UE: Introduction à l'apprentissage supervisé In [2]: #Importons les packages nécessaires . import numpy as np import pandas as pd import seaborn as sns import matplotlib.pyplot as plt from summarytools import dfSummary from sklearn.model\_selection import train\_test\_split from sklearn.tree import DecisionTreeClassifier from sklearn.inspection import DecisionBoundaryDisplay from sklearn.tree import plot\_tree from sklearn.model\_selection import train\_test\_split

Exercices sur l'apprentissage supervisé(Machine Learning)

MOISE EHIMIGAYE SENGHOR

from sklearn.metrics import accuracy\_score, recall\_score, precision\_score, f1\_score, roc\_auc\_score, confusion\_matrix

from sklearn.tree import DecisionTreeRegressor, plot\_tree from sklearn.metrics import mean\_absolute\_error as MAE, mean\_squared\_error as MSE from sklearn.ensemble import RandomForestRegressor from sklearn.model\_selection import train\_test\_split, cross\_val\_score from sklearn.preprocessing import LabelEncoder from sklearn.ensemble import BaggingClassifier from sklearn.ensemble import RandomForestClassifier from sklearn.model selection import GridSearchCV import warnings warnings.filterwarnings('ignore')#enlever les warnings **Projet 1: les pingouins** - Commencencons par importer les données : In [3]: data = pd.read\_csv('penguins\_size.csv', sep =',') Exploration du dataset et recapitulatif In [4]: data.info() <class 'pandas.core.frame.DataFrame'>

RangeIndex: 344 entries, 0 to 343 Data columns (total 7 columns): # Column Non-Null Count Dtype --- ----------

2 culmen\_length\_mm 342 non-null float64 3 culmen\_depth\_mm 342 non-null float64 flipper\_length\_mm 342 non-null float64 5 342 non-null body\_mass\_g float64 sex 334 non-null object dtypes: float64(4), object(3) memory usage: 18.9+ KB data.head()

Out[5]: culmen\_length\_mm culmen\_depth\_mm flipper\_length\_mm body\_mass\_g species 0 Adelie Torgersen 39.1 18.7 181.0 3750.0 3800.0 FEMALE Adelie Torgersen 39.5 17.4 186.0 2 Adelie Torgersen 40.3 18.0 195.0 3250.0 FEMALE

> **Data Frame Summary** data Dimensions: 344 x 7 Duplicates: 0

> > Freqs / (% of Valid)

152 (44.2%)

124 (36.0%)

68 (19.8%)

168 (48.8%)

124 (36.0%)

52 (15.1%)

164 distinct values

80 distinct values

55 distinct values

94 distinct values

168 (48.8%)

165 (48.0%)

10 (2.9%)

1 (0.3%)

sex

MALE

NaN

NaN

3450.0 FEMALE

Graph

Missing

(0.0%)

0

2

2

2

(0.6%)

(0.6%)

10

(2.9%)

(0.6%)

(0.6%)

(0.0%)

3 Adelie Torgersen NaN NaN NaN Adelie Torgersen 193.0 4 36.7 19.3

344 non-null object

344 non-null object

In [6]: data.shape

Out[6]:

(344, 7)

from sklearn.tree import export\_graphviz

from sklearn.metrics import roc\_curve

from sklearn.metrics import accuracy\_score

from sklearn.metrics import confusion\_matrix

from matplotlib.ticker import PercentFormatter

from sklearn.metrics import precision\_recall\_curve, auc

import graphviz

0 species 1 island

In [7]: data.groupby('species').size() Out[7]:

species Adelie 152 Chinstrap 68 Gentoo 124 dtype: int64

In [8]: # Utilisons le package summary tools pour faire un récapitulatif du jeu de données . dfSummary(data)

Out[8]:

Variable 1. Adelie species 2. Gentoo [object] 1. Biscoe island 2. Dream

Stats / Values No 3. Chinstrap 2 [object] culmen\_length\_mm 3 [float64] culmen\_depth\_mm

3. Torgersen Mean (sd): 43.9 (5.5) min < med < max: 32.1 < 44.5 < 59.6 IQR (CV): 9.3 (8.0) Mean (sd): 17.2 (2.0) min < med < max: 4 [float64] 13.1 < 17.3 < 21.5 IQR (CV): 3.1 (8.7) Mean (sd): 200.9 (14.1) flipper\_length\_mm min < med < max: [float64] 172.0 < 197.0 < 231.0 IQR (CV): 23.0 (14.3) Mean (sd): 4201.8 (802.0) body\_mass\_g min < med < max: [float64] 2700.0 < 4050.0 < 6300.0 IQR (CV): 1200.0 (5.2) 1. MALE 2. FEMALE

6 7 [object] 3. nan 4. . - Gestion des données manquantes :

In [9]: data.isnull().sum() Out[9]: species 0 island 0 2 culmen\_length\_mm culmen\_depth\_mm 2 2 2 10

flipper\_length\_mm body\_mass\_g sex dtype: int64 In [10]: data = data.dropna() data.shape Out[10]: (334, 7)In [11]: Out[11]: species 0 island 0 culmen\_length\_mm 0 culmen\_depth\_mm 0 flipper\_length\_mm 0 body\_mass\_g sex 0

data.isnull().sum() # verification de l'épuration des données manquantes . dtype: int64 Il nous est demandé d'essayer de prédire l'espèce des pingouins en fonction des caractéristiques de leurs culmens. Les variables explicatives seront le bill depth(profondeur du bec) et le bill lenght(longueur du bec). Nous allons commencer par extraire les données que nous utiliserons pour notre modéle. In [12]: pingouins = data.iloc[:, [0, 2, 3]] pingouins.head()#verification Out[12]: species culmen\_length\_mm culmen\_depth\_mm 0 Adelie

1 Adelie 2 Adelie 4 Adelie Adelie In [13]: pingouins.shape Out[13]: (334, 3)In [14]: In [15]:

39.1

39.5

40.3

36.7

39.3

18.7

17.4

18.0

19.3

20.6

Adelie Chinstrap Gentoo

> Adelie Chinstrap Gentoo

60

Divisons le jeu de données en deux entités une partie pour l'entrainement du modele et l'autre pour le tester. X= pingouins.iloc[:, [1,2]] # variables explicatives y = pingouins.iloc[:, [0]]# variable cible X\_train, X\_test, y\_train, y\_test = train\_test\_split(X, y, test\_size=0.3, random\_state=0) # les paremetres nous on été c 1 - Entrainez un arbre de profondeur 1 :

Creeons le modele et entrainons le: In [16]: modele = DecisionTreeClassifier( max\_depth=1)#modele modele.fit(X\_train, y\_train)#entrainement Out[16]: DecisionTreeClassifier DecisionTreeClassifier(max depth=1) 2- faire une visualisation basée sur les caracteristiques du culmen avec en differenciant les especes des pingouins. In [17]: sns.scatterplot(x = 'culmen\_length\_mm', y = 'culmen\_depth\_mm', hue = 'species', style= 'species', palette='bright', data plt.legend(bbox\_to\_anchor=(1.05, 1), loc='upper left') plt.tight\_layout() plt.xlabel("Longueur du culmen(mm)") plt.ylabel("Profondeur du culmen(mm)") plt.title("Especes des pingouins en fonction de la longueur et la profondeur du culmen.")

Especes des pingouins en fonction de la longueur et la profondeur du culmen. 20

plt.show() Profondeur du culmen(mm) 18 16

14 3 - faire le meme scatterplot en utilisant l'abre du modele : In [18]: DecisionBoundaryDisplay.from estimator( modele, X train, response method="predict") sns.scatterplot( data=pingouins, x = 'culmen\_length\_mm', y = 'culmen\_depth\_mm', hue='species', palette='bright', style=' plt.legend(bbox\_to\_anchor=(1.05, 1), loc="upper left") plt.xlabel("Longueur du culmen(mm)")

35

plt.ylabel("Profondeur du culmen(mm)")

22

20

18

16

14

differenciative que la variable longueur.

gini = 0.109

samples = 87

value = [4, 1, 82]

class = Gentoo

impacter la précision des prédictions.

In [20]:

In [21]:

deux.

In [22]:

Out[22]:

In [23]:

èce : ['Adelie']

Proba.plot.bar()#viz plt.ylabel("Probabilité")

e 35 mm et de 17 mm de profondeur")

du modèle favorisant les classes majoritaires.

In [19]:

Out[19]:

35

4 - Faire une visualisation de l'arbre de décision :

40

45

partition est faite sur l'axe des ordonnées. On peut l'interpreter en disant que la variable profondeur est plus significative et

plot\_tree(modele, feature\_names=X\_train.columns.tolist(), class\_names=modele.classes\_.tolist())

Text(0.25, 0.25, 'gini = 0.109\nsamples = 87\nvalue = [4, 1, 82]\nclass = Gentoo'), Text(0.75, 0.25, 'gini = 0.489\nsamples = 146\nvalue = [92, 49, 5]\nclass = Adelie')]

> culmen\_depth\_mm <= 16.45 gini = 0.645samples = 233value = [96, 50, 87]class = Adelie

des ajustements visant à mieux intégrer les classes minoritaires, comme Chinstrap.

Visualisation de la probabilité d'appartenance à une espece de l'oiseau cité plus haut :

Proba = pd.Series(proba[0], index=modele.classes\_)#serie pour la visualisation

0.6

0.5

0.2

0.1

0.0

16,45 mm, l'espèce associée à cette feuille est Adélie.

# Commencons par créer le modele et par l'entrainer : modele\_2 = DecisionTreeClassifier( max\_depth=2)#modele

plt.legend(bbox\_to\_anchor=(1.05, 1), loc="upper left")

22

20

18

16

14

35

culmen depth mm <= 16.45 gini = 0.645 samples = 233value = [96, 50, 87] class = Adelie

> culmen\_length\_mm <= 44.25 gini = 0.489

samples = 146

value = [92, 49, 5]

class = Adelie

gini = 0.259

samples = 55

value = [3, 47, 5]

Gentoo

id diagnosis radius\_mean texture\_mean perimeter\_mean area\_mean smoothness\_mean compactness\_mean concavit

1001.0

1326.0

1203.0

386.1

1297.0

122.80

132.90

130.00

77.58

135.10

0.11840

0.08474

0.10960

0.14250

0.10030

0.27760

0.07864

0.15990

0.28390

0.13280

gini = 0.043

samples = 91

value = [89, 2, 0]

Profondeur du culmen(mm)

print(f"{feature}: {importance}")

culmen\_length\_mm <= 39.

gini = 0.109

samples = 87value = [4, 1, 82]

class = Gentoo

gini = 0.024

value = [0, 1, 82]

1.0

0.8

0.6 Probabilité

0.4

0.2

0.0

]])

10.38

17.77

21.25

20.38

14.34

Non-Null Count Dtype

-----

int64

object

float64

float64

float64

float64

float64

float64 float64

float64

float64

float64

float64

float64 float64

float64

float64

float64 float64

float64

float64

float64

float64

float64

float64

float64

float64

float64

float64

float64

float64

float64 float64

> **Data Frame Summary** cancer Dimensions: 569 x 33 Duplicates: 0

> > Freqs / (% of Valid)

569 distinct values

456 distinct values

479 distinct values

522 distinct values

539 distinct values

474 distinct values

537 distinct values

537 distinct values

542 distinct values

432 distinct values

499 distinct values

540 distinct values

519 distinct values

533 distinct values

528 distinct values

547 distinct values

541 distinct values

533 distinct values

507 distinct values

498 distinct values

545 distinct values

457 distinct values

511 distinct values

514 distinct values

544 distinct values

411 distinct values

529 distinct values

357 (62.7%)

212 (37.3%)

Graph

Missing

0

0

0

0

0

0

0

0

(0.0%)

(0.0%)

0

0

(0.0%)

(0.0%)

0

(0.0%)

(0.0%)

0

(0.0%)

(0.0%)

0

(0.0%)

(0.0%)

0

(0.0%)

(0.0%)

0

(0.0%)

(0.0%)

0

(0.0%)

(0.0%)

(0.0%)

0

0

0

0

0

0

0

(0.0%)

569

concavity\_mean

0.3001

0.0869

0.1974

0.2414

0.1980

0.27760

0.07864

0.15990

0.28390

0.13280

(100.0%)

(0.0%)

(0.0%)

(0.0%)

(0.0%)

(0.0%)

(0.0%)

(0.0%)

(0.0%)

(0.0%)

(0.0%)

(0.0%)

(0.0%)

(0.0%)

(0.0%)

569 non-null

0 non-null

Stats / Values

min < med < max:

1. B

2. M

IQR (CV): 7943911.0 (0.2)

Mean (sd): 14.1 (3.5)

min < med < max:

7.0 < 13.4 < 28.1

IQR (CV): 4.1 (4.0)

Mean (sd): 19.3 (4.3)

min < med < max:

9.7 < 18.8 < 39.3

IQR (CV): 5.6 (4.5)

min < med < max:

43.8 < 86.2 < 188.5

IQR (CV): 28.9 (3.8)

Mean (sd): 654.9 (351.9) min < med < max:

143.5 < 551.1 < 2501.0

IQR (CV): 362.4 (1.9)

Mean (sd): 0.1 (0.0) min < med < max:

0.1 < 0.1 < 0.2

IQR (CV): 0.0 (6.9)

Mean (sd): 0.1 (0.1)

min < med < max:

IQR (CV): 0.1 (2.0)

Mean (sd): 0.1 (0.1)

min < med < max:

IQR (CV): 0.1 (1.1)

Mean (sd): 0.0 (0.0)

min < med < max:

IQR (CV): 0.1 (1.3)

Mean (sd): 0.2 (0.0)

min < med < max:

IQR (CV): 0.0 (6.6)

Mean (sd): 0.1 (0.0)

min < med < max:

IQR (CV): 0.0 (8.9)

Mean (sd): 0.4 (0.3)

min < med < max:

IQR (CV): 0.2 (1.5)

Mean (sd): 1.2 (0.6)

min < med < max:

IQR (CV): 0.6 (2.2)

Mean (sd): 2.9 (2.0)

min < med < max:

IQR (CV): 1.8 (1.4)

6.8 < 24.5 < 542.2

IQR (CV): 27.3 (0.9)

Mean (sd): 0.0 (0.0)

min < med < max:

IQR (CV): 0.0 (2.3)

Mean (sd): 0.0 (0.0)

min < med < max:

IQR (CV): 0.0 (1.4)

Mean (sd): 0.0 (0.0) min < med < max:

IQR (CV): 0.0 (1.1)

Mean (sd): 0.0 (0.0)

min < med < max:

IQR (CV): 0.0 (1.9)

Mean (sd): 0.0 (0.0) min < med < max:

IQR (CV): 0.0 (2.5)

Mean (sd): 0.0 (0.0)

min < med < max:

IQR (CV): 0.0 (1.4)

Mean (sd): 16.3 (4.8)

min < med < max:

7.9 < 15.0 < 36.0

IQR (CV): 5.8 (3.4)

Mean (sd): 25.7 (6.1)

min < med < max:

12.0 < 25.4 < 49.5

IQR (CV): 8.6 (4.2)

min < med < max:

50.4 < 97.7 < 251.2

IQR (CV): 41.3 (3.2)

Mean (sd): 880.6 (569.4) min < med < max:

185.2 < 686.5 < 4254.0

IQR (CV): 568.7 (1.5)

Mean (sd): 0.1 (0.0)

min < med < max:

IQR (CV): 0.0 (5.8)

Mean (sd): 0.3 (0.2) min < med < max:

IQR (CV): 0.2 (1.6)

Mean (sd): 0.3 (0.2)

min < med < max:

IQR (CV): 0.3 (1.3)

Mean (sd): 0.1 (0.1)

min < med < max:

IQR (CV): 0.1 (1.7)

Mean (sd): 0.3 (0.1)

min < med < max:

IQR (CV): 0.1 (4.7)

Mean (sd): 0.1 (0.0)

min < med < max:

IQR (CV): 0.0 (4.6)

0.1 < 0.1 < 0.2

1. nan

cancer\_sein= cancer.drop(['id','Unnamed: 32'], axis =1)

10.38

17.77

21.25

20.38

14.34

1. Divisez aléatoirement le jeu de données en partie train et partie test :

122.80

132.90

130.00

77.58

135.10

Choisissons les données prevues pour notre modele en selectionnant certaines variables du jeu de données :

17.99

20.57

19.69

11.42

20.29

0.2 < 0.3 < 0.7

0.0 < 0.1 < 0.3

0.0 < 0.2 < 1.3

0.1 < 0.1 < 0.2

0.0 < 0.2 < 1.1

Mean (sd): 107.3 (33.6)

0.0 < 0.0 < 0.0

0.0 < 0.0 < 0.1

0.0 < 0.0 < 0.1

0.0 < 0.0 < 0.1

0.0 < 0.0 < 0.4

0.0 < 0.0 < 0.0

Mean (sd): 40.3 (45.5) min < med < max:

0.8 < 2.3 < 22.0

0.4 < 1.1 < 4.9

0.1 < 0.3 < 2.9

0.0 < 0.1 < 0.1

0.1 < 0.2 < 0.3

0.0 < 0.0 < 0.2

0.0 < 0.1 < 0.4

0.0 < 0.1 < 0.3

Mean (sd): 92.0 (24.3)

Mean (sd): 30371831.4 (125020585.6)

8670.0 < 906024.0 < 911320502.0

569 non-null

gini = 0.0

samples = 4

value = [4, 0, 0]

modele\_2.fit(X\_train, y\_train)#entrainement

DecisionTreeClassifier

DecisionTreeClassifier(max depth=2)

Adelie

Probabilité 0.3

est 35 mm et dont la profondeur est de 17 mm.

Longueur du culmen(mm)

La variable utilisée pour partitionner le jeu de données est la variable culmen\_depth\_mm (profondeur du culmen) car on voit que la

[Text(0.5, 0.75, 'culmen\_depth\_mm <= 16.45\ngini = 0.645\nsamples = 233\nvalue = [96, 50, 87]\nclass = Adelie'),

gini = 0.489

samples = 146

value = [92, 49, 5]class = Adelie

Gini de 0.1, la classe Gentoo compte 82 individus sur un total de 87, témoignant d'une forte pureté de la feuille.

Les labels obtenus semblent cohérents, car on observe que, pour les feuilles avec un indice de Gini faible, la différence de population entre l'espèce majoritaire et les autres est nettement marquée. Par exemple, dans la feuille numéro deux, qui présente un indice de

En revanche, pour les feuilles où l'indice de Gini est plus élevé, les proportions entre les différentes espèces tendent à s'équilibrer, indiquant une diminution de la pureté. Cela reflète la difficulté du modèle à distinguer clairement les classes dans ces cas, ce qui peut

Par ailleurs, on note que la classe Chinstrap n'est pas représentée dans les feuilles obtenues, alors qu'elle constitue 20 % du jeu de données initial. Cela suggère que cette classe est sous-représentée dans le modèle ou qu'elle n'a pas été correctement prise en compte lors de l'apprentissage. Cette observation pourrait être due à un déséquilibre dans les données ou à un choix de paramètres

Ces résultats mettent néanmoins en évidence une certaine efficacité prédictive du modèle, particulièrement pour la classe dominante (Gentoo) dans certaines feuilles. Une analyse complémentaire serait utile pour mieux comprendre les limites du modèle et envisager

5 - Faisons une prédiction à partir du modèle pour un pingouin dont la longueur du bec

print("Le modèle prédit qu'un pingouin dont la longueur de culmen est de 35 mm et dont la profondeur est de 17 mm est d

Le modèle prédit qu'un pingouin dont la longueur de culmen est de 35 mm et dont la profondeur est de 17 mm est de l'esp

plt.title("Prediction de la probabilité d'appartenance à une espece pour un pingouin avec une longueur de culmen de 35

Text(0.5, 1.0, "Prediction de la probabilité d'appartenance à une espece pour un pingouin avec une longueur de culmen d

Prediction de la probabilité d'appartenance à une espece pour un pingouin avec une longueur de culmen de 35 mm et de 17 mm de profondeur

Chinstrap

La prévision correspond à une feuille de l'arbre de décision. En effet, le modèle prédit que si la profondeur du bec est supérieure à

6 - Refaisons la meme chose en faisant un autre arbre de décision avec une profondeur de

#Essayons aussi de faire le nécéssaire pour avoir une visualisation de la repartition des especes en fonction de modele

sns.scatterplot( data=pingouins, x = 'culmen\_length\_mm', y = 'culmen\_depth\_mm', hue='species', palette='bright', style='

45

Longueur du culmen(mm)

55

60

DecisionBoundaryDisplay.from\_estimator( modele\_2, X\_train, response\_method="predict", cmap="tab10", alpha=0.5)

Gentoo

Adelie Chinstrap Gentoo

proba = modele.predict\_proba(pd.DataFrame({"culmen\_length\_mm": [35], "culmen\_depth\_mm": [17]}))#probabilité

50

55

60

Profondeur du culmen(mm)

Longueur du culmen(mm)

plt.title("Especes des pingouins en fonction des caracteristisques du culmen et de l'arbre de décision")

Text(0.5, 1.0, "Especes des pingouins en fonction des caracteristisques du culmen et de l'arbre de décision")

Especes des pingouins en fonction des caracteristisques du culmen et de l'arbre de décision

plt.xlabel("Longueur du culmen(mm)") plt.ylabel("Profondeur du culmen(mm)") plt.title("Especes des pingouins en fonction des caracteristisques du culmen et de l'arbre de décision (de profondeur 2 Text(0.5, 1.0, "Especes des pingouins en fonction des caracteristisques du culmen et de l'arbre de décision (de profond eur 2)") Especes des pingouins en fonction des caracteristisques du culmen et de l'arbre de décision (de profondeur 2) In [24]: # Essayons de déterminer comment l'arbre de décision à fait les dichotomies : importances = modele\_2.feature\_importances\_ for feature, importance in zip(X\_train.columns, importances): culmen\_length\_mm: 0.46741020473018674 culmen\_depth\_mm: 0.5325897952698132 La visualisation met en évidence deux dichotomies successives. L'analyse des feature\_importances révèle que la première scission est effectuée sur la profondeur du culmen (culmen\_depth\_mm), suivie de la longueur du culmen (culmen\_length\_mm). In [25]: # Essayons de visualiser l'arbre de décision : plot\_tree(modele\_2, feature\_names=X\_train.columns.tolist(), class\_names=modele\_2.classes\_.tolist()) Out[25]: [Text(0.5, 0.833333333333334, 'culmen\_depth\_mm <= 16.45\ngini = 0.645\nsamples = 233\nvalue = [96, 50, 87]\nclass = Ad elie'), Text(0.25, 0.5, 'culmen\_length\_mm <= 39.7\ngini = 0.109\nsamples = 87\nvalue = [4, 1, 82]\nclass = Gentoo'), Text(0.375, 0.16666666666666666, 'gini = 0.024\nsamples = 83\nvalue = [0, 1, 82]\nclass = Gentoo'), Text(0.75, 0.5, 'culmen\_length\_mm <= 44.25\ngini = 0.489\nsamples = 146\nvalue = [92, 49, 5]\nclass = Adelie'), Text(0.625, 0.16666666666666666, 'gini = 0.043\nsamples = 91\nvalue = [89, 2, 0]\nclass = Adelie'), Text(0.875, 0.1666666666666666, 'gini = 0.259\nsamples = 55\nvalue = [3, 47, 5]\nclass = Chinstrap')] On constate une nette amélioration des performances du modèle. En effet, l'indice de Gini de trois des quatre feuilles est inférieur à 0,1, ce qui indique une amélioration significative de la pureté des feuilles par rapport au modèle avec une profondeur limitée à 1. De plus, on remarque que la classe Chinstrap, auparavant absente, est désormais représentée et prise en compte dans les prédictions, ce qui témoigne d'une meilleure capacité du modèle à intégrer toutes les classes présentes dans le jeu de données. In [26]: # Faisons une prédiction de l'especes d'un pingoin avec une profondeur de 17 et un longeur de 35 : proba = modele\_2.predict\_proba(pd.DataFrame({"culmen\_length\_mm": [35], "culmen\_depth\_mm": [17]}))#probabilité Proba = pd.Series(proba[0], index=modele\_2.classes\_)#serie pour la visualisation Proba.plot.bar()#viz plt.ylabel("Probabilité") plt.title("Prediction de la probabilité d'appartenance à une espece pour un pingouin avec une longueur de culmen de 35 Out[26]: Text(0.5, 1.0, "Prediction de la probabilité d'appartenance à une espece pour un pingouin avec une longueur de culmen d e 35 mm et de 17 mm de profondeur(modele 2)") Prediction de la probabilité d'appartenance à une espece pour un pingouin avec une longueur de culmen de 35 mm et de 17 mm de profondeur(modele 2) In [27]: modele\_2.predict\_proba(pd.DataFrame({"culmen\_length\_mm": [35], "culmen\_depth\_mm": [17]})) Out[27]: array([[0.97802198, 0.02197802, 0. Projet 2 : Cancer du sein Exercice 1: méthode du Hold Out Importons le dataset et faire une analyse exploratoire : cancer = pd.read\_csv('cancer.csv') cancer.head() Out[28]: 0 84300903 84348301 **4** 84358402 5 rows × 33 columns In [29]: cancer.shape Out[29]: (569, 33)In [30]: cancer.info() <class 'pandas.core.frame.DataFrame'> RangeIndex: 569 entries, 0 to 568 Data columns (total 33 columns): # 0 1 2 3 texture\_mean 4 perimeter\_mean 5 area\_mean 6 smoothness\_mean 7 8 9 10 11 fractal\_dimension\_mean 569 non-null 12 radius\_se 13 texture\_se 14 perimeter\_se 15 area\_se 16 smoothness\_se 17 compactness\_se 18 concavity\_se 19 concave points\_se 20 21 fractal\_dimension\_se radius\_worst 23 texture worst 24 perimeter\_worst area\_worst 25 26 27 28 29 30 fractal\_dimension\_worst 569 non-null 32 Unnamed: 32 dtypes: float64(31), int64(1), object(1) memory usage: 146.8+ KB In [31]: # Résumé statistique du jeu de données : dfSummary(cancer) Out[31]: 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30

842302

842517

Column -----

diagnosis

radius\_mean

compactness\_mean

concave points\_mean

concavity\_mean

symmetry\_mean

symmetry\_se

smoothness\_worst

compactness\_worst

concave points\_worst

concavity\_worst

symmetry\_worst

**Variable** 

id

[int64]

diagnosis

radius\_mean

texture mean

perimeter\_mean

 $smoothness\_mean$ 

compactness\_mean

concavity\_mean

[object]

[float64]

[float64]

[float64]

area\_mean

[float64]

[float64]

[float64]

[float64]

concave

[float64]

[float64]

[float64]

radius\_se

texture se

perimeter\_se

[float64]

[float64]

area\_se

[float64]

area\_worst

smoothness\_worst

compactness\_worst

concavity\_worst

[float64]

[float64]

[float64]

[float64]

concave

[float64]

[float64]

[float64]

[float64]

cancer\_sein.head()

Μ

M

Μ

Μ

5 rows × 31 columns

Unnamed: 32

31

32

33

In [32]:

In [33]:

Out[33]:

0

1

2

3

4

points\_worst

symmetry\_worst

fractal\_dimension\_w

radius\_worst

texture\_worst

perimeter\_worst

symmetry\_se

fractal\_dimension\_s

smoothness\_se

compactness\_se

concavity\_se

concave points\_se

[float64]

points\_mean

symmetry\_mean

fractal\_dimension\_

id

Μ

M

M

Μ

17.99

20.57

19.69

11.42

20.29

539 distinct values 492 distinct values 500 distinct values 535 distinct values 569 (100.0%) #Le dataset semble correct avec des données manquantes sur une seule variable et des données qui semblent vraisembles. # Nous avons choisi deux variables non pertinentes la variable 'id' et la 'Unnamed :32'(totalement vide) à élaguer pour diagnosis radius\_mean texture\_mean perimeter\_mean area\_mean smoothness\_mean compactness\_mean 1001.0 1326.0 1203.0 386.1 1297.0

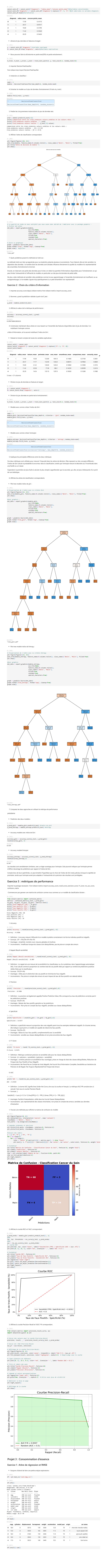
0.11840

0.08474

0.10960

0.14250

0.10030





0.02 - 0.00	othotase age sterase gillings	n albumin proteens	Ratio ailitubin Gender	
	plexité du modèle, hy		Ratio Gender Gen	
La grid search est une hyperparametres pour ombinatoire de la meill  2. Instancier la foret al	e methode particulière du param r optimiser notre modèle. Aussi d leure association de parametres léatoire	neter tuning qui cons appellé recherche ex parmi une grille de	iste à rechercher la meilleure combinaison haustive de parametres, elle fait une recherche parametres prédéfinis.	
3. Definir la grille d'hy			F=1, min_samples_split=2, random_state=1)  14, 0.16, 0.18]}	
[110]: modele wille = GridSearchCV Entrainement wille.fit(X_train, y		rid=parametres, s	coring='roc_auc')	
► RandomFore [111]: cint("les meilleurs	omForestClassifier estClassifier hyperparamètres : ", grille ramètres : {'max_depth': 2,		f': 0.18}	
6. Meilleur modele :  [112]: :illeur_modele = gri Représenter visuelleme	lle.best_estimator_ ent, par la méthode de votre cho			
rsults = pd.DataFram Transformation en to vot_table = results Visualisation Utiliser imshow pou ux = plt.imshow(pivo		th', columns='para ation ect='auto')	ons am_min_samples_leaf', values='mean_test_sco	re')
plt.text(j, i customisation de L t.title("Heatmap de t.xlabel("min_sampl t.ylabel("max_depth Legendes de L'axe d t.xticks(np.arange(	<pre>vot_table.shape[1]):    , f"{pivot_table.iloc[i, j]    a visualisation es Résultats de la Grid Sear</pre>	<pre>ch") es t_table.columns)</pre>	r', va='center', color='black', fontsize=12	)
t.show()	pivot_table.shape[0]), pivo  des Résultats de la Gri  0.720 0.719		- 0.723 - 0.722	
3 - 0.716	0.719 0.719	0.723	- 0.721 - 0.720 - 0.719 ONC - 0.718	
0.716	0.719 0.719 0.14 0.16 min_samples_leaf	0.723	- 0.717 - 0.716 - 0.715	
[]:				

Out[103]:

In [104]:

In [105]:

In [106]:

Out[106]:

In [107]:

plt.show()

dtype: float64

#visualisation

plt.xlabel("Variables", )

Age

Accuracy: 0.661 Recall: 0.912 AUC: 0.731

RandomForestClassifier

predictions1 = foret.predict(X\_test)# predictions du modele sur les données test

5. Estimer l'importance relative des différents features grace à rf. feature importances , en créant un vecteur

 $\verb|importance_vecteur.sort_values(ascending=False)| \textit{#tri décroissant}|$ 

sns.barplot(x=importance\_vecteur.index, y=importance\_vecteur.values, palette='inferno')
plt.title("Importance des variables dans le modele forêt aléatoire")

Importance des variables dans le modele forêt aléatoire

print('Accuracy: ',round(accuracy\_score(y\_test, predictions1),3))
print('Recall: ', round(recall\_score(y\_test, predictions1),3))

importance\_vecteur= pd.Series(importance, index=feature\_names)

0.136

print('AUC: ', round(roc\_auc\_score(y\_test, y\_prob1),3))

y\_prob1 = bagging\_model.predict\_proba(X\_test)[:, 1] #base de calcul pour AUC pour des probabilité positives

RandomForestClassifier(random\_state=42)

importance = foret.feature\_importances\_

feature\_names = X.columns

round(importance\_vecteur[:3], 3)

Aspartate\_Aminotransferase 0.145 Alkaline\_Phosphotase 0.144

plt.ylabel("Importance de la variables")
plt.xticks(rotation=45, ha='right')