'espace en sous-espaces où l'entropie est minimale



Puis, à partir des caractéristiques de l'échantillon étudié, l'algorithme pourra prédire sa classe, en s'orientant grâce à l'arbre

Il est possible d'élaguer des branches de l'arbre pour éliminer les moins pertinentes et ainsi améliorer la précision de l'arbre

3. Kmeans



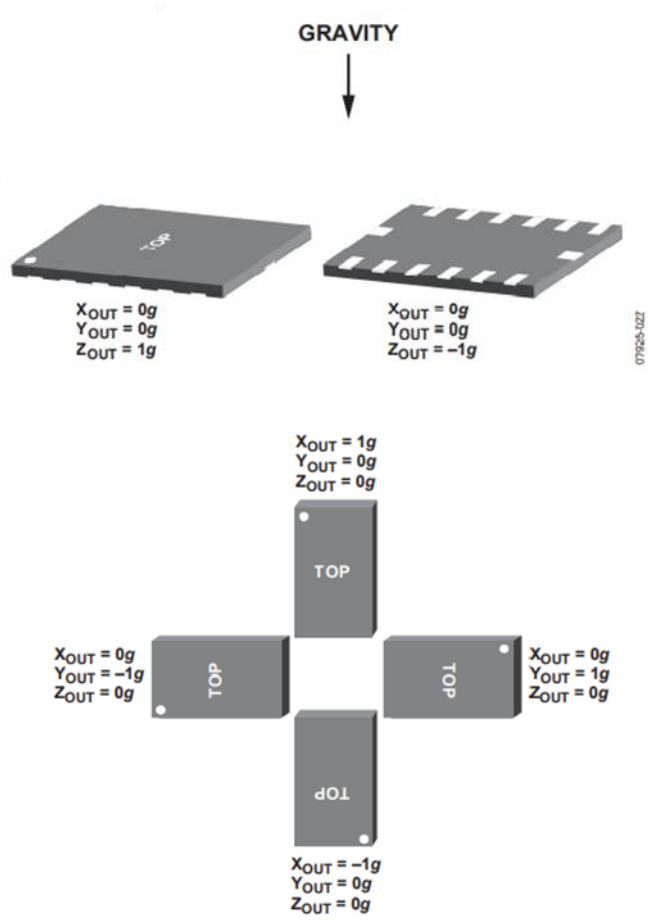


Analyse des données (sources)

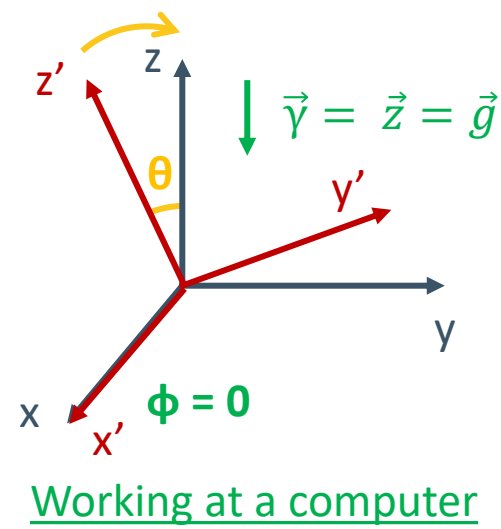
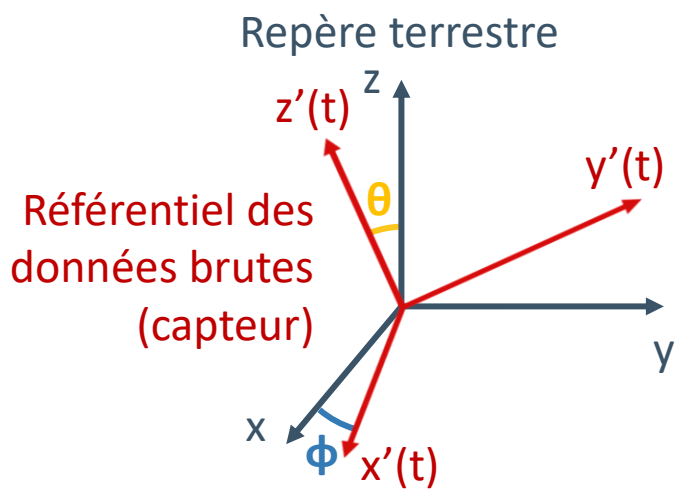
Analyse descriptive de notre ensemble de données

Calcul de l'angle Thêta pour corriger les valeurs

Un problème de repère ?



Output Response vs. Orientation to Gravity



$$(1) \begin{pmatrix} x \\ y \\ z \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & \cos \theta & -\sin \theta \\ 0 & \sin \theta & \cos \theta \end{pmatrix} \begin{pmatrix} x' \\ y' \\ z' \end{pmatrix}$$

$$(2) z = \sqrt{x'^2 + y'^2 + z'^2} \neq g$$

(1)+(2) → Θ

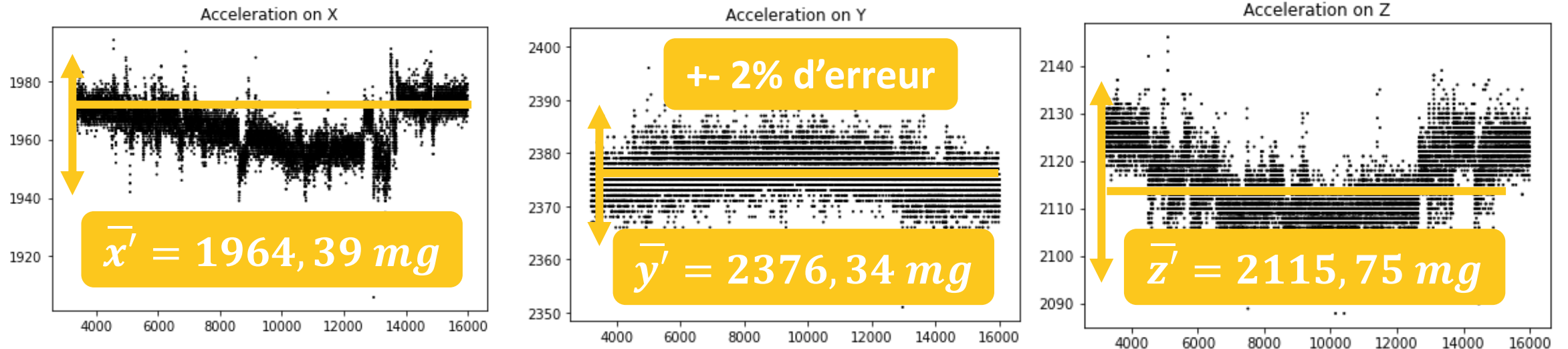
symPy

Analyse descriptive de notre ensemble de données

Calcul de l'angle Thêta pour corriger les valeurs

23

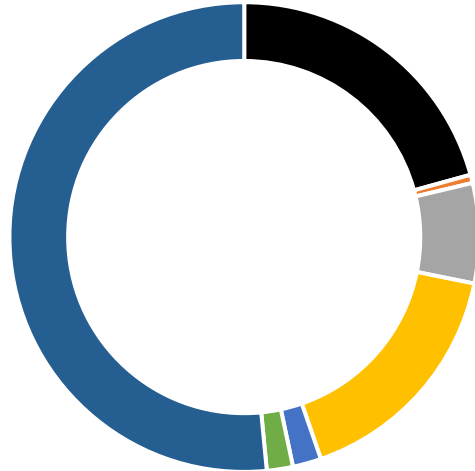
Isolation des valeurs d'accélération sur l'activité 1 (stable)



$$3739,28 = 2376,34 \sin(\bar{\theta}) + 2115,75 \cos(\bar{\theta})$$

Analyse des données - Activités dynamiques très peu représentées

Activités



- Les activités « statiques » sont représentées à + de 75%
- Certaines données sont très peu représentées

- 1 - Working at Computer
- 2 - Standing up
- 3 - Standing
- 4 - Walking
- 5 - Going up/down stairs
- 6 - Walking + talking with someone
- 7 - Talking while standing

Analyse descriptive de notre ensemble de données

Analyse des données brutes

Activité	Moyennes		
	Accélération x	Accélération y	Accélération z
1	1965	2373	2119
2	1892	2376	2067
3	1893	2383	2011
4	1884	2381	2051
5	1918	2373	2103
6	1888	2381	2031
7	1900	2383	2008

Pour toutes les activités, c'est l'accélération en y (resp. x) qui est la plus (resp. moins) élevée

Les activités « statiques » (1, 3, 7) ont des accélérations anormalement élevées

Activité	Ecart-types		
	Accélération x	Accélération y	Accélération z
1	28	22	24
2	33	50	49
3	20	24	33
4	45	88	49
5	49	89	76
6	46	79	59
7	20	9	39

Les activités « dynamiques » (2, 4, 5, 6) – les moins représentées - ont les écarts-types plus élevés



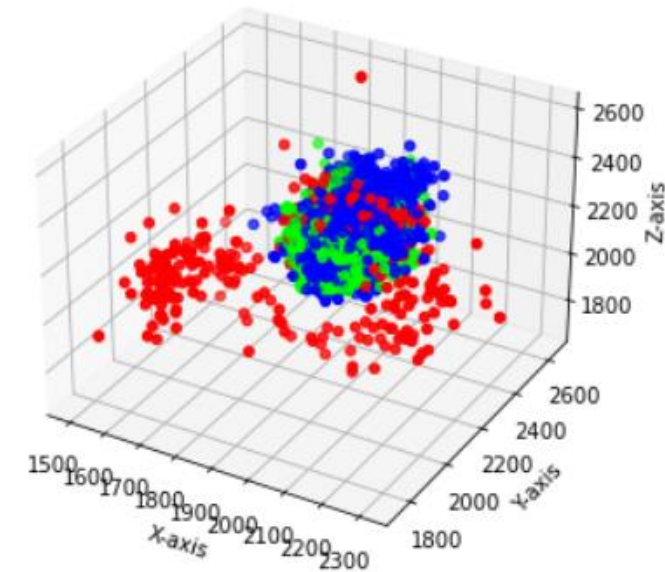
Analyse descriptive de notre ensemble de données

Analyse des données

26

	Accélération x	Accélération y	Accélération z
Moyenne	1910	2380	2041
Ecart-type	41	42	60
Min	1455	1697	1644
Max	1935	2386	2101

Analyse statistique sur l'ensemble des accélérations,
toutes activités confondues

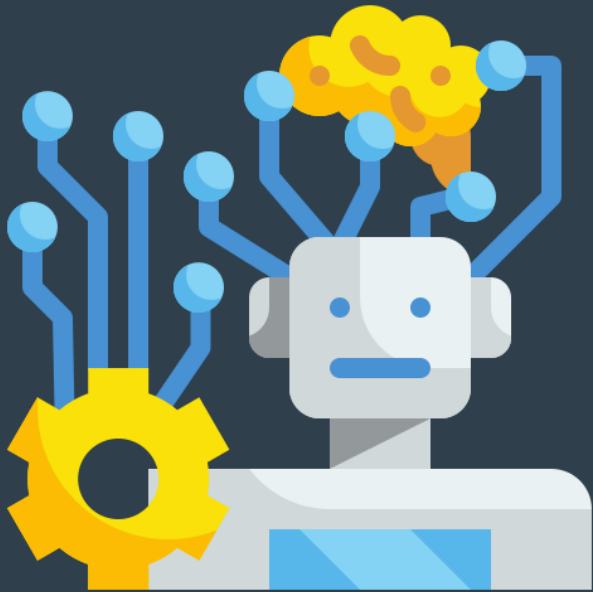


Analyse descriptive de notre ensemble de données

Preprocessing

27

Moyennes des accélérations			
Activité	En x	En y	En z
1	1965	2373	2119
2	1892	2376	2067
3	1893	2383	2011
4	1884	2381	2051
5	1918	2373	2103
6	1888	2381	2031
7	1900	2383	2008



Entraînement d'un modèle de ML avec KNN (snippets)

1. Import des données via Pandas

```
X=pd.read_csv('data/1.csv',header=None,delimiter=',',usecols=[1, 2, 3])  
y=pd.read_csv('data/1.csv',header=None,delimiter=',',usecols=[4])  
  
X = X.values  
y = y.values[:,0]
```

2. Entraînement et test avec exactement les mêmes données

```
KNeighborsClassifier(n_neighbors=10)
```

```
# If we were in regresssion  
# clf = neighbors.KNeighborsRegressor(n_neighbors=) # je précise la méthode et ses hyperparamètres
```

✓ 0.3s

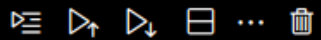
Python

```
# prevision  
clf.predict(X)
```

✓ 7.7s

Python

```
array([0.0000e+00, 1.0000e+00, 2.0000e+00, ..., 3.6985e+04, 1.5683e+05,  
       1.5807e+05])
```



```
clf.score(X,Y) # predict + calcul le score = la précision
```

✓ 7.1s

Python

0.10551384615384615 10% de précision ! C'est très faible ! Overfitting ?

2. Méthode KFold pour compléter les données manquantes

```
# Pour voir l'effet du partage ...

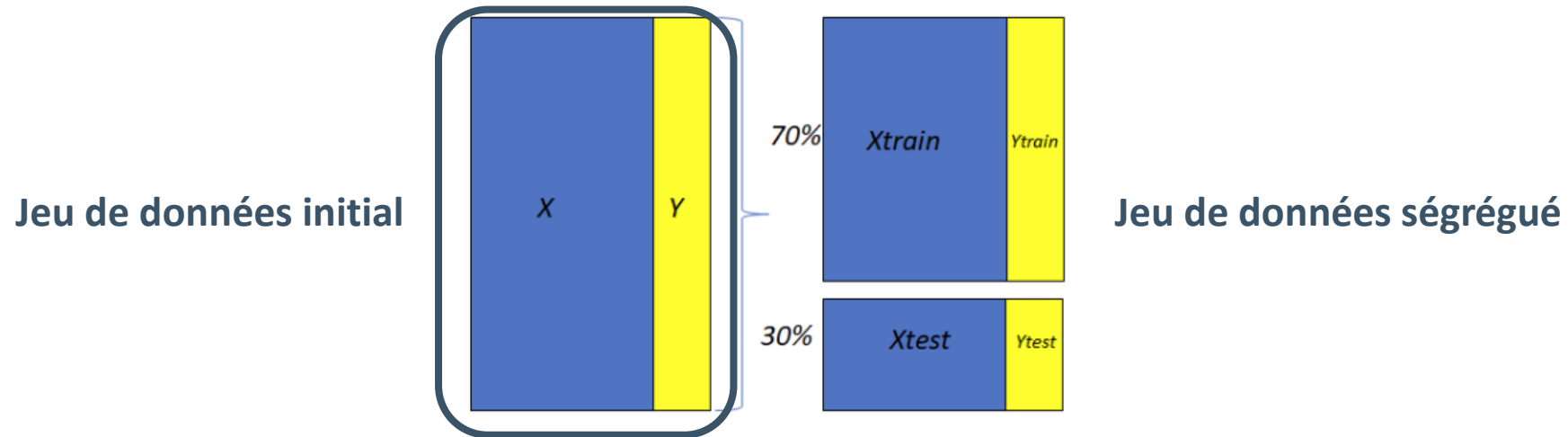
from sklearn.model_selection import KFold
kf=KFold(n_splits=4, shuffle=True) # partages de validation
XXX = X_train[:12,:] # je prends uniquement les 12 premières lignes
#print(XXX)
for learn,test in kf.split(XXX): # boucle sur différents partages de validation
    print("Learn")
    print(learn)
    print("test")
    print(test)
```

[18] ✓ 0.6s

Python

```
... Learn
[ 0  2  3  4  5  6  7  9 10]
test
[ 1  8 11]
Learn
[ 1  3  4  5  6  7  8  9 11]
test
[ 0  2 10]
Learn
[ 0  1  2  4  5  7  8 10 11]
```


3. Préparation des échantillons d'entraînement et de test



```
import random
from sklearn.model_selection import train_test_split
X_train,X_test,Y_train,Y_test = train_test_split(X,y,test_size=0.3,random_state=random.seed())
print(X_train.shape)
print(X_test.shape)
```

On construit cet échantillon sur la base d'un nombre pseudo-aléatoire

70% / 30% classique

[13]

✓ 0.7s

Python

```
... (113750, 2)
    (48750, 2)
```

3. Préparation des échantillons d'entraînement et de test

On entraîne notre modèle avec les données splittées en 70/30

```
n_neighbors=9
clf = neighbors.KNeighborsClassifier(n_neighbors) # je précise la méthode et ses hyperparamètres
clf.fit(X_train, Y_train) # je lance l'apprentissage
prev_test = clf.predict(X_test)
sc_train = clf.score(X_train, Y_train)
sc_test = clf.score(X_test, Y_test)
print(sc_train)
print(sc_test)
```

[15] ✓ 5.1s

Python

...

0.8255648351648351
0.8105641025641026

Le score est intéressant, mais pas non plus imbattable...

4. Evaluation du modèle avec la matrice de confusion

```
# Compute the confusion matrix between kmeans activity label and accelerations
from sklearn.metrics import confusion_matrix

cm=confusion_matrix(Y_test,prev_test)
# the confusion matrix is difficult to read as labels assigned by kmeans are arbitrary
print(cm)
```

[17] ✓ 0.1s

1.	[9382	2	5	217	25	2	386]
2	[23	2	7	133	5	0	117]
3	[11	5	688	267	15	1	2362]
4	[320	3	84	6468	131	13	935]
5	[62	1	4	517	212	0	148]
6	[29	1	9	571	18	8	239]
7	[1034	13	616	885	12	7	22755]

Activités
réelles

Activités correctes
(réel correspond à ce qui a été prédit)

Activités incorrectes
(réel correspond à autre chose que ce qui a été prédit)

1 2 3 4 5 6 7

Activités prédites par le modèle

1. Recherche de l'hyperparamètre K optimal

```
from sklearn.model_selection import KFold
kf=KFold(n_splits=3, shuffle=True) # partages de validation

from sklearn import neighbors
scores=[]
for k in range(1,10): # les différentes valeurs de k à tester
    score=0
    clf=neighbors.KNeighborsClassifier(k)
    for learn,test in kf.split(X_train): # boucle sur différents partages de validation
        X_app=X_train[learn]
        Y_app=Y_train[learn]
        clf.fit(X_app,Y_app)
        X_val=X_train[test]
        Y_val=Y_train[test]
        score+=clf.score(X_val,Y_val)
    scores.append(score)
print(scores)
#plt(scores)
k_opt=scores.index(max(scores)) + 1 # valeur optimale de k
print(k_opt)
```

Python

```
[2.430260421697527, 2.4124207978205026, 2.5024149145857715, 2.5292163036099855, 2.5417375292364217,
2.55765088301482, 2.561503642421665, 2.5670314489655865, 2.567910907329197]
```

9 Index du score maximal

Score maximal



Entraînement d'un modèle d'arbre de décisions

1. Construction simple

Même chose qu'avec KNN

```
import random
from sklearn.model_selection import train_test_split
from sklearn.metrics import confusion_matrix

X=pd.read_csv('data/3.csv',header=None,delimiter=',',usecols=[1, 2, 3]).values
Y=pd.read_csv('data/3.csv',header=None,delimiter=',',usecols=[4]).values[:,0]

X_train,X_test,Y_train,Y_test = train_test_split(X,Y,test_size=0.3,random_state=random.seed())
clf=tree.DecisionTreeClassifier(max_leaf_nodes=199,criterion='entropy')
clf.fit(X_train,Y_train)
prev_test = clf.predict(X_test)
#print(prev)
scoretree = clf.score(X_test,Y_test)
# Compute the confusion matrix between kmean label and iris types

cm=confusion_matrix(Y_test,prev_test)
# the confusion matrix is difficult to read as labels assigned by kmeans are arbitrary
print(cm)
print(scoretree)
```

On peut construire notre DTC selon différents critères

On configure le nombre de feuilles maximal (ici au doigt mouillé)

Même chose qu'avec KNN

✓ 1.6s

Python

2. Recherche de l'hyperparamètre K (nombre de feuilles) optimal

```
from sklearn import neighbors
scores=[]
for k in range(190,210): # les différentes valeurs de k à tester
    score=0
    clf=tree.DecisionTreeClassifier(max_leaf_nodes=k,criterion='entropy')
    for learn,test in kf.split(X_train): # boucle sur différents partages de validation
        X_app=X_train[learn]
        Y_app=Y_train[learn]
        clf.fit(X_app,Y_app)
        X_val=X_train[test]
        Y_val=Y_train[test]
        score+=clf.score(X_val,Y_val)
    scores.append(score)
print(scores)
#plt(scores)
k_opt=scores.index(max(scores)) + 1 # valeur optimale de k
print(k_opt)
```

✓ 6.2s

Python

```
[2.5648120535533594, 2.563806944251248, 2.560624199841145, 2.5633881329256107, 2.56853901822021, 2.565398323471006,
2.565691397927982, 2.5633461639338417, 2.563555691477194, 2.5664033556113415, 2.5655238534058125, 2.56732454095933,
2.5637231097346387, 2.566193915751826, 2.5682039449589613, 2.5657332985263572, 2.5659846898408736,
2.569292811615906, 2.5637232324920105, 2.56481198340629]
```

18

Le nombre de feuilles optimal vaut donc $190+18 = 208$