

# Analyses et résultats des données de l'égypte : Le fichier Egypte\_period\_prediction.csv

chahrazed labba

le 17 Mai 2023

## 1 Prédire la période temporelle avec 17 étiquettes

Cette section décrit les données utilisées et présente les résultats des prédictions obtenues en utilisant les algorithmes de machine learning suivants : RF, svm et xGboost.

### 1.1 Les données

- Nombre d'individus : 315
- Les caractéristiques (features) : 'classe\_age', 'Position\_corps', 'Orientation\_tete', 'Orientation\_Face', 'Position\_mains', 'Nb\_actes', 'Mobilier', 'Nb\_Mobilier', 'mob\_tete', 'mob\_membre\_inf', 'mob\_thorax', 'mob\_dos', 'mob\_devant\_corps', 'mob\_contre\_jarre', 'mob\_dans\_jarre', 'mob\_exterieur\_coffre', 'mob\_dans\_coffre', 'mob\_contre\_paroie', 'mob\_dans\_remplissage', 'Signalisation', 'pelle', 'Ornement', 'Nb\_Ornement', 'orn\_Porte', 'orn\_Nonporte', 'Coquillage', 'Aspatharia\_unio', 'Perles', 'Fard', 'Outils', 'Prosterne', 'Decoupe'

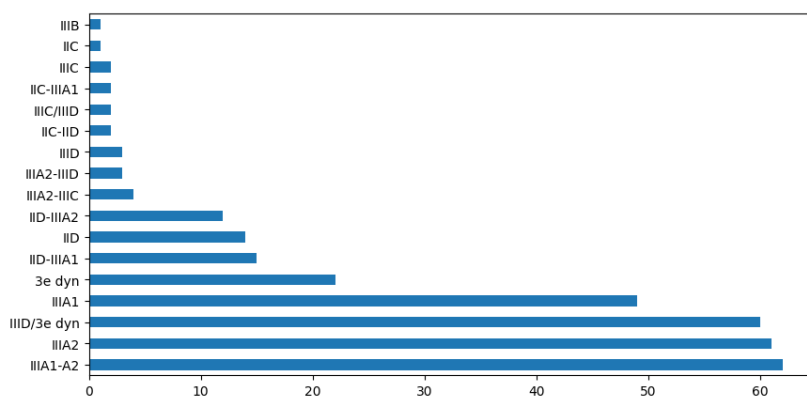


Figure 1: Nombre des individus par class label

La répartition des données entre les différentes étiquettes de classe présente un énorme déséquilibre. Certaines étiquettes de classe ne contiennent qu'un seul échantillon.

	A totale	0	1	2	3	4	5	6	7	8	9	10	11	12	13	15	16
RF (%)	37.14	75.0	53.22	42.62	22.44	9.09	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
SVM (%)	33.01	58.33	45.16	29.50	20.40	27.27	6.66	21.42	25.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
Xgboost(%)	35.87	66.66	43.54	36.06	32.65	22.72	0.0	21.42	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0

Table 1: Accuracy globale et accuracy par class label

les étiquettes de classe et leur signification

- IID/3e dyn:0
- 'IIIA1-A2':1
- 'IIIA2':2
- 'IIIA1':3
- 3e dyn':4
- IID-IIIA1':5
- 'IID':6
- 'IID-IIIA2':7
- IIIA2-IIIC':8
- 'IIID':9
- 'IIIA2-IIID':10
- 'IIC':11
- 'IIC/IIID':12
- 'IIC-IID':13
- IIC-IIIA1':14
- 'IIC':15
- 'IIIB':16

Pour les trois modèles, les étiquettes de classe de 8 à 16 n'ont jamais été prédites. Ce problème est dû à un problème de représentation des données (sous-représentation par rapport au reste des étiquettes de classe).

- df\_results\_SVM\_17\_labels\_LOO.csv (contient les prédictions faites par SVM)
- df\_results\_RF\_17\_labels\_LOO.csv (contient les prédictions faites par RF)
- df\_results\_xgboost\_17\_labels.csv (contient les prédictions faites par Xgboost)

## 2 Prédire la période temporelle avec 5 étiquettes

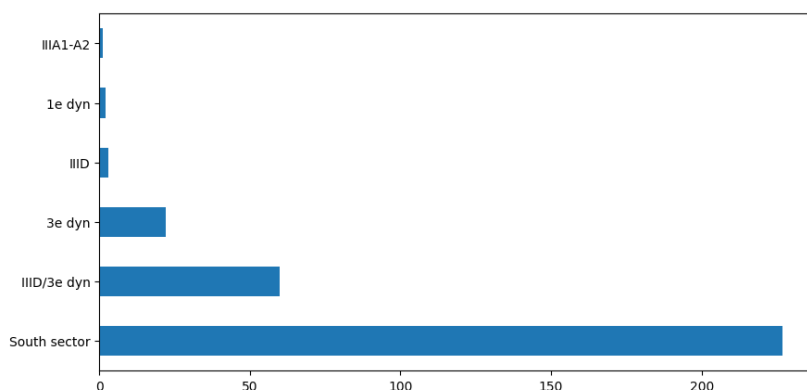


Figure 2: Nombre des individus par class label

La répartition des données entre les différentes étiquettes de classe présente un énorme déséquilibre. Certaines étiquettes de classe ne contiennent qu'un seul échantillon.

- 'South sector':0
- 'IIID/3e dyn':1
- '3e dyn':2
- 'IIID':3
- '1e dyn':4
- 'IIIA1-A2' : 5

	Accuracy	0	1	2	3	4	5
RF	78.41	99.55	35.0	0.0	0.0	0.0	0.0
SVM	77.46	92.95	45.0	27.27	0.0	0.0	0.0
Xgboost	78.41	94.27	46.66	22.72	0.0	0.0	0.0

Table 2: Accuracy globale et accuracy par class label

## 3 Résultat de clustering en utilisant k-means

Le clustering donne lieu à deux clusters.

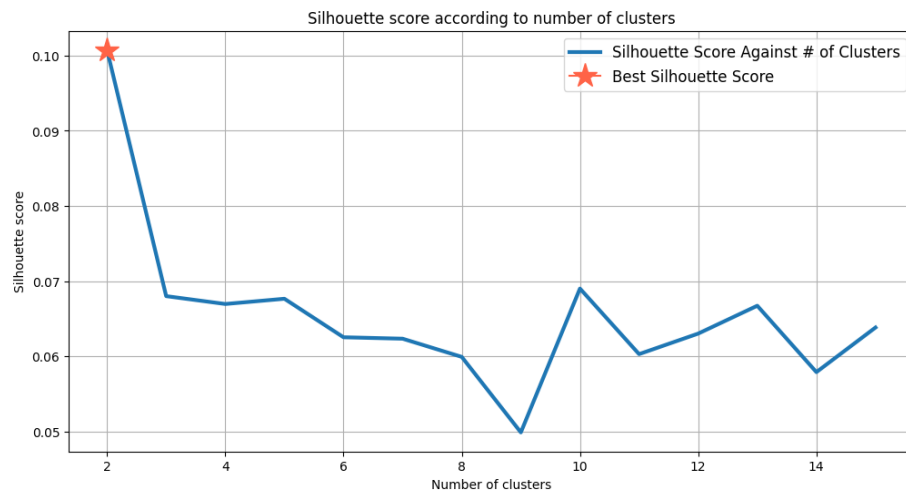


Figure 3: Nombre de clusters

Comme K-means est une boîte noire et que nous ne pouvons pas savoir quelles caractéristiques ont été utilisées pour réaliser le clustering, j'ai appliqué un algorithme de classification (Random Forest), puis la lib Shap pour expliquer le clustering.

Le fichier Egypt\_Cluster\_Data\_315\_V.csv contient les résultats de clustering.

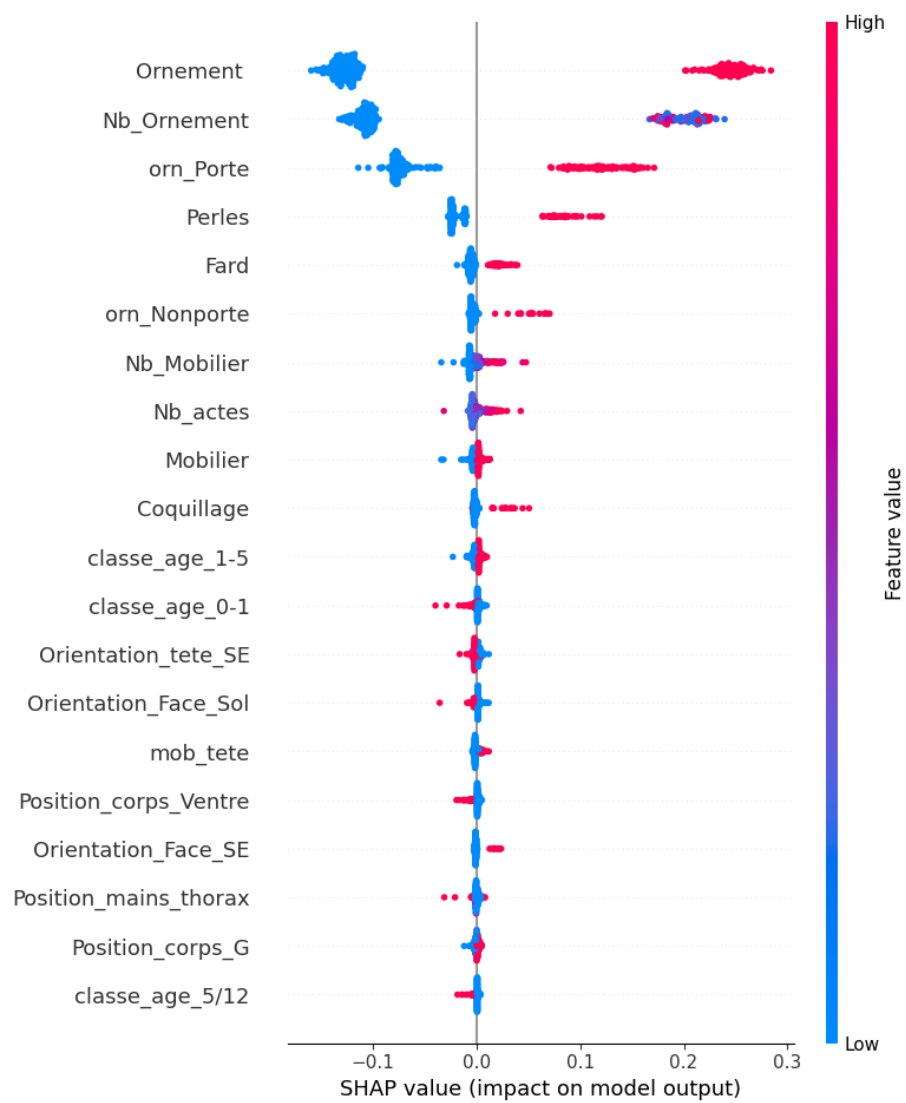


Figure 4: Cluster 0

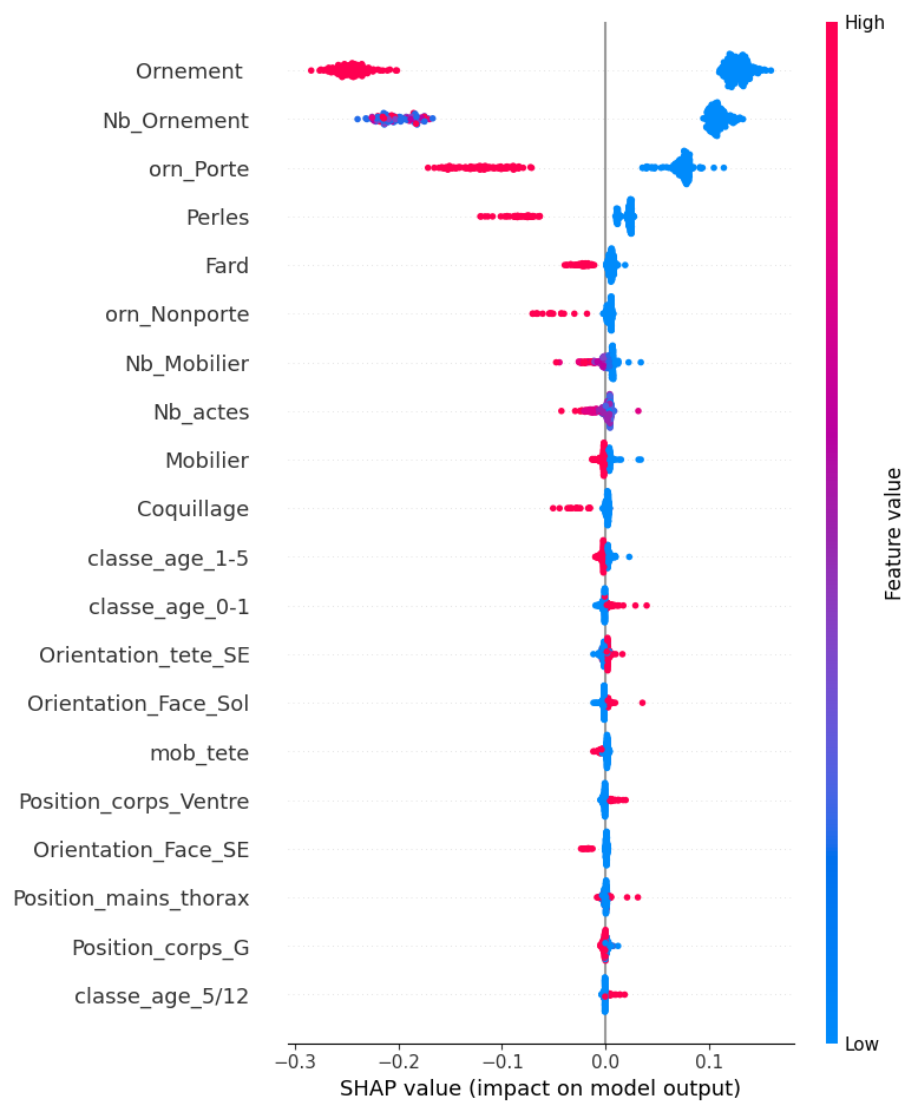


Figure 5: Cluster 1