

Supervised and Unsupervised Learning in Machine Learning

Machine Learning is the science of making computers learn and act like humans by feeding data and information without being explicitly programmed.

Machine learning algorithms are trained with training data. When new data comes in, they can make predictions and decisions accurately based on past data.

For example, whenever you ask Siri to do something, a powerful speech recognition converts the audio into its corresponding textual form. This is sent to the Apple servers for further processing where language processing algorithms are run to understand the user's intent. Then finally, Siri tells you the answer.

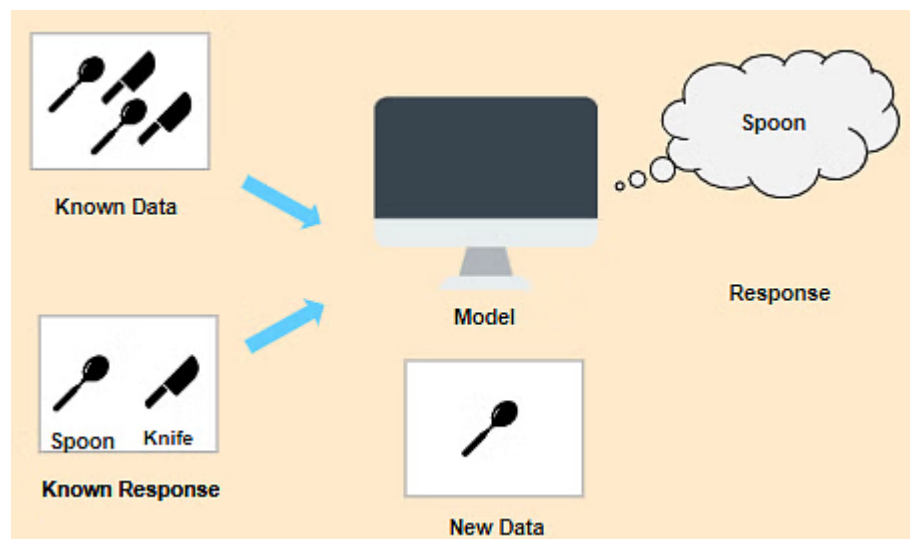
There are **two types** of machine learning:

1. Supervised Learning
2. Unsupervised Learning

Entrée [1]: 1 `from IPython.display import Image`

What is Supervised Learning?

In Supervised Learning, the machine learns under supervision. It contains a model that is able to predict with the help of a labeled dataset. A labeled dataset is one where you already know the target answer.



In this case, we have images that are labeled a spoon or a knife. This known data is fed to the machine, which analyzes and learns the association of these images based on its features such as shape, size, sharpness, etc. Now when a new image is fed to the machine without any label, the machine is able to predict accurately that it is a spoon with the help of the past data.

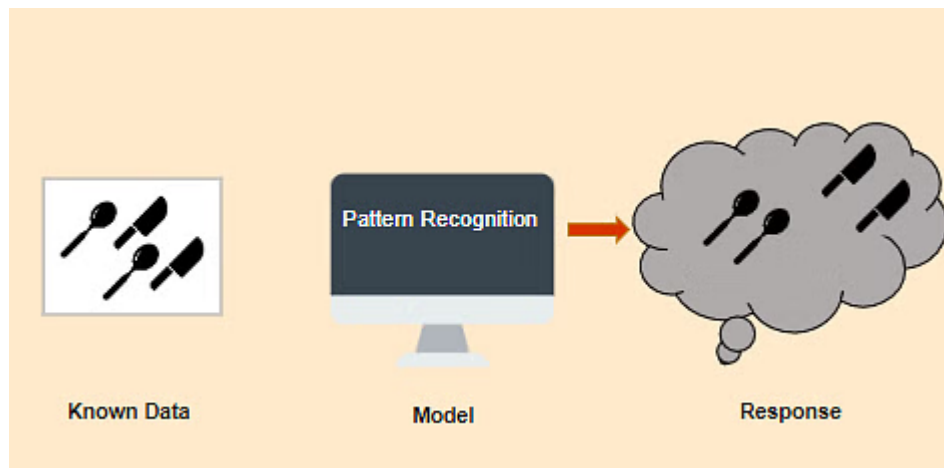
Supervised learning can be further divided into two types:

1. Classification

2. Regression

What is Unsupervised Learning?

In Unsupervised Learning, the machine uses unlabeled data and learns on itself without any supervision. The machine tries to find a pattern in the unlabeled data and gives a response.



Unsupervised learning can be further grouped into types:

1. Clustering
2. Association

Difference Between Supervised and Unsupervised Learning

Supervised Learning

It uses known and labeled data as input

It has a feedback mechanism

The most commonly used supervised learning algorithms are:

- Decision tree
- Logistic regression
- Support vector machine

Unsupervised Learning

It uses unlabeled data as input

It has no feedback mechanism

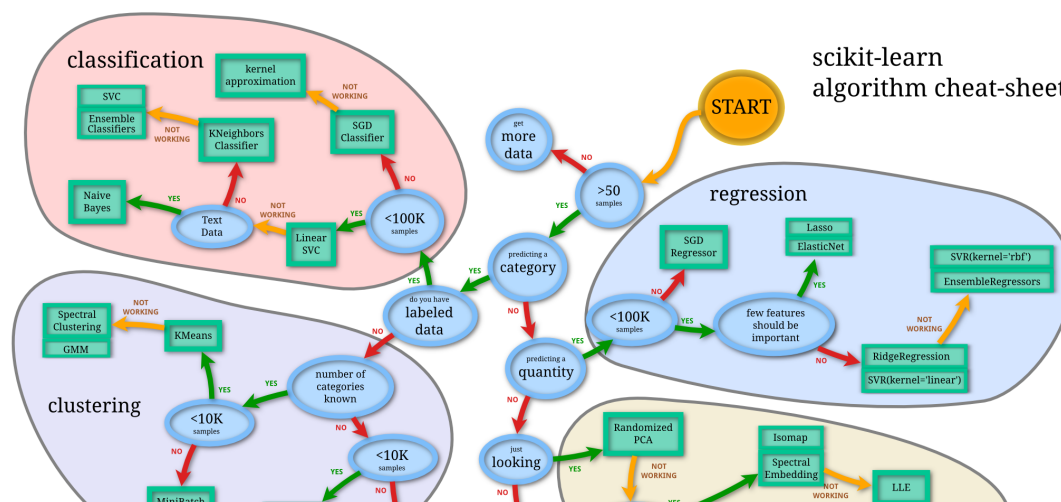
The most commonly used unsupervised learning algorithms are:

- K-means clustering
- Hierarchical clustering
- Apriori algorithm

Choosing the right estimator

Often the hardest part of solving a machine learning problem can be finding the right estimator for the job.

Different estimators are better suited for different types of data and different problems.



In machine learning, Classification is used to split data into categories. But after cleaning and preprocessing the data and training our model, how do we know if our classification model performs well? That is where a confusion matrix comes into the picture.

A confusion matrix is used to measure the performance of a classifier in depth. In this simple guide to Confusion Matrix, we will get to understand and learn confusion matrices better.

MNIST dataset and performance measures

In this notebook, we would like to go through some of the ML algorithms. Precisely, we would like to evaluate and compare the following performance metrics :

- Confusion matrix
- Recall
- Precision
- FP Rate
- Specificity
- ROC curve

The measures will be taken on classification tasks on handwritten data.

The metrics will be implemented from scratch and will be compared to the metrics offered by the standard librairies (*scikit-learn*).

 KNN will be coded from scratch !

Importing librairies

```
Entrée [55]: 1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 from sklearn.model_selection import train_test_split
6 from sklearn import metrics
7 from sklearn.model_selection import cross_val_score
8 from sklearn.neighbors import KNeighborsClassifier
9 from sklearn.neural_network import MLPClassifier
10 from sklearn import tree
11 from sklearn.tree import DecisionTreeClassifier
12 from sklearn.tree import plot_tree
13 from sklearn.naive_bayes import GaussianNB, CategoricalNB
14 from sklearn import svm
15 from sklearn.tree import DecisionTreeClassifier
16 plt.style.use('seaborn')
17 from sklearn.svm import SVC
18 from sklearn import svm
19 import warnings
20 warnings.filterwarnings('ignore')
```

Lecture des fichiers de données

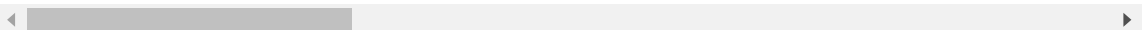
Pour ce TP, nous allons lire les données à partir d'un fichier csv.

```
Entrée [7]: 1 Data = pd.read_csv('Data/breast.csv')
2 Data
```

Out[7]:

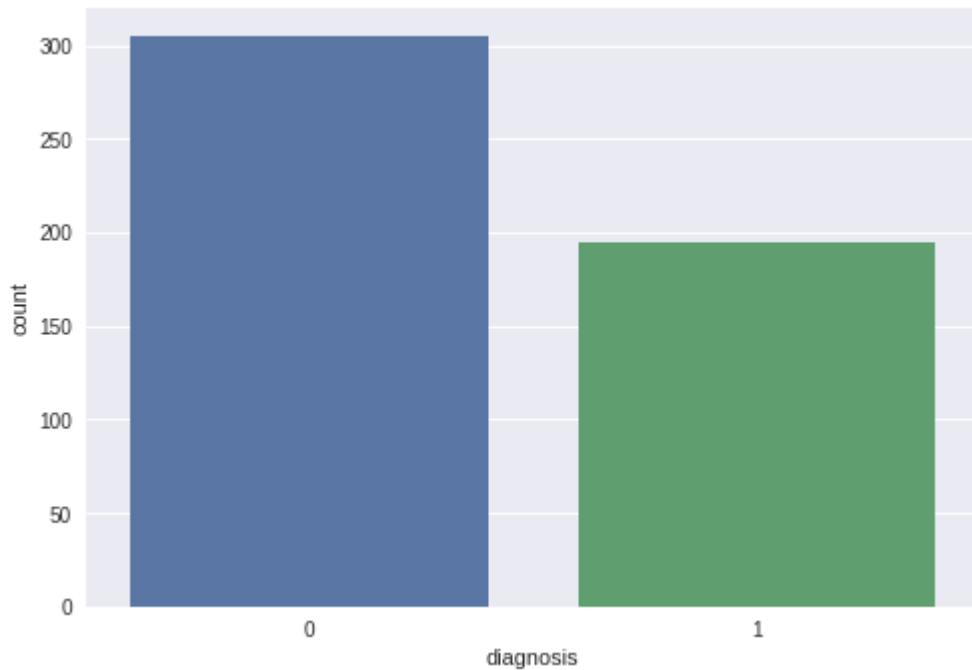
	radius_mean	texture_mean	perimeter_mean	area_mean	smoothness_mean	compactness_mean
0	17.99	10.38	122.80	1001.0	0.11840	0.26630
1	20.57	17.77	132.90	1326.0	0.08474	0.26688
2	19.69	21.25	130.00	1203.0	0.10960	0.26688
3	11.42	20.38	77.58	386.1	0.14250	0.26688
4	20.29	14.34	135.10	1297.0	0.10030	0.26688
...
495	14.87	20.21	96.12	680.9	0.09587	0.26688
496	12.65	18.17	82.69	485.6	0.10760	0.26688
497	12.47	17.31	80.45	480.1	0.08928	0.26688
498	18.49	17.52	121.30	1068.0	0.10120	0.26688
499	20.59	21.24	137.80	1320.0	0.10850	0.26688

500 rows × 19 columns



Entrée [10]: 1 sns.countplot(x='diagnosis', data=Data)

Out[10]: <AxesSubplot:xlabel='diagnosis', ylabel='count'>



Entrée [11]: 1 X = Data.drop("diagnosis", axis = 1)
2 Y = Data['diagnosis']

Normalisation

Entrée [12]:

```

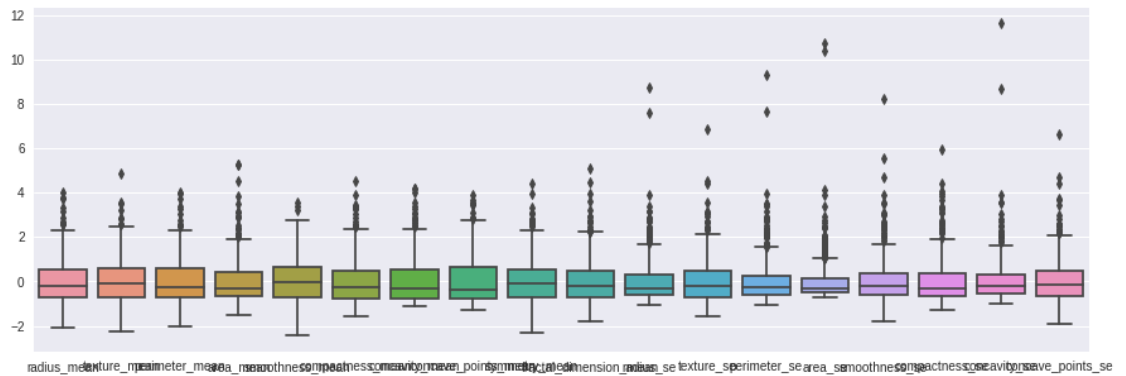
1 Newdata2 = X.copy()
2 for column in X.columns:
3     if Newdata2[column].dtype == 'float64' or Newdata2[column].dtype == 'int64':
4         Newdata2[column] = (X[column] - X[column].mean()) / X[column].std()
5     M=[]
6 l=Newdata2.shape[0]
7 for i in range(l):
8     c=0
9     for column in Newdata2.columns:
10        if Newdata2[column].dtype == 'float64':
11            c=c+1
12        M.append(Newdata2[column].iloc[i])
13 newar=np.asarray(M)
14 newdat=newar.reshape(l,c)
15 X = pd.DataFrame(newdat,columns = [column for column in Newdata2.columns])
16 
```

Box-Plot

Entrée [13]:

```
1 plt.figure(figsize=(15,5))
2 sns.boxplot(data=X)
```

Out[13]: <AxesSubplot:>



The training set (data & labels)



Data

Entrée [14]:

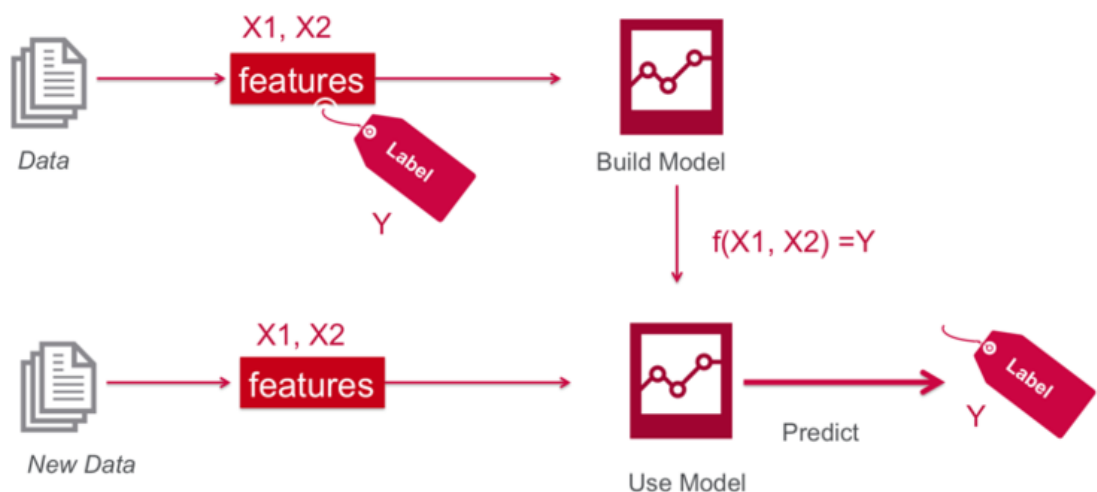
```
1 x_train, x_test, y_train ,y_test = train_test_split(X,Y,test_si
```



Labels

In machine learning, data labeling is the process of identifying raw data (images, text files, videos, etc.) and adding one or more meaningful and informative labels to provide context so that a machine learning model can learn from it.

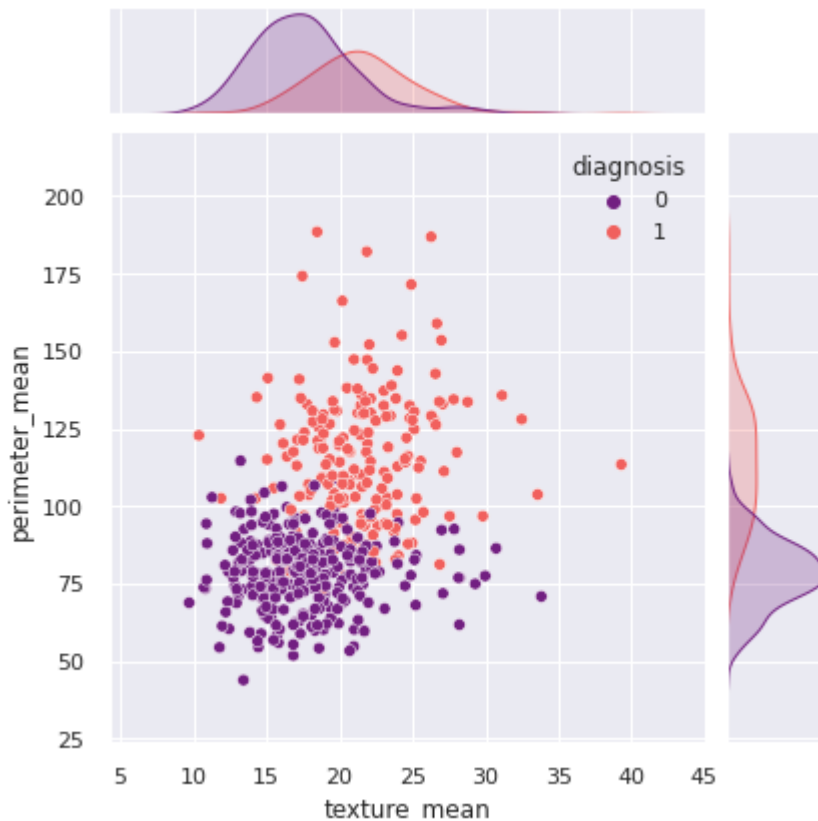
In our case, we have a csv file comprising of labels (from 1 to 10).



Let's dive into our `labels.csv` file and see how our labels look like :

```
Entrée [17]: 1 plt.figure(figsize=(20,10))
2 sns.set_theme()
3 sns.jointplot(data=Data,x='texture_mean',y='perimeter_mean',hue
```

```
Out[17]: <seaborn.axisgrid.JointGrid at 0x7f5fc1796c70>
```



By convention, **label 10** corresponds to **0**. This is used in order to facilitate heavy computations.



Performance measures

In this part, we're going to define the general performance measures all along with **their implementation from scratch**.

What Are Confusion Matrices, and Why Do We Need Them?

Classification Models have multiple categorical outputs. Most error measures will calculate the total error in our model, but we cannot find individual instances of errors in our model. The model might misclassify some categories more than others, but we cannot see this using a standard accuracy measure.

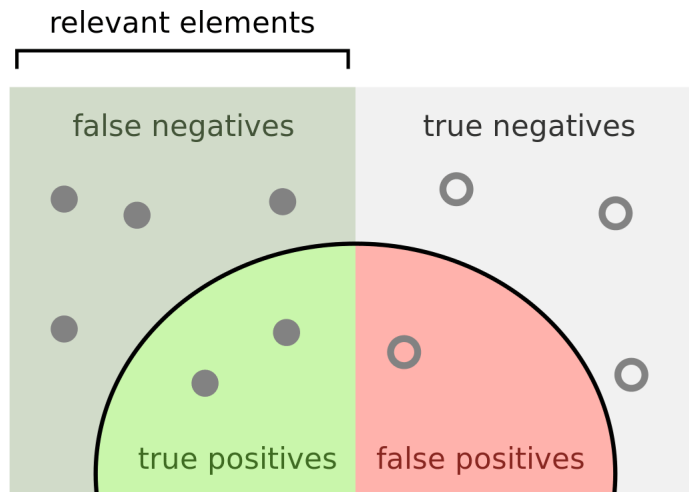
Furthermore, suppose there is a significant class imbalance in the given data. In that case, i.e., a class has more instances of data than the other classes, a model might predict the majority class for all cases and have a high accuracy score; when it is not predicting the minority classes. This is where confusion matrices are useful.

A confusion matrix presents a table layout of the different outcomes of the prediction and results of a classification problem and helps visualize its outcomes.

It plots a table of all the predicted and actual values of a classifier.

		Predicted	
		Negative (N) -	Positive (P) +
Actual	Negative -	True Negatives (TN)	False Positives (FP) Type I error
	Positive +	False Negatives (FN) Type II error	True Positives (TP)

False Positives (FP-Type 1 error) vs False Negatives (FN-Type 2 error)



● Confusion Matrix

In our case, we have 8 classes, so the matrix is 8*8.

The matrix will look like this one :

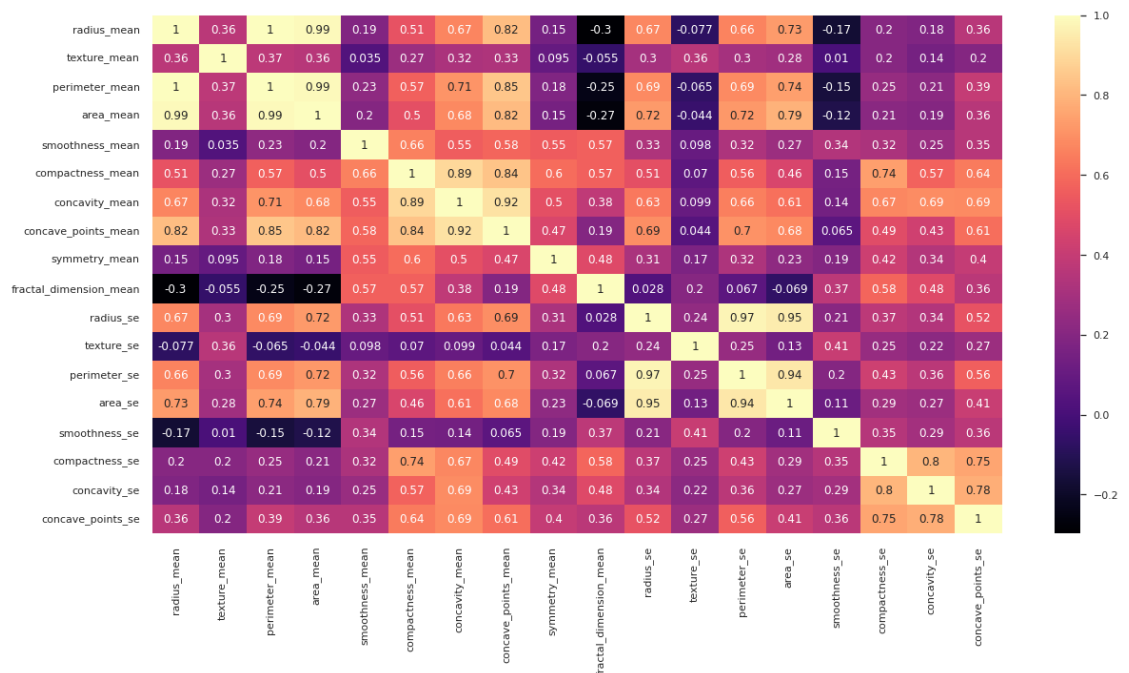
Entrée [18]:

```

1 plt.figure(figsize=(20,10))
2 corr = X.corr() #utilisé pour trouver la corrélation par paires
3 sns.heatmap(corr, annot=True, cmap="magma")

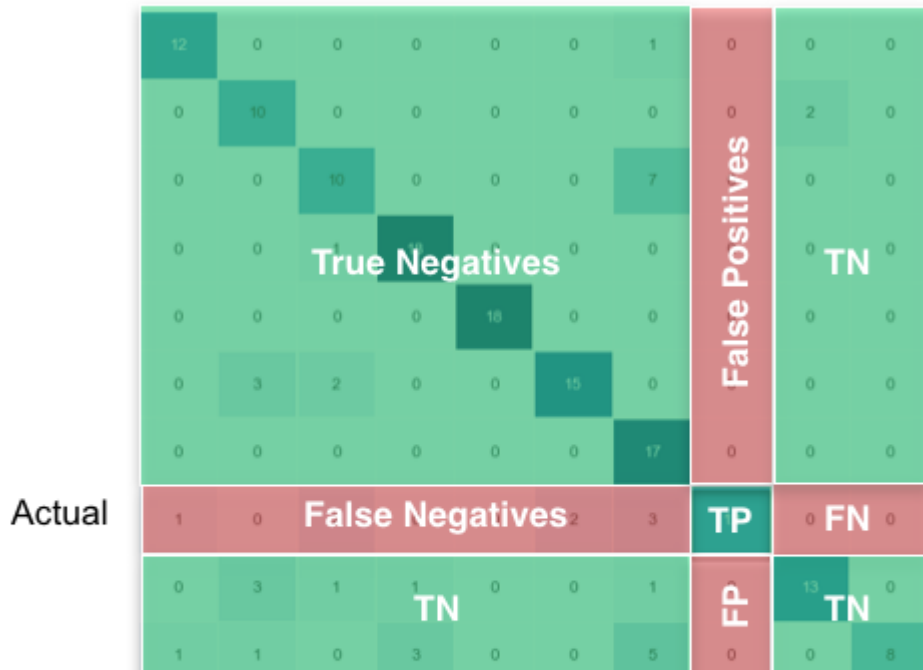
```

Out[18]: <AxesSubplot:>



Now, we got our confusion matrix.

In order to compute the main performance measure, we need to extract from that matrix the FP, FN, TP and TN.



Accuracy

How many values did we predict correctly? How many true predictions out of all samples there are?

$$\text{Accuracy} = \frac{\text{True Positives} + \text{True Negatives}}{\text{All Samples}}$$

Classification

Méthode 1 : KNN

Ici il faudra implémenter la méthode, puis la tester et vérifier les métriques en variant le nombre K

Entrée [19]:

```
1 KNN=KNeighborsClassifier(n_neighbors=3)
2 KNN.fit(x_train,y_train)
```

Out[19]:

```
KNeighborsClassifier
KNeighborsClassifier(n_neighbors=3)
```

Entrée [20]:

```

1 from prettytable import PrettyTable
2
3 x = PrettyTable(["Model", "Train SCORE", "Test SCORE"])
4 z=str(int(KNN.score(x_train,y_train)*100))+ "%"
5 v=str(int(KNN.score(x_test,y_test)*100))+ "%"
6 x.add_row(["KNN",z,v])
7 print(x)

```

```

+-----+-----+-----+
| Model | Train SCORE | Test SCORE |
+-----+-----+-----+
| KNN   | 96%         | 90%         |
+-----+-----+-----+

```

Confusion Matrices KNN

Entrée [21]:

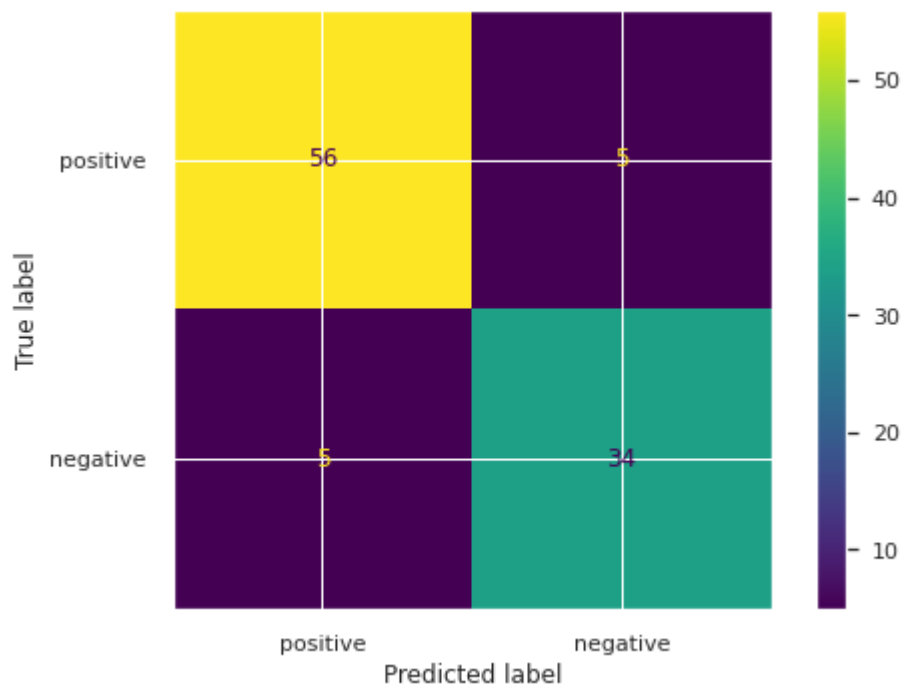
```

1 y_pred=KNN.predict(x_test)
2 print(metrics.classification_report(y_test,y_pred))
3 display=metrics.ConfusionMatrixDisplay((metrics.confusion_matri
4 display.plot()

```

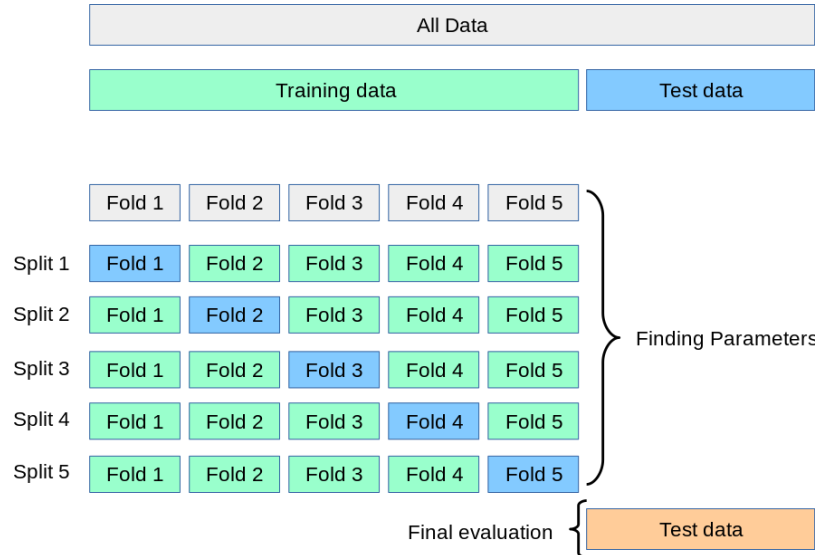
	precision	recall	f1-score	support
0	0.92	0.92	0.92	61
1	0.87	0.87	0.87	39
accuracy			0.90	100
macro avg	0.89	0.89	0.89	100
weighted avg	0.90	0.90	0.90	100

Out[21]: <sklearn.metrics._plot.confusion_matrix.ConfusionMatrixDisplay at 0x7f5fc1737f70>



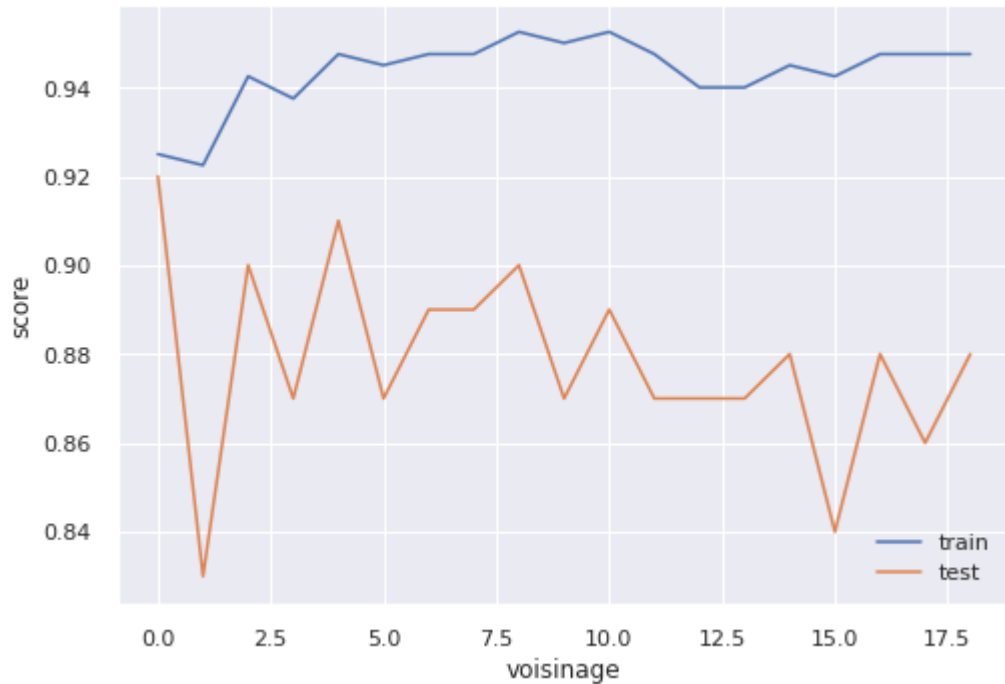
Cross-validation: evaluating estimator performance

La mesure de performance rapportée par la validation croisée k-fold est alors la moyenne des valeurs calculées dans la boucle. Cette approche peut être coûteuse en termes de calcul, mais ne gaspille pas trop de données (comme c'est le cas lors de la correction d'un ensemble de validation arbitraire), ce qui constitue un avantage majeur dans des problèmes tels que l'inférence inverse où le nombre d'échantillons est très petit.



```
Entrée [22]: 1 val_scortrain=[]
2 val_scorstest=[]
3 for k in range(1,20):
4     score=cross_val_score(KNeighborsClassifier(k),x_train,y_tra
5     val_scortrain.append(score)
6     score=cross_val_score(KNeighborsClassifier(k),x_test,y_test
7     val_scorstest.append(score)
8 plt.plot(val_scortrain,label='train')
9 plt.plot(val_scorstest,label='test')
10 plt.ylabel('score')
11 plt.xlabel('voisinage')
12 plt.legend()
```

Out[22]: <matplotlib.legend.Legend at 0x7f5fc037ae50>



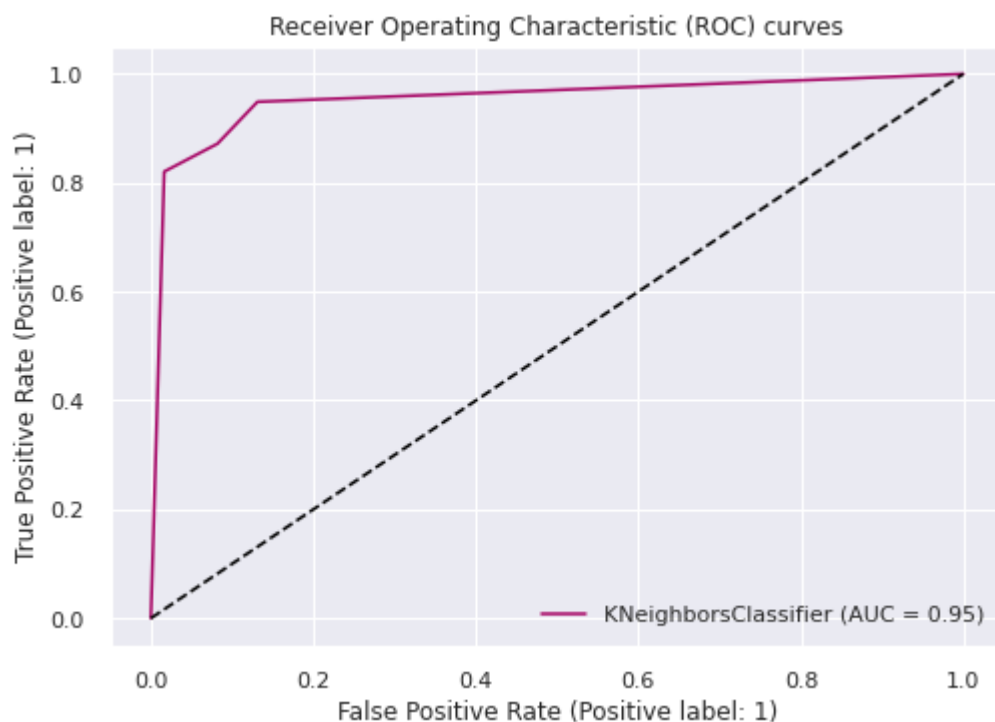
la courbe (ROC AUC)

Il s'agit d'une fonction générale, étant donné des points sur une courbe. Pour calculer l'aire sous la courbe ROC

```

Entrée [23]: 1 # Make predictions on the test set
2 y_pred = KNN.predict(x_test)
3 from sklearn.preprocessing import label_binarize
4 from sklearn.metrics import roc_curve, auc
5 from sklearn.metrics import DetCurveDisplay, RocCurveDisplay
6 y_test_binarized = label_binarize(y_test, classes=np.unique(y_t
7 if len(y_test_binarized[0])<=1:
8     fig, ax_roc = plt.subplots()
9     RocCurveDisplay.from_estimator(KNN, x_test, y_test, ax=ax_r
10     plt.plot([0,1],[0,1], '--', c='black')
11     ax_roc.set_title("Receiver Operating Characteristic (ROC) c
12     #ax_roc.grid(linestyle="--", c='black')
13     plt.legend()
14     plt.show()
15 elif (len(y_test_binarized[0])>1):
16     C=["#4B2991", "#952EA0", "#D44292", "#F66D7A", "#F6A97A"]
17     Classes=np.unique(y_test)
18     pred_prob=SVM.predict_proba(X_test)
19     fpr = {}
20     tpr = {}
21     thresh = {}
22     roc_auc = dict()
23     n_classe = len(Classes)
24     for i in range (n_classe):
25         fpr[i], tpr[i], thresh[i] = roc_curve(y_test_binarized[
26         roc_auc[i] = auc (fpr[i], tpr[i])
27     # traçage
28     plt.plot(fpr[i], tpr[i], label='%s vs Rest (AUC=%0.2f)'%
29     plt.plot([0,1], [0,1], linestyle='--', c='black')
30     plt.xlim([0,1])
31     plt.ylim([0,1.05])
32     plt.title('Receiver Operating Characteristic(ROC) curves')
33     plt.xlabel('False Rate')
34     plt.ylabel('True Rate')
35     plt.legend()
36     plt.show()

```



Méthode 2:Naive Bayes

Entrée [24]:

```
1 Naive=GaussianNB()
2 Naive.fit(x_train,y_train)
```

Out[24]:

```
▼ GaussianNB
GaussianNB()
```

Entrée [25]:

```
1 from prettytable import PrettyTable
2
3 x = PrettyTable(["Model", "Train SCORE", "Test SCORE"])
4 z=str(int(Naive.score(x_train,y_train)*100))+ "%"
5 v=str(int(Naive.score(x_test,y_test)*100))+ "%"
6 x.add_row(["Naive Bayes",z,v])
7 print(x)
```

```
+-----+-----+-----+
|   Model   | Train SCORE | Test SCORE |
+-----+-----+-----+
| Naive Bayes |      91%    |      88%    |
+-----+-----+-----+
```

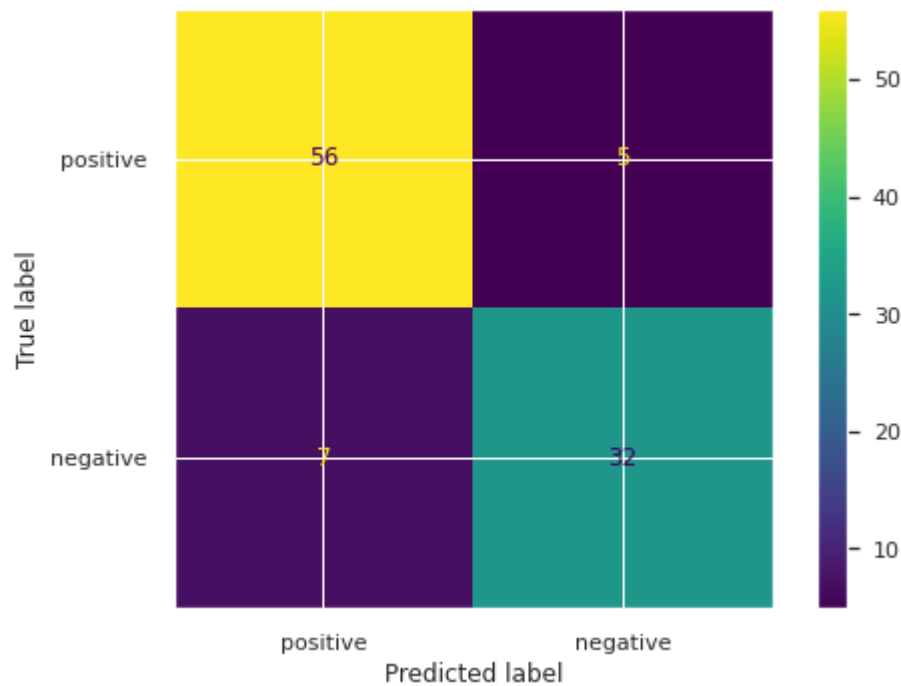
Confusion Matrices Naive Bayes

Entrée [26]:

```
1 y_pred=Naive.predict(x_test)
2 print(metrics.classification_report(y_test,y_pred))
3 display=metrics.ConfusionMatrixDisplay((metrics.confusion_matri
4 display.plot()
```

	precision	recall	f1-score	support
0	0.89	0.92	0.90	61
1	0.86	0.82	0.84	39
accuracy				100
macro avg				100
weighted avg				100

Out[26]: <sklearn.metrics._plot.confusion_matrix.ConfusionMatrixDisplay at 0x7f5fc0385af0>



Cross-validation: evaluating estimator performance

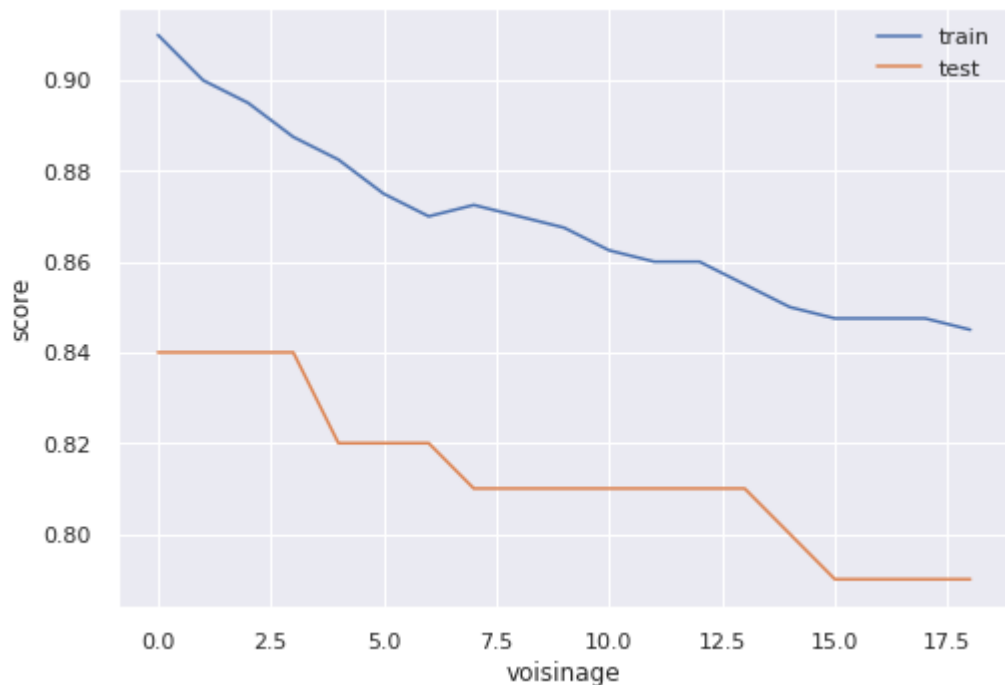
Entrée [27]:

```

1 val_scortrain=[]
2 val_scorctest=[]
3 for k in range(1,20):
4     m=float(k/4)
5     score=cross_val_score(GaussianNB(var_smoothing=m),x_train,y_
6     val_scortrain.append(score)
7     score=cross_val_score(GaussianNB(var_smoothing=m),x_test,y_
8     val_scorctest.append(score)
9 plt.plot(val_scortrain,label='train')
10 plt.plot(val_scorctest,label='test')
11 plt.ylabel('score')
12 plt.xlabel('voisinage')
13 plt.legend()

```

Out[27]: <matplotlib.legend.Legend at 0x7f5fc01f5b50>



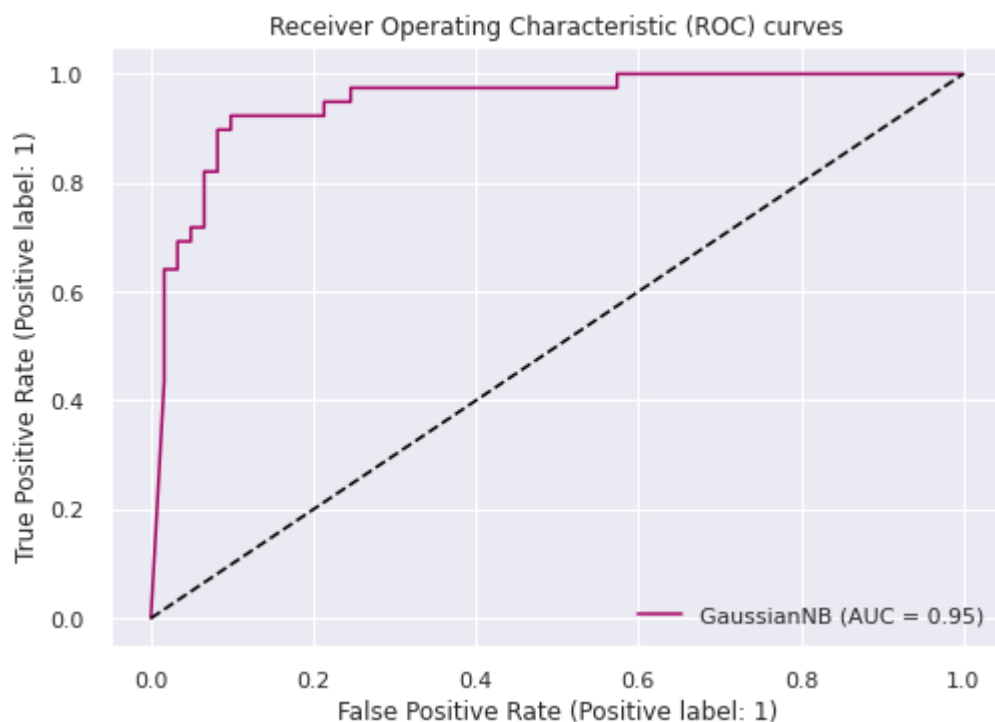
la courbe (ROC AUC)

Il s'agit d'une fonction générale, étant donné des points sur une courbe. Pour calculer l'aire sous la courbe ROC

```

Entrée [28]: 1 # Make predictions on the test set
2 y_pred = Naive.predict(x_test)
3 from sklearn.preprocessing import label_binarize
4 from sklearn.metrics import roc_curve, auc
5 from sklearn.metrics import DetCurveDisplay, RocCurveDisplay
6 y_test_binarized = label_binarize(y_test, classes=np.unique(y_t
7 if len(y_test_binarized[0])<=1:
8     fig, ax_roc = plt.subplots()
9     RocCurveDisplay.from_estimator(Naive, x_test, y_test, ax=ax
10     plt.plot([0,1],[0,1], '--', c='black')
11     ax_roc.set_title("Receiver Operating Characteristic (ROC) c
12     #ax_roc.grid(linestyle="--", c='black')
13     plt.legend()
14     plt.show()
15 elif (len(y_test_binarized[0])>1):
16     C=["#4B2991", "#952EA0", "#D44292", "#F66D7A", "#F6A97A"]
17     Classes=np.unique(y_test)
18     pred_prob=Naive.predict_proba(X_test)
19     fpr = {}
20     tpr = {}
21     thresh = {}
22     roc_auc = dict()
23     n_classe = len(Classes)
24     for i in range (n_classe):
25         fpr[i], tpr[i], thresh[i] = roc_curve(y_test_binarized[
26         roc_auc[i] = auc (fpr[i], tpr[i])
27     # traçage
28     plt.plot(fpr[i], tpr[i], label='%s vs Rest (AUC=%0.2f)'%
29     plt.plot([0,1], [0,1], linestyle='--', c='black')
30     plt.xlim([0,1])
31     plt.ylim([0,1.05])
32     plt.title('Receiver Operating Characteristic(ROC) curves')
33     plt.xlabel('False Rate')
34     plt.ylabel('True Rate')
35     plt.legend()
36     plt.show()

```



Méthode 3:DecisionTreeClassifier

Entrée [37]:

```

1 #entropy ,gini
2 AD=DecisionTreeClassifier(criterion = 'entropy',max_depth=10)
3 AD.fit(x_train,y_train)

```

Out[37]:

```

▼
DecisionTreeClassifier
DecisionTreeClassifier(criterion='entropy', max_depth=10)

```

Entrée [39]:

```

1 from prettytable import PrettyTable
2
3 x = PrettyTable(["Model", "Train SCORE", "Test SCORE"])
4 z=str(int(AD.score(x_train,y_train)*100))+ "%"
5 v=str(int(AD.score(x_test,y_test)*100))+ "%"
6 x.add_row([" DecisionTree",z,v])
7 print(x)

```

```

+-----+-----+-----+
|      Model      | Train SCORE | Test SCORE |
+-----+-----+-----+
| DecisionTree    |      100%   |      90%   |
+-----+-----+-----+

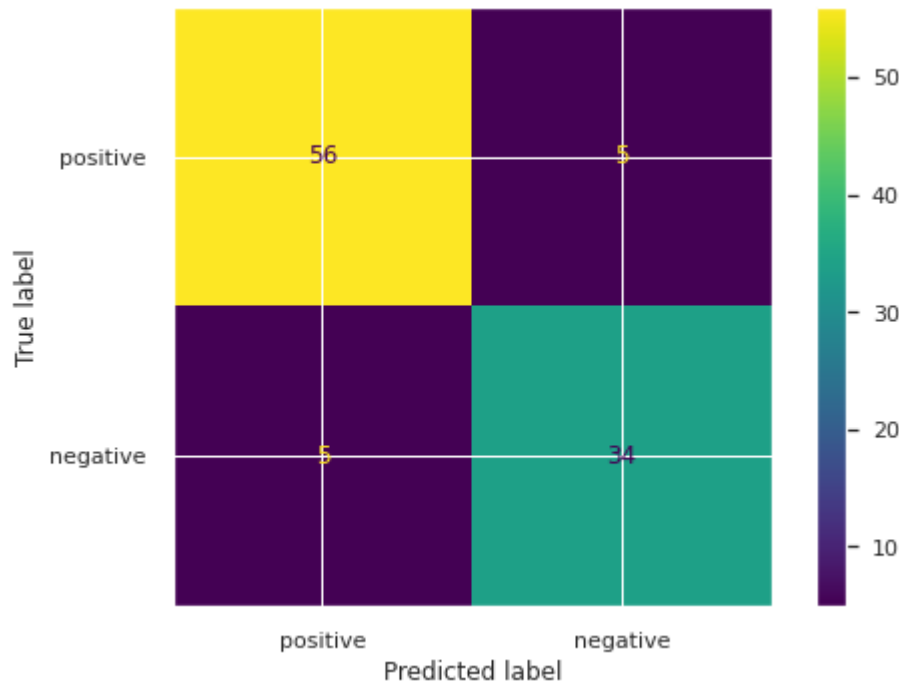
```

Confusion Matrices DecisionTreeClassifier

Entrée [40]:

```
1 y_pred =AD.predict(x_test)
2 print(metrics.classification_report(y_test,y_pred))
3 display=metrics.ConfusionMatrixDisplay((metrics.confusion_matri
4 display.plot()
5 plt.show()
```

	precision	recall	f1-score	support
0	0.92	0.92	0.92	61
1	0.87	0.87	0.87	39
accuracy			0.90	100
macro avg	0.89	0.89	0.89	100
weighted avg	0.90	0.90	0.90	100



Cross-validation: evaluating estimator performance

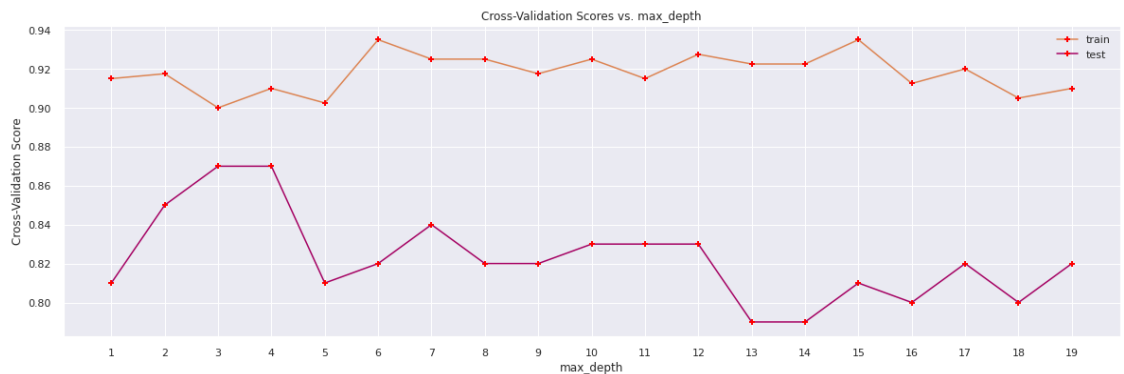
Entrée [45]:

```

1 max_depth_values = range(1, 20)
2 cross_val_scores = []
3 cross_val_scorestest=[]
4 for depth in max_depth_values:
5     clf = DecisionTreeClassifier(criterion = 'entropy',max_depth=depth)
6     scores_train = cross_val_score(clf, x_train, y_train, cv=5)
7     cross_val_scores.append(np.mean(scores_train))
8     scores_test = cross_val_score(clf, x_test, y_test, cv=5)
9     cross_val_scorestest.append(np.mean(scores_test))
10
11 plt.figure(figsize=(20, 6))
12 plt.plot(max_depth_values, cross_val_scores, marker='P',c="C1",label='train')
13 plt.plot(max_depth_values, cross_val_scorestest, marker='P',c='C2',label='test')
14
15 plt.title('Cross-Validation Scores vs. max_depth')
16 plt.xlabel('max_depth')
17 plt.ylabel('Cross-Validation Score')
18 plt.xticks(max_depth_values)
19 plt.grid(True)
20 plt.legend()
21 plt.show

```

Out[45]: <function matplotlib.pyplot.show(close=None, block=None)>

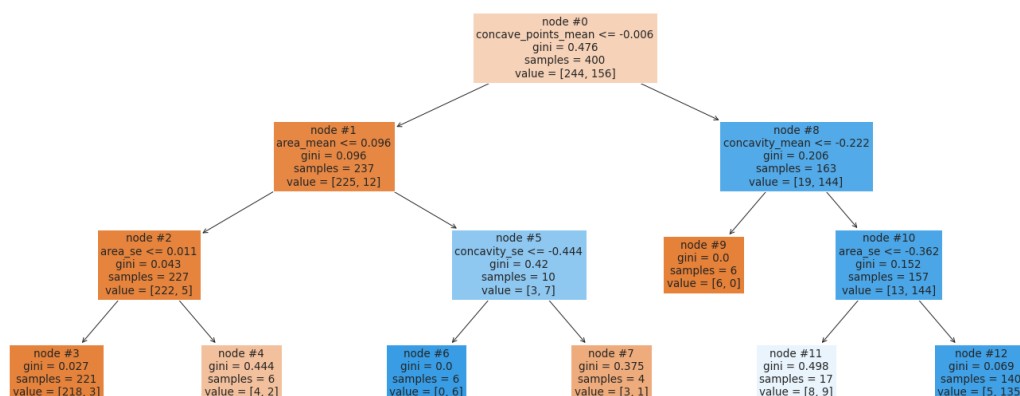


DecisionTreeClassifier

Entrée [43]:

```

1 AD=DecisionTreeClassifier(max_depth = 3)
2 AD.fit(x_train,y_train)
3 plt.figure(figsize=(25,10))
4 plot_tree(AD,feature_names = x_test.columns,filled = True,node_
5 plt.show()
```



la courbe (ROC AUC)

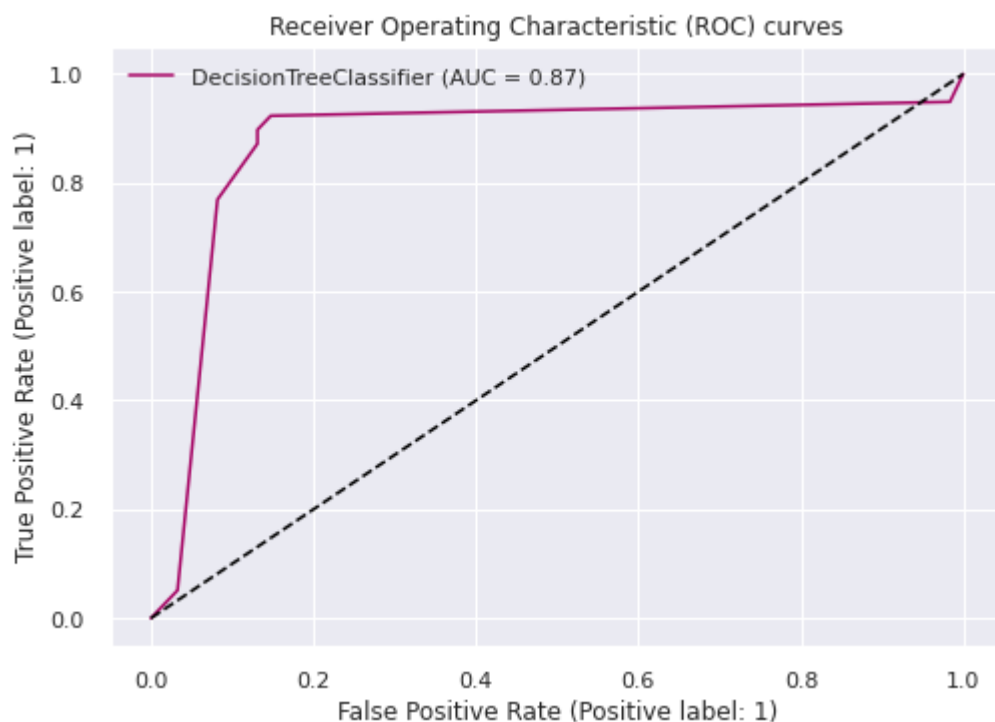
Il s'agit d'une fonction générale, étant donné des points sur une courbe. Pour calculer l'aire sous la courbe ROC

Entrée [46]:

```

1 # Make predictions on the test set
2 y_pred = AD.predict(x_test)
3 from sklearn.preprocessing import label_binarize
4 from sklearn.metrics import roc_curve, auc
5 from sklearn.metrics import DetCurveDisplay, RocCurveDisplay
6 y_test_binarized = label_binarize(y_test, classes=np.unique(y_t
7 if len(y_test_binarized[0])<=1:
8     fig, ax_roc = plt.subplots()
9     RocCurveDisplay.from_estimator(AD, x_test, y_test, ax=ax_ro
10    plt.plot([0,1],[0,1], '--', c='black')
11    ax_roc.set_title("Receiver Operating Characteristic (ROC) c
12    #ax_roc.grid(linestyle="--", c='black')
13    plt.legend()
14    plt.show()
15 elif (len(y_test_binarized[0])>1):
16     C=["#4B2991", "#952EA0", "#D44292", "#F66D7A", "#F6A97A"]
17     Classes=np.unique(y_test)
18     pred_prob=AD.predict_proba(X_test)
19     fpr = {}
20     tpr = {}
21     thresh = {}
22     roc_auc = dict()
23     n_classe = len(Classes)
24     for i in range (n_classe):
25         fpr[i], tpr[i], thresh[i] = roc_curve(y_test_binarized[
26         roc_auc[i] = auc (fpr[i], tpr[i])
27     # traçage
28     plt.plot(fpr[i], tpr[i], label='%s vs Rest (AUC=%0.2f)'%
29     plt.plot([0,1], [0,1], linestyle='--', c='black')
30     plt.xlim([0,1])
31     plt.ylim([0,1.05])
32     plt.title('Receiver Operating Characteristic(ROC) curves')
33     plt.xlabel('False Rate')
34     plt.ylabel('True Rate')
35     plt.legend()
36     plt.show()

```



Méthode 4:Reseau Neurones

Entrée [47]:

```

1 # Create an instance of MLPClassifier
2 mlf = MLPClassifier(hidden_layer_sizes=(45), max_iter=1000)
3 # Train the model
4 mlf.fit(x_train, y_train)

```

Out[47]:

```

▼
MLPClassifier
MLPClassifier(hidden_layer_sizes=45, max_iter=1000)

```

Entrée [48]:

```

1 from prettytable import PrettyTable
2
3 x = PrettyTable(["Model", "Train SCORE", "Test SCORE"])
4 z=str(int(mlf.score(x_train,y_train)*100))+ "%"
5 v=str(int(mlf.score(x_test,y_test)*100))+ "%"
6 x.add_row([" Reseau Neurones",z,v])
7 print(x)

```

```

+-----+-----+-----+
|      Model      | Train SCORE | Test SCORE |
+-----+-----+-----+
| Reseau Neurones |      99%    |      92%    |
+-----+-----+-----+

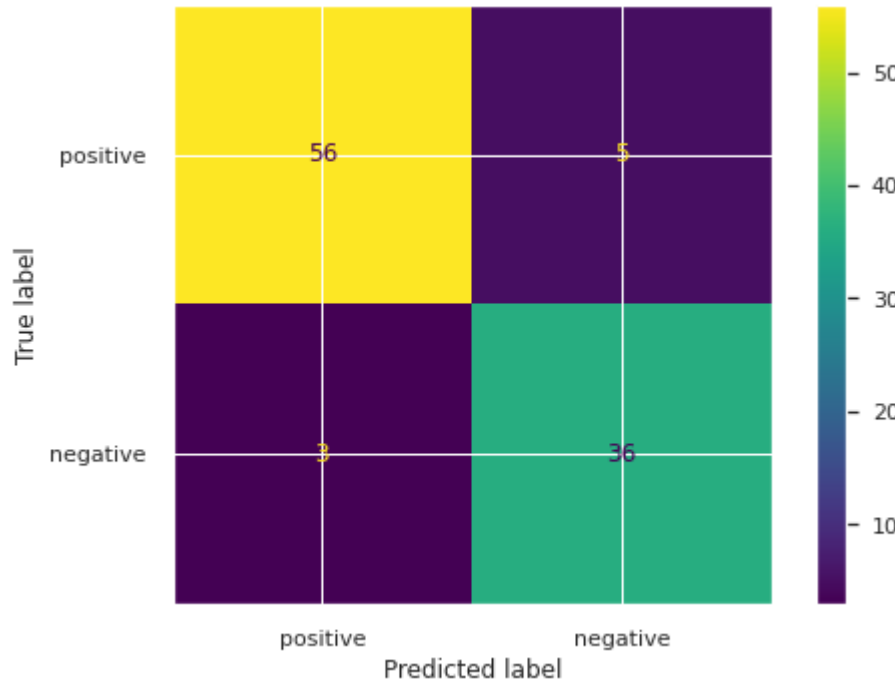
```


Confusion Matrices Réseau Neurones

Entrée [49]:

```
1 y_pred =mlf.predict(x_test)
2 print(metrics.classification_report(y_test,y_pred))
3 display=metrics.ConfusionMatrixDisplay((metrics.confusion_matri
4 display.plot()
5 plt.show()
6
```

	precision	recall	f1-score	support
0	0.95	0.92	0.93	61
1	0.88	0.92	0.90	39
accuracy			0.92	100
macro avg	0.91	0.92	0.92	100
weighted avg	0.92	0.92	0.92	100



Cross-validation: evaluating estimator performance

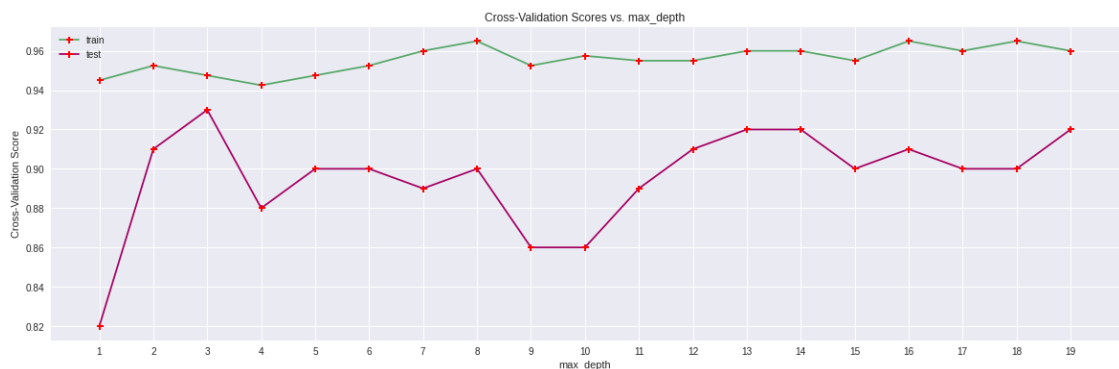
Entrée [59]:

```

1 max_depth_values = range(1, 20)
2 cross_val_scores = []
3 cross_val_scorestest=[]
4 for depth in max_depth_values:
5     clf = mlf = MLPClassifier(hidden_layer_sizes=(depth), max_i
6     scores_train = cross_val_score(clf, x_train, y_train, cv=5)
7     cross_val_scores.append(np.mean(scores_train))
8     scores_test = cross_val_score(clf, x_test, y_test, cv=5)
9     cross_val_scorestest.append(np.mean(scores_test))
10
11 plt.figure(figsize=(20, 6))
12 plt.plot(max_depth_values, cross_val_scores, marker='P',c="C1",)
13 plt.plot(max_depth_values, cross_val_scorestest, marker='P',c=':
14
15 plt.title('Cross-Validation Scores vs. max_depth')
16 plt.xlabel('max_depth')
17 plt.ylabel('Cross-Validation Score')
18 plt.xticks(max_depth_values)
19 plt.grid(True)
20 plt.legend()
21 plt.show

```

Out[59]: <function matplotlib.pyplot.show(close=None, block=None)>



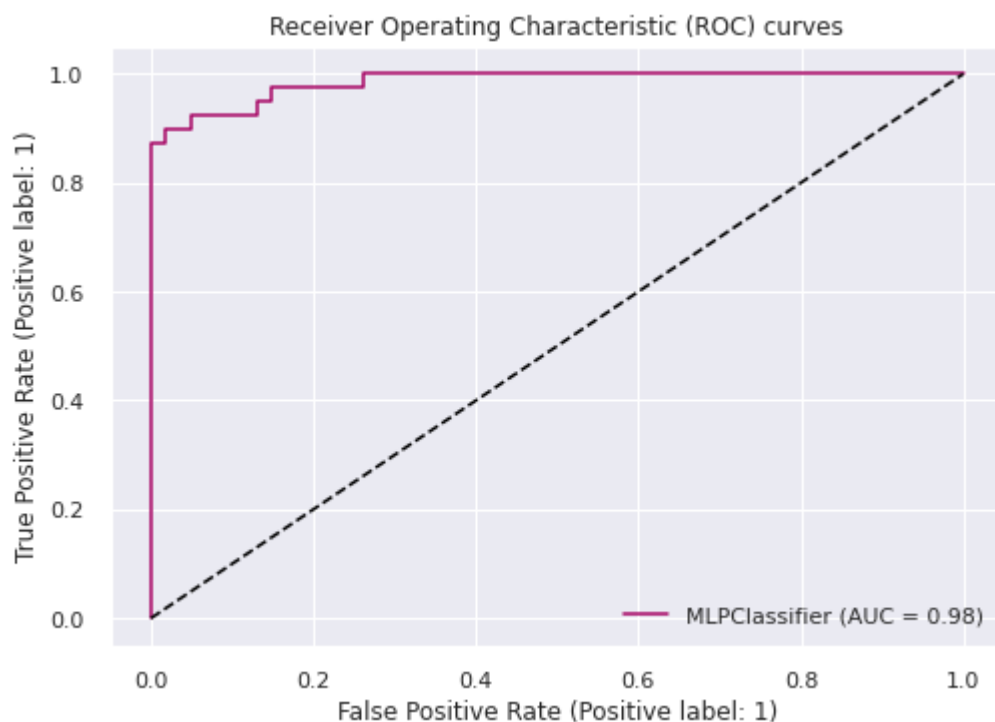
la courbe (ROC AUC)

Il s'agit d'une fonction générale, étant donné des points sur une courbe. Pour calculer l'aire sous la courbe ROC

```

Entrée [50]: 1 # Make predictions on the test set
2 y_pred = mlf.predict(x_test)
3 from sklearn.preprocessing import label_binarize
4 from sklearn.metrics import roc_curve, auc
5 from sklearn.metrics import DetCurveDisplay, RocCurveDisplay
6 y_test_binarized = label_binarize(y_test, classes=np.unique(y_t
7 if len(y_test_binarized[0])<=1:
8     fig, ax_roc = plt.subplots()
9     RocCurveDisplay.from_estimator(mlf, x_test, y_test, ax=ax_r
10     plt.plot([0,1],[0,1], '--', c='black')
11     ax_roc.set_title("Receiver Operating Characteristic (ROC) c
12     #ax_roc.grid(linestyle="--", c='black')
13     plt.legend()
14     plt.show()
15 elif (len(y_test_binarized[0])>1):
16     C=["#4B2991", "#952EA0", "#D44292", "#F66D7A", "#F6A97A"]
17     Classes=np.unique(y_test)
18     pred_prob=mlf.predict_proba(X_test)
19     fpr = {}
20     tpr = {}
21     thresh = {}
22     roc_auc = dict()
23     n_classe = len(Classes)
24     for i in range (n_classe):
25         fpr[i], tpr[i], thresh[i] = roc_curve(y_test_binarized[
26         roc_auc[i] = auc (fpr[i], tpr[i])
27     # traçage
28     plt.plot(fpr[i], tpr[i], label='%s vs Rest (AUC=%0.2f)'%
29     plt.plot([0,1], [0,1], linestyle='--', c='black')
30     plt.xlim([0,1])
31     plt.ylim([0,1.05])
32     plt.title('Receiver Operating Characteristic(ROC) curves')
33     plt.xlabel('False Rate')
34     plt.ylabel('True Rate')
35     plt.legend()
36     plt.show()

```



Méthode 5:SVM

Entrée [51]:

```

1  # Create an instance of MLPClassifier
2
3  svm = svm.SVC(kernel='rbf',C=0.1)
4  # Train the model
5  svm.fit(x_train, y_train)

```

Out[51]:

```

▼ SVC
SVC(C=0.1)

```

Entrée [52]:

```

1  from prettytable import PrettyTable
2
3  x = PrettyTable(["Model", "Train SCORE", "Test SCORE"])
4  z=str(int(svm.score(x_train,y_train)*100))+ "%"
5  v=str(int(svm.score(x_test,y_test)*100))+ "%"
6  x.add_row([" SVM",z,v])
7  print(x)

```

```

+-----+-----+-----+
| Model | Train SCORE | Test SCORE |
+-----+-----+-----+
|  SVM |      94%    |      90%    |
+-----+-----+-----+

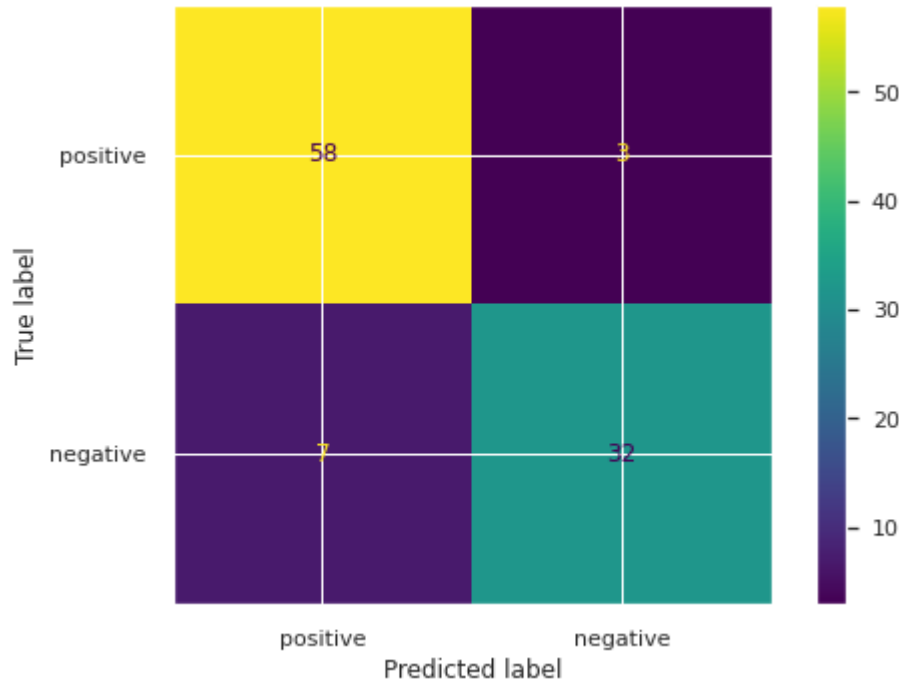
```

Confusion Matrices SVM

Entrée [53]:

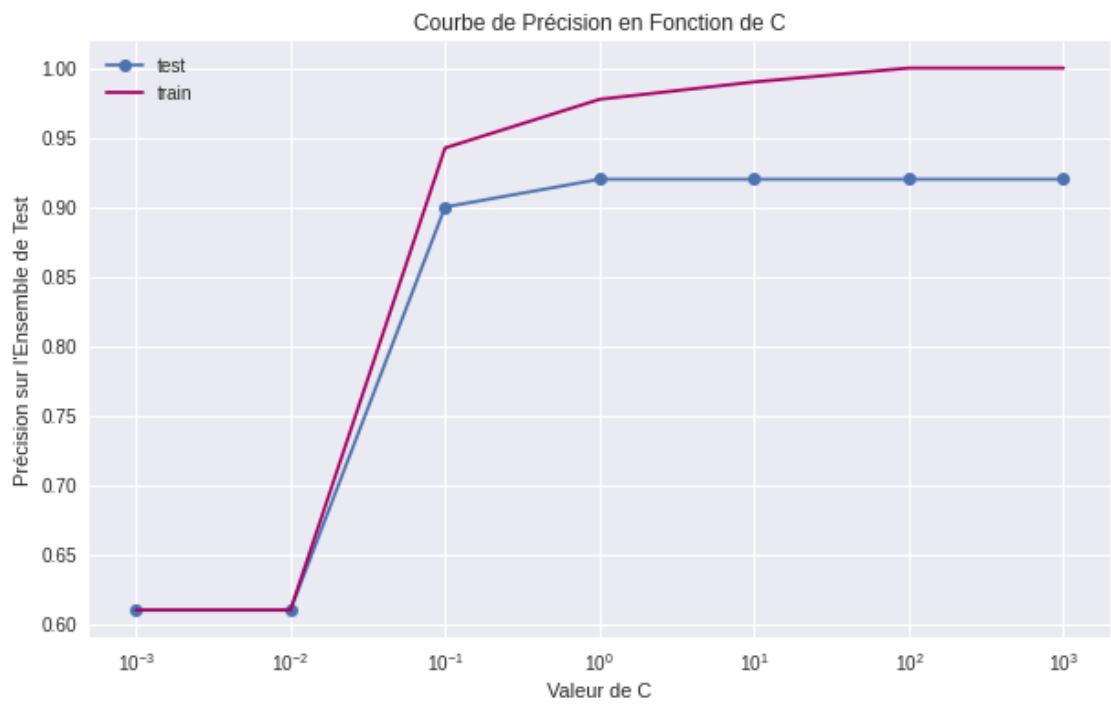
```
1 y_pred =svm.predict(x_test)
2 print(metrics.classification_report(y_test,y_pred))
3 display=metrics.ConfusionMatrixDisplay((metrics.confusion_matri
4 display.plot()
5 plt.show()
```

	precision	recall	f1-score	support
0	0.89	0.95	0.92	61
1	0.91	0.82	0.86	39
accuracy			0.90	100
macro avg	0.90	0.89	0.89	100
weighted avg	0.90	0.90	0.90	100



Cross-validation: evaluating estimator performance

```
Entrée [56]: 1 #Définir les valeurs de C que vous souhaitez tester
2 from sklearn.metrics import accuracy_score
3 C_values = np.logspace(-3, 3, 7)
4
5 # Entraîner le modèle pour différentes valeurs de C et enregistrer
6 accuracies = []
7 train=[]
8 for C in C_values:
9     # Créer un classificateur SVM avec noyau linéaire et la validation croisée
10     svm_model = svm.SVC(kernel='rbf', C=C)
11
12     # Entraîner le modèle sur l'ensemble d'entraînement
13     svm_model.fit(x_train, y_train)
14
15     # Faire des prédictions sur l'ensemble de test
16     y_pred = svm_model.predict(x_test)
17     y_predt = svm_model.predict(x_train)
18
19     # Calculer la précision et l'enregistrer
20     accuracy = accuracy_score(y_test, y_pred)
21     accuracyt = accuracy_score(y_train, y_predt)
22     accuracies.append(accuracy)
23     train.append(accuracyt)
24
25 # Tracer la courbe de précision en fonction de C
26 plt.figure(figsize=(10, 6))
27 plt.semilogx(C_values, accuracies, marker='o',label='test')
28 plt.semilogx(C_values, train, c='#A50062',label='train',marker='x')
29 plt.title('Courbe de Précision en Fonction de C')
30 plt.xlabel('Valeur de C')
31 plt.ylabel('Précision sur l\'Ensemble de Test')
32 plt.grid(True)
33 plt.legend()
34 plt.show()
35
36 # Trouver la meilleure valeur de C
37 best_C = C_values[np.argmax(accuracies)]
38 print(f"Meilleur paramètre C : {best_C}")
39
```



Meilleur paramètre C : 1.0

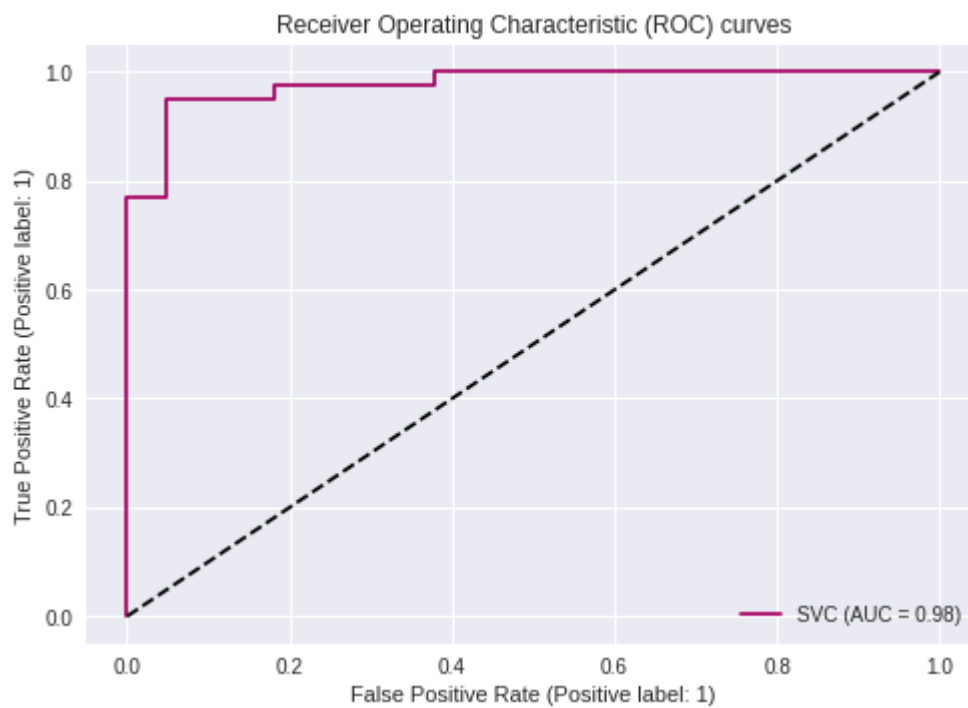
la courbe (ROC AUC)

Entrée [57]:

```

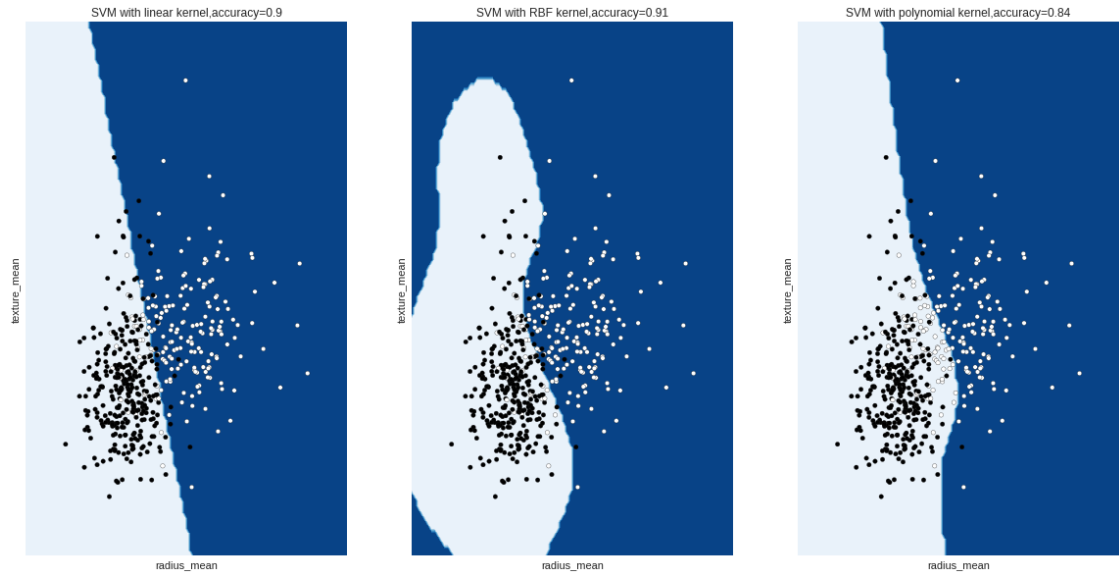
1 SVM = svm.SVC(probability=True)
2 SVM.fit(x_train, y_train)
3
4 # Make predictions on the test set
5 y_pred = SVM.predict(x_test)
6 from sklearn.preprocessing import label_binarize
7 from sklearn.metrics import roc_curve, auc
8 from sklearn.metrics import DetCurveDisplay, RocCurveDisplay
9 y_test_binarized = label_binarize(y_test, classes=np.unique(y_t
10 if len(y_test_binarized[0])<=1:
11     fig, ax_roc = plt.subplots()
12     RocCurveDisplay.from_estimator(SVM, x_test, y_test, ax=ax_r
13     plt.plot([0,1],[0,1], '--', c='black')
14     ax_roc.set_title("Receiver Operating Characteristic (ROC) c
15     #ax_roc.grid(linestyle="--", c='black')
16     plt.legend()
17     plt.show()
18 elif (len(y_test_binarized[0])>1):
19     C=["#4B2991", "#952EA0", "#D44292", "#F66D7A", "#F6A97A"]
20     Classes=np.unique(y_test)
21     pred_prob=SVM.predict_proba(x_test)
22     fpr = {}
23     tpr = {}
24     thresh = {}
25     roc_auc = dict()
26     n_classe = len(Classes)
27     for i in range (n_classe):
28         fpr[i], tpr[i], thresh[i] = roc_curve(y_test_binarized[
29         roc_auc[i] = auc (fpr[i], tpr[i])
30     # traçage
31     plt.plot(fpr[i], tpr[i], label='%s vs Rest (AUC=%0.2f)' %
32     plt.plot([0,1], [0,1], linestyle='--', c='black')
33     plt.xlim([0,1])
34     plt.ylim([0,1.05])
35     plt.title('Receiver Operating Characteristic(ROC) curves')
36     plt.xlabel('False Rate')
37     plt.ylabel('True Rate')
38     plt.legend()
39     plt.show()

```

Comparison des 3 kernel de SVM

```
Entrée [58]: 1 import matplotlib.pyplot as plt
2
3 from sklearn import datasets, svm
4 from sklearn.inspection import DecisionBoundaryDisplay
5
6 XX = X.iloc[:, 0:2]
7 YY = Y
8
9 # we create an instance of SVM and fit out data. We do not scale
10 # data since we want to plot the support vectors
11 C = 1.0 # SVM regularization parameter
12 models = (
13     svm.SVC(kernel="linear", C=C),
14     svm.SVC(kernel="rbf", gamma=0.7, C=C),
15     svm.SVC(kernel="poly", degree=3, gamma="auto", C=C),
16 )
17 models = (clf.fit(XX, YY) for clf in models)
18
19 # title for the plots
20 titles = (
21     "SVM with linear kernel",
22     "SVM with RBF kernel",
23     "SVM with polynomial kernel",
24 )
25
26
27 # Set-up 2x2 grid for plotting.
28 fig, sub = plt.subplots(1, 3, figsize=(20, 10))
29
30
31 #plt.subplots_adjust(wspace=0.4, hspace=0.4)
32
33 X0, X1 = XX.iloc[:, 0], XX.iloc[:, 1]
34
35 for clf, title, ax in zip(models, titles, sub.flatten()):
36     disp = DecisionBoundaryDisplay.from_estimator(
37         clf,
38         XX,
39         response_method="predict",
40         cmap=plt.cm.Blues,
41         alpha=1,
42         ax=ax,
43     )
44     m=round(clf.score(XX,YY),2)
45     ax.scatter(X0, X1, c=YY, s=20, cmap=plt.cm.gist_gray, edgeco
46
47
48     ax.set_xticks(())
49     ax.set_yticks(())
50     ax.set_title(f"{title},accuracy={m}")
51
52 plt.show()
```



Entrée []:

1	
---	--