

Rapport TP Introduction à l'IA: Cas de classification

Harouna KINDO

September 2024

Contents

1	Introduction	3
2	Préparation du dataset	3
2.1	Chargement et inspection des données	3
2.2	Séparation des features et des labels	5
2.3	Standardisation des données	6
3	Exploration des données	6
3.1	Distribution des classes	6
3.2	Distribution des features	7
3.3	Visualisation des données	8
4	Application des algorithmes de classification	8
4.1	Régression Logistique	8
4.1.1	1er modèle:	8
4.1.2	Nouveau Modèle avec Top 1,2, 3 features:	9
4.1.3	Etude de cas	12
4.2	SVM (Support Vector Machine)	12
4.2.1	1er modèle	12
4.2.2	Test du modèle avec top1,2,3 feautres	13
4.2.3	Etude de cas	14
5	Partie autonomie	14
5.1	K-NN	14
5.2	Decision tree et Random Forrest	15
6	Conclusion	16

1 Introduction

Dans ce rapport, nous appliquons des méthodes de classification supervisée sur un jeu de données portant sur l'expression des gènes du cancer du côlon. L'objectif est de prédire si un échantillon appartient à la classe "normal" ou "tumoral". Donc c'est un cas de classification binaire alors nous allons utiliser des modèles de classification pour nos prédictions. Mais avant d'en arriver aux prédictions, nous allons d'abord procéder à l'analyse exploratoire de notre jeu de données.

2 Préparation du dataset

Les étapes de préparation des données incluent :

- Chargement des données à l'aide de **pandas**.
- Vérification et nettoyage des données.
- Séparation des données en ensembles d'entraînement et de test.

Pour cette phase de préparation et aussi pour l'exploration de nos données notamment avec des graphiques, nous aurons besoin des bibliothèques **pandas**, **matplotlib** et **seaborn**.

```
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
```

Figure 1: Importation des bibliothèques pour préparation et exploration des données

2.1 Chargement et inspection des données

Comme notre dataset est un fichier csv, nous sommes parvenu à la lire grâce à la fonction `read_csv` de **pandas**. Le paramètre `sep=","` indique à la fonction la nature du séparateur utilisé dans la dataset pour distingé les lignes.

```
[2]: dataset = pd.read_csv("/kaggle/input/cancer/colon_cancer.csv", sep=";")

dataset.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 804 entries, 0 to 803
Data columns (total 62 columns):
 #   column      Non-Null Count  Dtype
---  -
 0   id_sample    804 non-null      object
 1   ADH1C        804 non-null      float64
 2   DHRS11       804 non-null      float64
```

Figure 2: Chargement puis inspection des données

Cette partie du code nous montre la lecture de la dataset puis des informations sur le contenu de cette dataset. En effet, notre dataset compte 804 observations 60 features qui sont toutes numériques.

```
[4]: dataset.head(10)
```

	IS	DAO	NIBAN1	PRUNE2	FOX2	...	BSPRY	CDHR2	ERRF1	CLICS	PLL	GAL	CRYL1	YBX2	ANGPTL4	tissue_status
38	8.465133	1.174665	0.975086	3.028995	...	5.528372	3.500730	2.893760	5.130743	1.485110	1.830484	6.379933	3.919293	3.031413		normal
34	8.159814	5.959414	1.625518	2.984629	...	5.671788	4.658790	4.715374	5.438104	2.131466	2.097157	6.863173	0.322829	3.978531		normal

Figure 3: Affichage des données

Cette affichage des données nous permet de confirmer que les fetures sont toutes numériques cependant la target est catégoricielle donc nous allons devoir l'encoder en binaire.

```
[5]: #Vérification de valeurs manquantes
dataset.isnull().sum()

[5]: id_sample    0
ADH1C            0
DHRS11           0
UGP2             0
SLC7A5           0

[6]: #Vérification de données dupliquées
sum(dataset.duplicated(subset = 'id_sample')) == 0

[6]: True
```

Figure 4: Vérification des données manquantes et dupliquées

Nous remarquons qu'il n'y a ni données manquantes ni données dipluquées dans notre dataset.

```
[8]: new_data=dataset.copy()

[9]: new_data['status'] = le().fit_transform(dataset['tissue_status'])

[10]: new_data.head(5)
```

	TSS	DAO	NIBAN1	PRUNE2	FOXF2	...	CDHR2	ERRF1	CLIC3	PLP	GAL	CRYL1	YBX2	ANGPTL4	tissue_status	status
9368	8.465133	1.174665	0.975086	3.028995	...	3.500730	2.893760	5.130743	1.485110	1.830484	6.379933	3.919293	3.031413		normal	0
3334	8.159814	5.959414	1.625518	2.984629	...	4.658790	4.715374	5.438104	2.131466	2.097157	6.863173	0.322829	3.978531		normal	0

Figure 5: Encodage de la target

Grâce à la fonction `LabelEncoder` de la bibliothèque `Sklearn`, nous sommes parvenu à transformer la target qui était une variable catégorielle en variable numérique. Nous rappelons que la fonction `LabelEncoder` remplace l'ensemble des classe de la variable par des nombre en commençant par 0. Vu que nous avons que deux classes alors nous avons que des 0 et 1 dans notre target maintenant.

Mais avant de faire cette transformation, nous nous assurer de copier notre dataset dans une nouvelle variable. Ainsi, nous aurons toujours une version de la dataset initiale.

2.2 Séparation des features et des labels

Les données sont séparées entre les features (X) et la variable cible (Y) :

```
[12]: from sklearn.model_selection import train_test_split
      Y=new_data["status"]
      X=new_data.drop(["id_sample", "tissue_status", "status"], axis=1)

[13]: X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=0.2, random_state=42)
```

Figure 6: séparation de la dataset

Nous venons ainsi de séparer d'une part nos données en features et en target. Vu que la target est constituée des labels ici qui est représenté par le status des tissus "normal(0)/tumoral(1)" alors notre target est constituée de la colonne "status" qui est la colonne "tissue_status" encodée. Puis les features sont constituées des autres colonnes sauf la colonne des Id.

D'autres part nous avons divisé nos données en données d'entraînement(80%) et de test(20%) grâce à la fonction `train_test_split()`. Avec la définition

d'un random state ici 42, je pourrais avoir les mêmes données donc les mêmes résultats à chaque fois que j'exécute mon code.

2.3 Standardisation des données

La standardisation des données est une étape importante dans la préparation des données pour le Machine Learning, surtout lorsque les features ont des unités ou des échelles très différentes. Elle facilite l'apprentissage au modèle et permet à l'algorithme de converger plus rapidement.

```
[17]: from sklearn.preprocessing import StandardScaler
      scaler = StandardScaler()
      X_train_scaled = scaler.fit_transform(X_train)

[18]: X_test_scaled = scaler.transform(X_test)
```

Figure 7: standardisation

Un point important avec la standardisation est que ses métriques (Moyenne et Ecart-type) se calculent uniquement sur les données d'entraînement puis on l'applique aux données d'entraînement et de test. Ici à travers `fit_transform()`, nous calculons les métriques puis nous l'appliquons aux données d'entraînement et à travers `transform()`, à l'appliquons aux données de test. Tout ceci pour que notre modèle ne puisse pas avoir des informations sur nos données de test. Cela pourrait biaiser la performance du modèle.

3 Exploration des données

Nous analysons la distribution des classes et réalisons des visualisations des données. Cette analyse aide à comprendre la structure des données avant d'appliquer des modèles.

3.1 Distribution des classes

Nous avons utilisé `describe()` pour analyser la répartition des classes "normal" et "cancer" :

```
[7]: dataset["tissue_status"].describe()

[7]: count      804
     unique      2
     top      normal
     freq      402
     Name: tissue_status, dtype: object
```

Figure 8: Analyse des classes

A travers cette observation, nous remarquons notre target compte 804 observations puis elle contient vraiment 02 classes et aussi on remarque que le mode (la classe la plus fréquente) est la classe normal avec une fréquence de 402. Ce qui signifie que la classe tumoral a également une fréquence de 402 vu qu'il n'y a pas de données manquantes ni de données dupliquées. A priori, nous savons donc que nos classes sont équilibrées.

3.2 Distribution des features

```
> # Statistique descriptive des variables
   for column in X.columns:
     print(X[column].describe())

count      804.000000
mean        5.588943
std         2.889901
min         0.610504
25%         3.174250
50%         5.467024
75%         8.444756
max        10.263846
Name: ADHIC, dtype: float64
```

Figure 9: Analyse des features

La fonction `describe()` nous permet également d'analyser le contenu de nos features. Comme nos features sont numériques, alors la fonction nous donne des informations différentes de celles de la variable "tissue_status" qui est catégorielle. En effet en plus de `count` (nombre d'observations), elle nous renvoie le `mean` (la moyenne arithmétique des valeurs de la variable), le `std` (l'écart-type de la variable), le `min` puis le `max` des valeurs de la variable et enfin elle renvoie les valeurs du 1er, de 2è et du 3è quartile. On rappelle que 25% des données sont inférieures au 1er quartile, 50% au 2ème et 75% au 3ème.

3.3 Visualisation des données

Nous utilisons des histogrammes pour visualiser la distribution des gènes entre les deux classes.

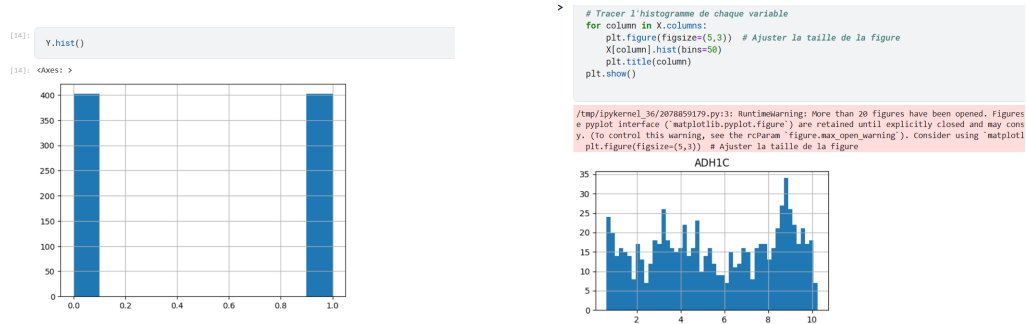


Figure 10: Visualisation

Avec les visualisations, nous pouvons confirmer les analyses statistiques que nous avons eu à faire plus haut. Ici pour les classes, par exemple, on remarque bien qu'elles sont équilibrées à travers l'histogramme.

4 Application des algorithmes de classification

Nous appliquerons deux méthodes de classification pour commencer : la Régression Logistique et le SVM.

4.1 Régression Logistique

4.1.1 1er modèle:

Nous entraînons un modèle de régression logistique sur l'ensemble de nos données standardisées:


```
[19]: lRegr = LogisticRegression()
      lRegr.fit(X_train_scaled, Y_train)

[19]: * LogisticRegression
      LogisticRegression()

>
      prediction = lRegr.predict(X_test_scaled)

      # Évaluation de la performance du modèle Logistic Regression
      accuracy2 = accuracy_score(Y_test, prediction)
      print("accuracy du modèle scaled :", accuracy2)

      # Calcul de la matrice de confusion pour le modèle Logistic Regression
      matrice_confusion = confusion_matrix(Y_test, prediction)
      #sns.heatmap(matrice_confusion, cmap="Blues", xticklabels = [])
      print("Matrice de confusion :")
      print(matrice_confusion)

accuracy du modèle scaled : 1.0
Matrice de confusion :
[[80  0]
 [ 0 81]]
```

Figure 11: Prédiction avec la regression logistique

Une fois l'entraînement effectué, vient la phase de prédiction puis la phase d'évaluation du modèle avec les données de test. Ainsi on compare les prédictions faite par le modèle avec les soties de test dont nous disposons. Puis les resultats sont regroupés dans la matrice de confusion Pour mesurer la performance de notre modèle, nous avons calculer son accuracy(qui calcule le rapport de bonne prédictions sur l'ensemble des données). Ici nous avons obtenu une accuracy de 1 on a plus besoin de calculer d'autres comme le recall ou encore le precesion et Fscore. Toutefois une accuracy de 1 nous interpelle beaucoup. En effet soit notre modèle a fait du overfitting, soit nos données étaient très clean pour permettre une telle performance.

4.1.2 Nouveau Modèle avec Top 1,2, 3 features:

L'intérêt de trouver les top features du modèle c'est de savoir quels sont les features qui interviennent plus dans la prise de décision dans mon modèle. Ainsi je pourrais reduire le nombre de features tout en ayant de bonnes performances.

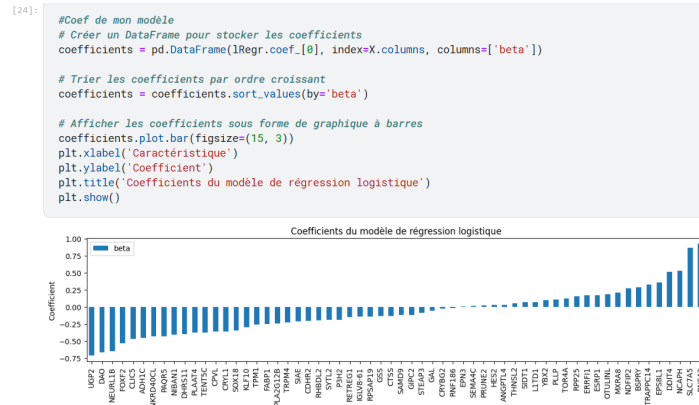


Figure 12: Recherche de top features

Nous remarquons que les tops 4 features sont SLC7A5, RNF43, UGP2 et DAO. Cependant, il est important de vérifier la corrélation entre ces features avant de les sélectionner. En effet si 02 features sont fortement corrélés, cela signifie qu'ils apportent quasiment la même information. Donc choisir les 02 aura le même impact que si l'on en choisissait seulement 1. Alors il serait plus intéressant de choisir parmi les tops features celle qui ne sont pas fortement corrélés.

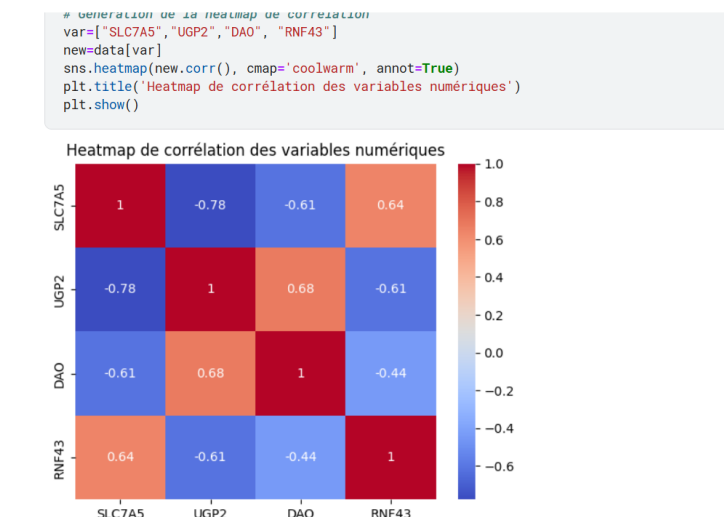


Figure 13: Corrélation entre top4 features

Nous remarquons que les features SLC7A5 et UGP2 sont fortement corrélés alors comme top3 features, nous choisirons SLC7A5, RNF43 et DAO.

```

feature_sets = [
    ["SLC7A5"],
    ["SLC7A5", "RNF43"],
    ["SLC7A5", "RNF43", "DAO"]
]

# Itérer sur chaque ensemble de caractéristiques
for idx, features in enumerate(feature_sets, 1):
    # Sélectionner les caractéristiques actuelles
    X = data[features]

    # Diviser les données en ensembles d'entraînement et de test
    X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=0.2, random_state=42)

    # Appliquer le scaling
    X_train_scaled = scaler.fit_transform(X_train)
    X_test_scaled = scaler.transform(X_test)

    # Initialiser et entraîner le modèle de régression logistique
    lr = LogisticRegression()
    lr.fit(X_train_scaled, Y_train)

    # Faire des prédictions sur l'ensemble de test
    predictions = lr.predict(X_test_scaled)

    # Évaluer la performance du modèle
    accuracy = accuracy_score(Y_test, predictions)
    conf_matrix = confusion_matrix(Y_test, predictions)

```

Figure 14: Test du modèle avec top1,2,3 feautres

Nous allons entrainer notre mmodèle maintenant sur le top1, le top2 puis le top3 features donc nous avons defini les 03 en tant que dataset. Puis à travers une boucle for, on obtient les 03 resultats.

```

# Afficher les résultats
print(f"\nRésultats pour Xtop{idx} ({', '.join(features)}):")
print(f"Accuraty du modèle : {accuracy:.4f}")
print("Matrice de confusion :")
print(conf_matrix)
print("-" * 40)

```

```

Résultats pour Xtop1 (SLC7A5):
Accuraty du modèle : 0.9255
Matrice de confusion :
[[74  6]
 [ 6 75]]
-----

Résultats pour Xtop2 (SLC7A5, RNF43):
Accuraty du modèle : 0.9565
Matrice de confusion :
[[79  1]
 [ 6 75]]
-----

Résultats pour Xtop3 (SLC7A5, RNF43, DAO):
Accuraty du modèle : 0.9938
Matrice de confusion :
[[79  1]
 [ 0 81]]
-----

```

Figure 15: Performance du modèle avec top1,2,3 features

Nous remarquons que plus nous ajoutons de features, plus la performance du modèle augmente. Avec les 03 tops features, nous avons une accuracy de 0.99 ce qui est assez excellent. Donc au lieu de travailler sur les 60features, avec seulement 03features nous pouvons avoir les mêmes résultats.

4.1.3 Etude de cas

```
[48]: new_patient = { 'SLC7A5': [4.18], 'RNF43': [4.68], 'DAO': [7.59]
new=pd.DataFrame(new_patient)
new.head()

[48]:   SLC7A5  RNF43  DAO
0      4.1    4.68  7.59

[50]: new_patient_scaled=scaler.transform(new)

[52]: predictions = lr.predict(new_patient_scaled)
print(predictions)

[0]

[54]: #La probabilité prédite pour ce patient d'avoir un cancer du colon
lr.predict_proba(new_patient_scaled)

[54]: array([[0.67515934, 0.32484066]])
```

Figure 16: Cas

Nous avons ensuite fait une prédiction pour un nouveau cas avec notre dernier modèle.

D'abord nous avons transformé les données du nouveau cas en dataframe avec pandas puis nous l'avons standardiser avant de procéder à la prédiction. Finalement notre modèle a prédit le nouveau cas comme normal avec une probabilité de 67%. Ce qui n'est pas trop fameux pour une prédiction de cancer. En effet cela signifie que le patient a 37% de chance d'avoir un cancer or nous savons qu'avec les cancers, le traitement est efficace lorsque le cancer est diagnostiqué tôt. Alors pour plus de prudence, nous recommanderons plus de test pour notre patient.

4.2 SVM (Support Vector Machine)

4.2.1 1er modèle

```

[46]: new_patient = { 'SLC7A5': [4.10], 'RNF43': [4.68], 'DAO': [7.59]}
      new = pd.DataFrame(new_patient)
      new.head()

[48]:
   SLC7A5  RNF43  DAO
0      4.1    4.68  7.59

[50]: new_patient_scaled = scaler.transform(new)

[52]:
      predictions = lr.predict(new_patient_scaled)
      print(predictions)

[0]

[54]:
      #La probabilité prédite pour ce patient d'avoir un cancer du colon
      lr.predict_proba(new_patient_scaled)

[54]: array([[0.67515934, 0.32484066]])

```

Figure 17: Prédiction avec le modèle SVM

Notre modèle SVM contient 02 paramètres: `random_state` qui nous permettra d'avoir la même distribution des données à chaque fois que nous exécutons notre code puis `probability=True` qui nous permettra de récupérer la probabilité de prédiction des différentes classes.

4.2.2 Test du modèle avec top1,2,3 feautres

```

# Afficher les résultats
print(f"\nRésultats pour Xtop{idx} ({', '.join(features)}):")
print(f"Accuracy du modèle : {accuracy:.4f}")
print("Matrice de confusion :")
print(conf_matrix)
print("-" * 40)

Résultats pour Xtop1 (SLC7A5):
Accuracy du modèle : 0.9130
Matrice de confusion :
[[76  4]
 [10 71]]
-----

Résultats pour Xtop2 (SLC7A5, RNF43):
Accuracy du modèle : 0.9441
Matrice de confusion :
[[78  2]
 [ 7 74]]
-----

Résultats pour Xtop3 (SLC7A5, RNF43, DAO):
Accuracy du modèle : 0.9938
Matrice de confusion :
[[79  1]
 [ 0 81]]

```

Figure 18: Prédiction avec le modèle SVM

Comparativement à la regression logistiques, nous voyons que avec seule-

ment les 03 tops features, nous avons les informations nécessaires pour réaliser nos prédictions.

4.2.3 Etude de cas

```
[33]: predictions = svm_model.predict(new_patient_scaled)
      print(predictions)

[0]

[34]: #La probabilité prédite pour ce patient d'avoir un cancer du colon
      svm_model.predict_proba(new_patient_scaled)

[34]: array([[0.73984611, 0.26015389]])
```

Figure 19: Prédiction d'un cas particulier

Nous remarquons que avec le modèle SVM, notre patient est normal à 73%. Ce qui est mieux par rapport aux 67% de la regression logistique.

5 Partie autonomie

5.1 K-NN

```
[34]: from sklearn.neighbors import KNeighborsClassifier
      # Initialiser le modèle K-NN
      knn_model = KNeighborsClassifier()
      # Entraînement du modèle sur les données d'entraînement
      knn_model.fit(X_train_scaled, Y_train)

      # Prédiction sur l'ensemble de test_scaled avec SVM
      predict = knn_model.predict(X_test_scaled)
      # Evaluation de la performance du modèle SVM
      Accuracy = accuracy_score(Y_test, predict)
      print('accuracy du modèle :', Accuracy)
      # Calcul de la matrice de confusion pour le modèle SVM
      matrice_confusion = confusion_matrix(Y_test, predict)
      print("Matrice de confusion :")
      print(matrice_confusion)

accuracy du modèle : 1.0
Matrice de confusion :
[[80  0]
 [ 0 81]]
```

Figure 20: Implémentation du K-NN

Nous avons implémenté le K- NN avec les hyperparamètres par défauts puis on remarque la même accuracy de 1. Ce qui est bon signe. Toute fois

nous allons chercher le nombre de voisin optimal pour notre modèle car par défaut il en prend 5.

```
[37]: from sklearn.model_selection import GridSearchCV
# Définir une grille d'hyperparamètres à tester
param_grid = {
    'n_neighbors': [1, 2, 3, 4, 10, 20], # nombre de voisins à tester
}
# Utiliser GridSearchCV pour trouver les meilleurs hyperparamètres
grid_search = GridSearchCV(knn_model, param_grid, cv=5, scoring='accuracy')
grid_search.fit(X_train_scaled, Y_train)

# Afficher les meilleurs hyperparamètres
print("Best hyperparameters for K-NN:", grid_search.best_params_)

Best hyperparameters for K-NN: {'n_neighbors': 3}
```

Figure 21: Nombre de voisins optimale

Ainsi nous observer grâce à une gridSearch que le nombre de voisin optimale est 3. C'est-à-dire qu'en obervant seulement 03 voisins le modèle est capable de prédire la classe d'une nouvelle observation. Ce qui plus performant.

5.2 Decision tree et Random Forrest

```
[54]: from sklearn.tree import DecisionTreeClassifier
dt_model = DecisionTreeClassifier()
# Entraînement du modèle sur les données d'entraînement
dt_model.fit(X_train_scaled, Y_train)

[54]: DecisionTreeClassifier
DecisionTreeClassifier()

# Prédiction sur l'ensemble de test_scaled avec Decision tree
predict = dt_model.predict(X_test_scaled)
# Evaluation de la performance du modèle Decision tree
Accuracy = accuracy_score(Y_test, predict)
print("accuracy du modèle :", Accuracy)
# Calcul de la matrice de confusion pour le modèle Decision tree
matrice_confusion = confusion_matrix(Y_test, predict)
print("Matrice de confusion :")
print(matrice_confusion)

accuracy du modèle : 1.0
Matrice de confusion :
[[ 0  0]
 [ 0  81]]
```

```
[58]: from sklearn.ensemble import RandomForestClassifier
# Initialiser le modèle Random Forest
rf_model = RandomForestClassifier()
# Entraînement du modèle sur les données d'entraînement
rf_model.fit(X_train_scaled, Y_train)

[58]: RandomForestClassifier
RandomForestClassifier()

# Prédiction sur l'ensemble de test_scaled avec Decision tree
predict = rf_model.predict(X_test_scaled)
# Evaluation de la performance du modèle Decision tree
Accuracy = accuracy_score(Y_test, predict)
print("accuracy du modèle :", Accuracy)
# Calcul de la matrice de confusion pour le modèle Decision tree
matrice_confusion = confusion_matrix(Y_test, predict)
print("Matrice de confusion :")
print(matrice_confusion)

accuracy du modèle : 1.0
Matrice de confusion :
[[ 0  0]
 [ 0  81]]
```

Figure 22: Implémentation des 02 modèles

Comme l'on pourrait s'en douter, avec ces 02 modèles aussi nous avons une accuracy de 1. Toutefois nous pouvons également faire une GridSearch pour rechercher les hyperparamètres optimaux pour ces 02 modèles. Cela sera intéressant car on pourrait notamment économiser des ressources. En attendant, nous allons chercher les top features qui interviennent plus dans la prise de décision de ces modèles.

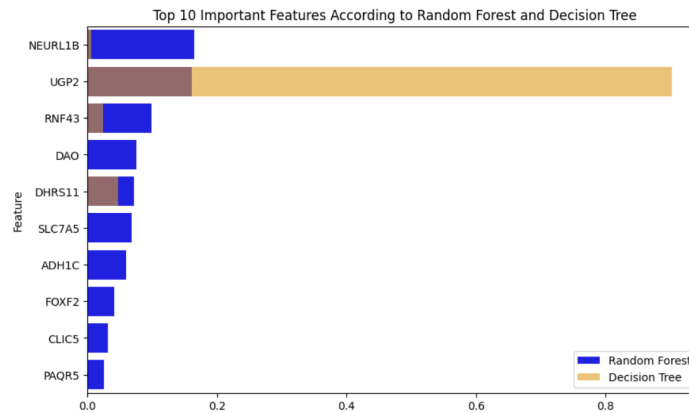


Figure 23: feature importance

Grâce à ce graphe, nous pouvons observer que pour la decision tree, il n'ya que 04 features (NEURL1B, UGP2, RNF43, DHRS11) qui interviennent dans la prise de décision. Donc on pourrait juste entrainer notre modèle sur ces 04 features puis avoir les mêmes resultats. Par contre pour le random forrest, plusieurs features interviennent dans la prise de décision cependant il ya des tops features avec lesquels on pourrait travailler puis avoir de bons résultats.

6 Conclusion

Ce rapport montre l'application des algorithmes de classification supervisée sur un dataset de cancer du côlon. Nous avons évalué la performance des modèles en utilisant plusieurs métriques et fait une prédiction pour un nouveau cas clinique. Particulièrement, notre dataset était bien clean et a surement été conçu sur mesure vu les performances de nos modèles. Toutefois nous restons conscient que dans la réalités il ne sera pas aisé d'avoir de telle performance dès les premières implémentations. En ce moment, nous verrons vraiment l'importance des gridSearch et des tops features entre autre pour améliorer et optimiser nos résultats.