

Mieux comprendre le fonctionnement d'un classifieur

Comme nous l'avons vu dans les premières classifications, le comportement des classifieurs peut être très différent. Nous présentons à présent les régions de décisions dans lesquelles le classifieur recherche les valeurs pour pouvoir prédire. L'objectif est donc ici de mieux comprendre le fonctionnement d'un classifieur, les raisons d'une bonne ou mauvaise classification et avoir une idée de l'impact des hyperparamètres.

Dans un problème de classification à deux classes, une région de décision ou surface de décision est une hypersurface qui partitionne l'espace vectoriel sous-jacent en deux ensembles : un pour chaque classe. Le classificateur classera tous les points d'un côté de la limite de décision comme appartenant à une classe et tous ceux de l'autre côté comme appartenant à l'autre classe.

Les illustrations sont faites à partir du jeu de données IRIS et nous retenons celui disponible dans scikit learn et qui a été introduit à la fin du notebook Ingénierie des données.

Installation

Avant de commencer, il est nécessaire de déjà posséder dans son environnement toutes les librairies utiles. Dans la seconde cellule nous importons toutes les librairies qui seront utiles à ce notebook. Il se peut que, lorsque vous lanciez l'exécution de cette cellule, une soit absente. Dans ce cas il est nécessaire de l'installer. Pour cela dans la cellule suivante utiliser la commande :

```
! pip install nom_librairie
```

Attention : il est fortement conseillé lorsque l'une des librairies doit être installer de relancer le kernel de votre notebook.

Remarque : même si toutes les librairies sont importées dès le début, les librairies utiles pour des fonctions présentées au cours de ce notebook sont ré-importées de manière à indiquer d'où elles viennent et ainsi faciliter la réutilisation de la fonction dans un autre projet.

```
In [1]: # utiliser cette cellule pour installer les librairies manquantes
# pour cela il suffit de taper dans cette cellule : !pip install no
m_librairie_manquante
# d'exécuter la cellule et de relancer la cellule suivante pour voi
r si tout se passe bien
# recommencer tant que toutes les librairies ne sont pas installées
...

#!pip install ..

# ne pas oublier de relancer le kernel du notebook
```

```
In [2]: # Importation des différentes librairies utiles pour le notebook

#Sickit learn met régulièrement à jour des versions et
#indique des futurs warnings.
#ces deux lignes permettent de ne pas les afficher.
import warnings
warnings.filterwarnings("ignore", category=FutureWarning)

# librairies générales
import pickle # pour charger le modèle
import pandas as pd
import string
from random import randint
import re
from tabulate import tabulate
import time
import numpy as np
import base64
import sys

# librairie affichage
import matplotlib.pyplot as plt
import seaborn as sns
from mlxtend.plotting import plot_decision_regions
import matplotlib.gridspec as gridspec
import itertools

# librairies scikit learn
from sklearn import datasets
from sklearn.metrics import accuracy_score
from sklearn.metrics import confusion_matrix
from sklearn.metrics import classification_report
from sklearn.model_selection import train_test_split
from sklearn.base import BaseEstimator, TransformerMixin
from sklearn.pipeline import Pipeline
from sklearn.svm import SVC
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score
from sklearn.linear_model import LogisticRegression
from sklearn.naive_bayes import GaussianNB
from sklearn.tree import DecisionTreeClassifier
from sklearn.svm import LinearSVC
```

Pour pouvoir lire et sauvegarder sur votre répertoire Google Drive, il est nécessaire de fournir une autorisation. Pour cela il suffit d'exécuter la ligne suivante et de saisir le code donné par Google.

```
In [3]: from google.colab import drive
drive.mount('/content/gdrive/')
```

Mounted at /content/gdrive/

Corriger éventuellement la ligne ci-dessous pour mettre le chemin vers un répertoire spécifique dans votre répertoire google drive :

```
In [4]: my_local_drive='/content/gdrive/My Drive/Colab Notebooks/ML_FDS'
# Ajout du path pour les librairies, fonctions et données
sys.path.append(my_local_drive)
# Se positionner sur le répertoire associé
%cd $my_local_drive

%pwd
```

```
/content/gdrive/My Drive/Colab Notebooks/ML_FDS
```

```
Out[4]: '/content/gdrive/My Drive/Colab Notebooks/ML_FDS'
```

Pour pouvoir afficher l'espace de décision, le plus simple est de se mettre dans un espace en 2 dimensions. Par la suite nous ne considérerons que les 2 attributs *sepal length* et *sepal width* dans notre jeu de données.

```
In [5]: print ('Lecture du fichier iris\n')
iris = datasets.load_iris()
#sélection des deux attributs sepals
X = iris.data[:, :2]
y = iris.target
```

```
Lecture du fichier iris
```

Affichage des valeurs des attributs afin de voir comment elles se répartissent.

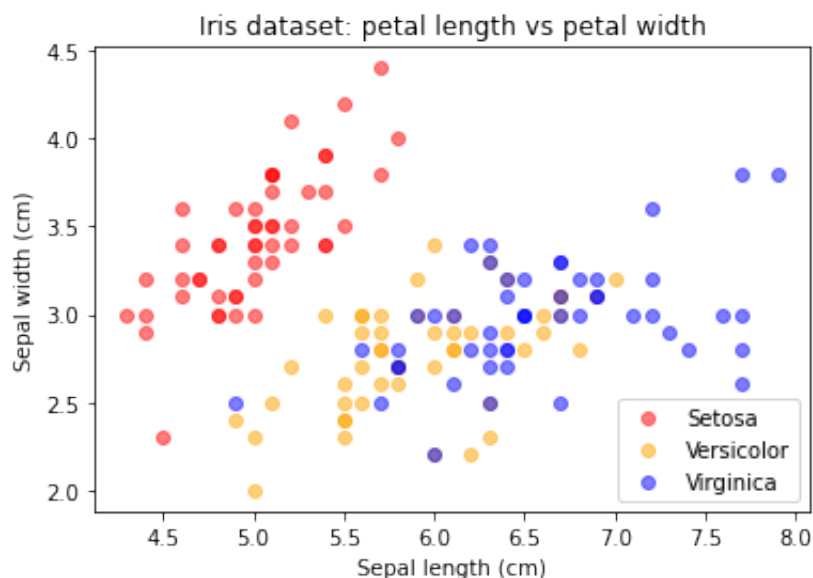
```
In [6]: #Passage par un dataframe par soucis de simplification
iris_df = pd.DataFrame(iris['data'], columns=iris['feature_names'])
iris_df['species'] = iris['target']

colours = ['red', 'orange', 'blue']
species = ['Setosa', 'Versicolor', 'Virginica']

for i in range(0, 3):
    species_df = iris_df[iris_df['species'] == i]
    plt.scatter(
        species_df['sepal length (cm)'],
        species_df['sepal width (cm)'],
        cmap=plt.cm.coolwarm,
        color=colours[i],
        alpha=0.5,
        label=species[i]
    )

plt.xlabel('Sepal length (cm)')
plt.ylabel('Sepal width (cm)')
plt.title('Iris dataset: petal length vs petal width')
plt.legend(loc='lower right')

plt.show()
```



Comme nous pouvons le constater Setosa est très séparée des autres classes et doit pouvoir facilement être séparé. Il est évident que la séparation entre Versicolor et Virginica va être plus difficile pour un classifieur.

Création d'un jeu de données d'apprentissage et de test.

```
In [7]: validation_size=0.3 #30% du jeu de données pour le test

testsize= 1-validation_size
seed=30
X_train,X_test,y_train,y_test=train_test_split(X, y,
                                                train_size=validation_size,
                                                test_size=testsize,
                                                random_state=seed)
```

Test de 4 classifieurs différents pour voir comment ils se comportent.

```
In [8]: lr=LogisticRegression(random_state=1,
                               solver='newton-cg',
                               multi_class='multinomial')

gnb = GaussianNB()
deci= DecisionTreeClassifier(random_state=1)
svm = SVC(gamma='auto')
```

Pour afficher les régions de décisions, nous utilisons la librairie mlxtend (cf notebook règles d'association) qui offre de nombreuses facilités pour afficher la zone.

Le principe consiste à parcourir la liste des classifieurs, de faire le fit, de faire la prédiction et d'afficher la zone de décision. Cette zone correspond aux différentes valeurs dans lesquelles pour une classe, le classifieur va chercher ses valeurs. La zone est définie par rapport à l'ensemble des données. La fonction plot_decision_regions va récupérer les valeurs min et max de tous les attributs et ensuite en fonction de la classe va plutôt étendre ou modifier la zone.

Ci-dessous une fonction qui appelle plot_decision_regions.

```

In [9]: def call_decision_regions_iris (labels,list_clf,X_train,y_train,X_test,y_test):
    # pour afficher les points du jeu de test plus clair
    scatter_kwargs = {'s': 120, 'edgecolor': None, 'alpha': 0.7}
    contourf_kwargs = {'alpha': 0.2}
    scatter_highlight_kwargs = {'s': 120, 'label': 'Jeu de test', 'alpha': 0.7}

    #pour afficher les 4 valeurs sur 2 colonnes et 2 lignes
    gs = gridspec.GridSpec(2, 2)

    fig = plt.figure(figsize=(12,12))

    for clf, label, grd in zip(list_clf,
                               labels,
                               itertools.product([0, 1], repeat=2)):
        #fit du modele
        clf.fit(X_train, y_train)

        #prediction sur le jeu de test
        result=clf.predict(X_test)
        acc=accuracy_score(result, y_test)

        #affichage de la zone de décision
        ax = plt.subplot(gs[grd[0], grd[1]])
        fig=plot_decision_regions(X, y, clf=clf, legend=2,
                                X_highlight=X_test,
                                scatter_kwargs=scatter_kwargs,
                                contourf_kwargs=contourf_kwargs,
                                scatter_highlight_kwargs=scatter_highlight_kwargs,
                                args)

        L = plt.legend()
        L.get_texts()[0].set_text('Setosa')
        L.get_texts()[1].set_text('Versicolor')
        L.get_texts()[2].set_text('Virginica')
        accu='%0.3f'%acc
        plt.xlabel('sepal length [cm]')
        plt.ylabel('petal length [cm]')
        label=label+ " (" +accu+') '
        plt.title(label, size=11)
    plt.show()

```

Vous pouvez constater qu'en fonction de la stratégie de l'algorithme les zones sont très différentes.

```
In [10]: labels = ['Logistic Regression',
                  'Naive Bayes',
                  'Decision Tree',
                  'SVM']
list_clf=[lr,gnb,deci,svm]

call_decision_regions_iris (labels,list_clf,X_train,y_train,X_test,
y_test)
```

```
/usr/local/lib/python3.7/dist-packages/mlxtend/plotting/decision_r
egions.py:244: MatplotlibDeprecationWarning: Passing unsupported k
eyword arguments to axis() will raise a TypeError in 3.3.
```

```
ax.axis(xmin=xx.min(), xmax=xx.max(), y_min=yy.min(), y_max=yy.m
ax())
```

```
/usr/local/lib/python3.7/dist-packages/mlxtend/plotting/decision_r
egions.py:313: MatplotlibDeprecationWarning: Using a string of sin
gle character colors as a color sequence is deprecated. Use an exp
licit list instead.
```

```
**scatter_highlight_kwargs)
```

```
/usr/local/lib/python3.7/dist-packages/mlxtend/plotting/decision_r
egions.py:244: MatplotlibDeprecationWarning: Passing unsupported k
eyword arguments to axis() will raise a TypeError in 3.3.
```

```
ax.axis(xmin=xx.min(), xmax=xx.max(), y_min=yy.min(), y_max=yy.m
ax())
```

```
/usr/local/lib/python3.7/dist-packages/mlxtend/plotting/decision_r
egions.py:313: MatplotlibDeprecationWarning: Using a string of sin
gle character colors as a color sequence is deprecated. Use an exp
licit list instead.
```

```
**scatter_highlight_kwargs)
```

```
/usr/local/lib/python3.7/dist-packages/mlxtend/plotting/decision_r
egions.py:244: MatplotlibDeprecationWarning: Passing unsupported k
eyword arguments to axis() will raise a TypeError in 3.3.
```

```
ax.axis(xmin=xx.min(), xmax=xx.max(), y_min=yy.min(), y_max=yy.m
ax())
```

```
/usr/local/lib/python3.7/dist-packages/mlxtend/plotting/decision_r
egions.py:313: MatplotlibDeprecationWarning: Using a string of sin
gle character colors as a color sequence is deprecated. Use an exp
licit list instead.
```

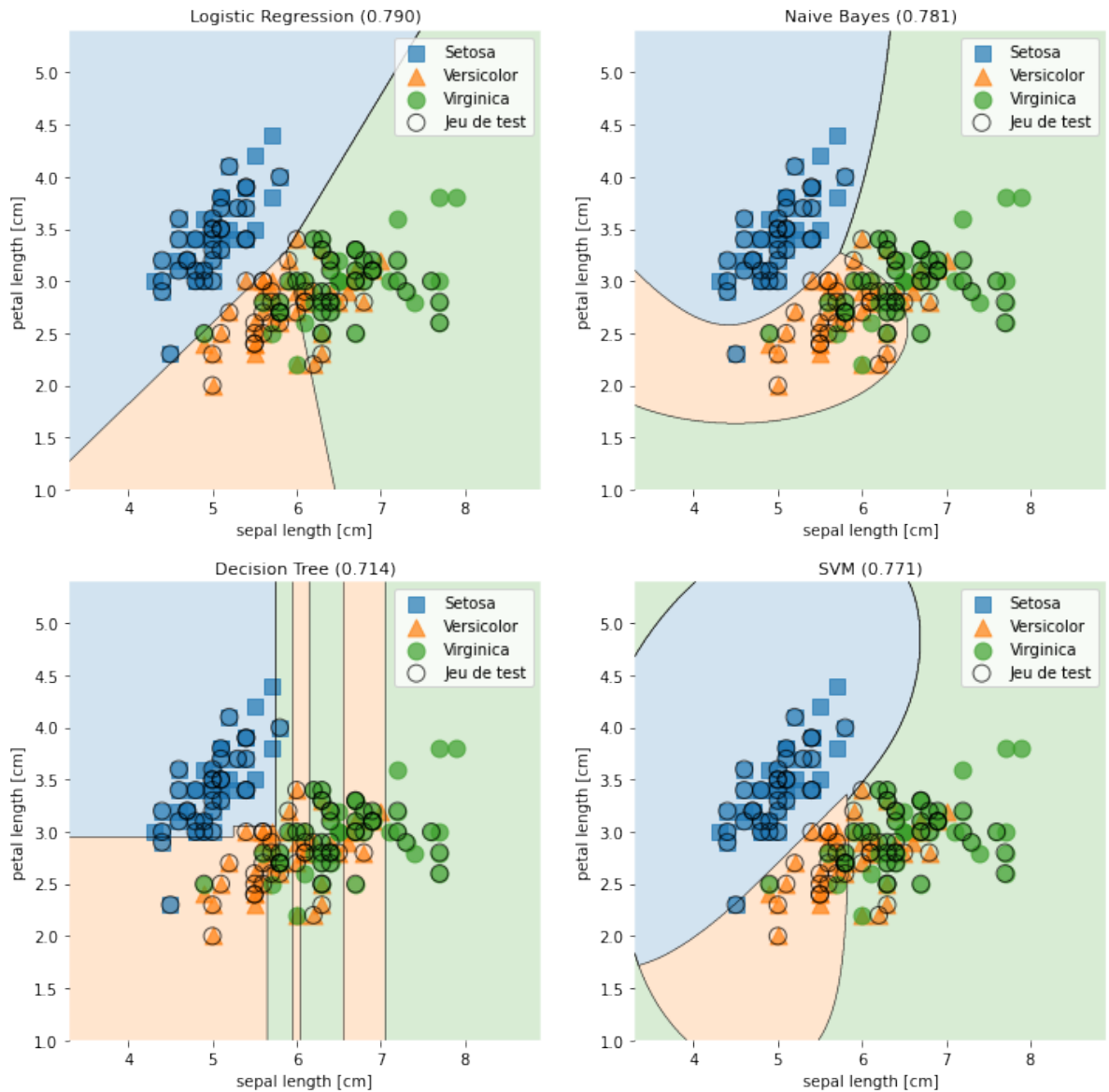
```
**scatter_highlight_kwargs)
```

```
/usr/local/lib/python3.7/dist-packages/mlxtend/plotting/decision_r
egions.py:244: MatplotlibDeprecationWarning: Passing unsupported k
eyword arguments to axis() will raise a TypeError in 3.3.
```

```
ax.axis(xmin=xx.min(), xmax=xx.max(), y_min=yy.min(), y_max=yy.m
ax())
```

```
/usr/local/lib/python3.7/dist-packages/mlxtend/plotting/decision_r
egions.py:313: MatplotlibDeprecationWarning: Using a string of sin
gle character colors as a color sequence is deprecated. Use an exp
licit list instead.
```

```
**scatter_highlight_kwargs)
```

Nous allons, à présent, tester le même classifieur mais des hyperparamètres différents. L'objectif ici est de montrer l'importance de ces choix dans un classifieur.

Dans l'exemple nous prenons le classifieur SVM. Tout d'abord avec un kernel (noyau) linéaire. Dans ce cas, la séparation entre les différentes classes se fait à l'aide de droites. Le second est LinearSVC. Il est un peu similaire au précédent mais l'implémentation est différente notamment sur la prise en compte du choix des pénalités et des fonctions de pertes. Il se comporte mieux que le précédent dans le cas de plus gros volumes de données. Le troisième est SVM avec un noyau de type 'rbf'. Il s'agit d'une Radial Basis Function qui prend comme paramètre gamma (le paramètre qui permet de spécifier la région de décision) et C (la pénalité pour les données mal classées. Si C est petit le classifieur est ok pour les points mal classés). Enfin le dernier considère un polynôme de degré trois.

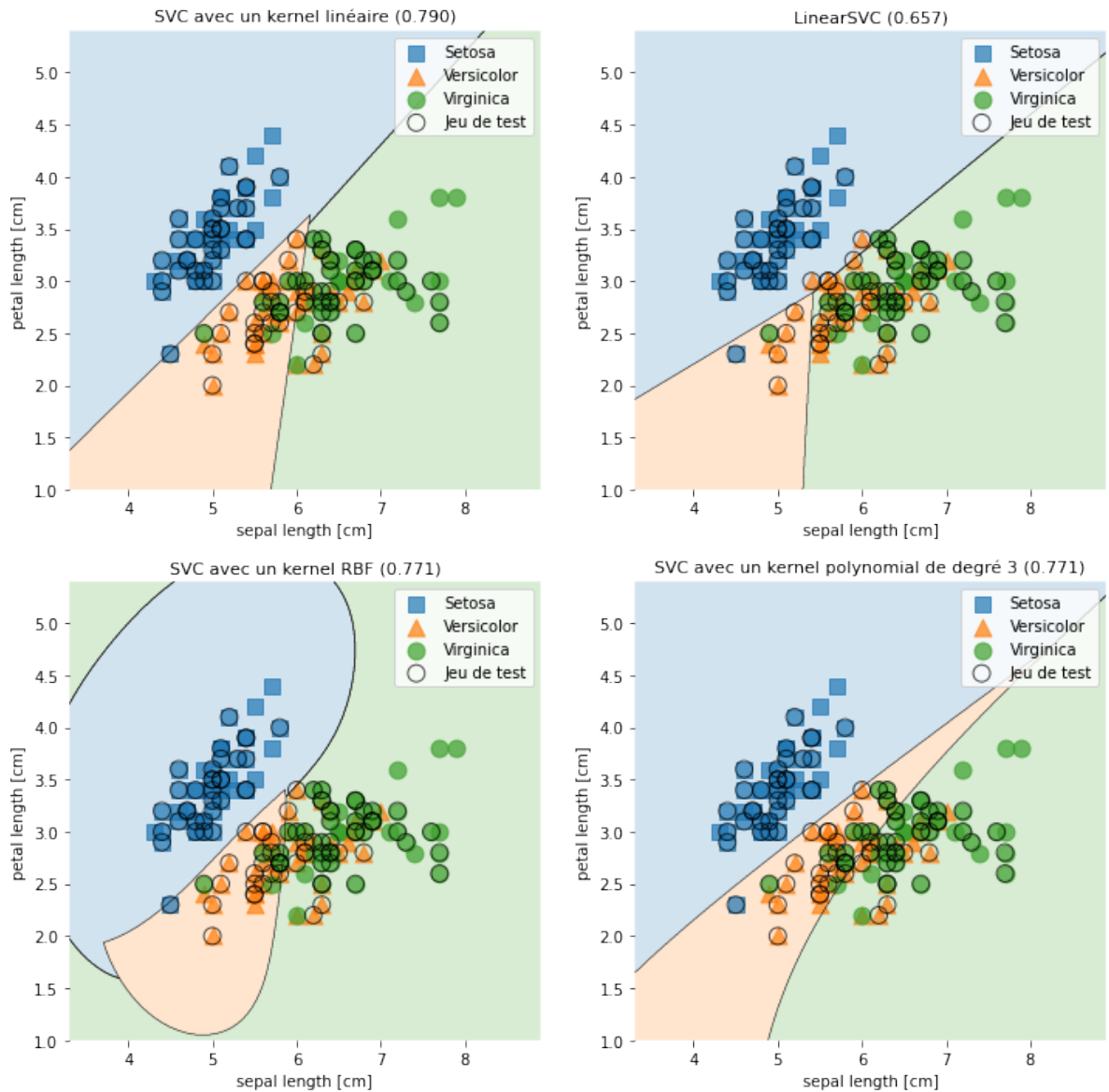
```
In [11]: C = 1.0  # valeur de pénalité
svc = SVC(kernel='linear', C=C)
# LinearSVC (linear kernel)
lin_svc = LinearSVC(max_iter=2000,C=C)
# SVC avec noyau RBF
rbf_svc = SVC(kernel='rbf', gamma=0.7, C=C)
# SVC avec noyau polynomial de degre 3
poly_svc = SVC(kernel='poly', degree=3, C=C)
```

```
In [12]: labels = ['SVC avec un kernel linéaire',
                   'LinearSVC',
                   'SVC avec un kernel RBF',
                   'SVC avec un kernel polynomial de degré 3']
list_clf=[svc,lin_svc,rbf_svc,poly_svc]
call_decision_regions_iris (labels,list_clf,X_train,y_train,X_test,
y_test)
```

```

/usr/local/lib/python3.7/dist-packages/mlxtend/plotting/decision_r
egions.py:244: MatplotlibDeprecationWarning: Passing unsupported k
eyword arguments to axis() will raise a TypeError in 3.3.
    ax.axis(xmin=xx.min(), xmax=xx.max(), y_min=yy.min(), y_max=yy.m
ax())
/usr/local/lib/python3.7/dist-packages/mlxtend/plotting/decision_r
egions.py:313: MatplotlibDeprecationWarning: Using a string of sin
gle character colors as a color sequence is deprecated. Use an exp
licit list instead.
    **scatter_highlight_kwargs)
/usr/local/lib/python3.7/dist-packages/mlxtend/plotting/decision_r
egions.py:244: MatplotlibDeprecationWarning: Passing unsupported k
eyword arguments to axis() will raise a TypeError in 3.3.
    ax.axis(xmin=xx.min(), xmax=xx.max(), y_min=yy.min(), y_max=yy.m
ax())
/usr/local/lib/python3.7/dist-packages/mlxtend/plotting/decision_r
egions.py:313: MatplotlibDeprecationWarning: Using a string of sin
gle character colors as a color sequence is deprecated. Use an exp
licit list instead.
    **scatter_highlight_kwargs)
/usr/local/lib/python3.7/dist-packages/mlxtend/plotting/decision_r
egions.py:244: MatplotlibDeprecationWarning: Passing unsupported k
eyword arguments to axis() will raise a TypeError in 3.3.
    ax.axis(xmin=xx.min(), xmax=xx.max(), y_min=yy.min(), y_max=yy.m
ax())
/usr/local/lib/python3.7/dist-packages/mlxtend/plotting/decision_r
egions.py:313: MatplotlibDeprecationWarning: Using a string of sin
gle character colors as a color sequence is deprecated. Use an exp
licit list instead.
    **scatter_highlight_kwargs)
/usr/local/lib/python3.7/dist-packages/mlxtend/plotting/decision_r
egions.py:244: MatplotlibDeprecationWarning: Passing unsupported k
eyword arguments to axis() will raise a TypeError in 3.3.
    ax.axis(xmin=xx.min(), xmax=xx.max(), y_min=yy.min(), y_max=yy.m
ax())
/usr/local/lib/python3.7/dist-packages/mlxtend/plotting/decision_r
egions.py:313: MatplotlibDeprecationWarning: Using a string of sin
gle character colors as a color sequence is deprecated. Use an exp
licit list instead.
    **scatter_highlight_kwargs)

```



Dans l'exemple suivant nous étudions différentes valeurs de gamma pour voir l'impact de SVM avec un kernel rbf.

```
In [13]: C = 1.0

# SVC avec des valeurs différentes de gamma
rbf_svc1 = SVC(kernel='rbf', gamma=1, C=C)
rbf_svc2 = SVC(kernel='rbf', gamma=10, C=C)
rbf_svc3 = SVC(kernel='rbf', gamma=50, C=C)
rbf_svc4 = SVC(kernel='rbf', gamma=100, C=C)
```

```
In [14]: labels = ['SVC RBF gamma=1',
                  'SVC RBF gamma=10',
                  'SVC RBF gamma=50',
                  'SVC RBF gamma=100']
list_clf=[rbf_svc1,rbf_svc2,rbf_svc3,rbf_svc4]
call_decision_regions_iris (labels,list_clf,X_train,y_train,X_test,
y_test)
```

```
/usr/local/lib/python3.7/dist-packages/mlxtend/plotting/decision_r
egions.py:244: MatplotlibDeprecationWarning: Passing unsupported k
eyword arguments to axis() will raise a TypeError in 3.3.
```

```
ax.axis(xmin=xx.min(), xmax=xx.max(), y_min=yy.min(), y_max=yy.m
ax())
```

```
/usr/local/lib/python3.7/dist-packages/mlxtend/plotting/decision_r
egions.py:313: MatplotlibDeprecationWarning: Using a string of sin
gle character colors as a color sequence is deprecated. Use an exp
licit list instead.
```

```
**scatter_highlight_kwargs)
```

```
/usr/local/lib/python3.7/dist-packages/mlxtend/plotting/decision_r
egions.py:244: MatplotlibDeprecationWarning: Passing unsupported k
eyword arguments to axis() will raise a TypeError in 3.3.
```

```
ax.axis(xmin=xx.min(), xmax=xx.max(), y_min=yy.min(), y_max=yy.m
ax())
```

```
/usr/local/lib/python3.7/dist-packages/mlxtend/plotting/decision_r
egions.py:313: MatplotlibDeprecationWarning: Using a string of sin
gle character colors as a color sequence is deprecated. Use an exp
licit list instead.
```

```
**scatter_highlight_kwargs)
```

```
/usr/local/lib/python3.7/dist-packages/mlxtend/plotting/decision_r
egions.py:244: MatplotlibDeprecationWarning: Passing unsupported k
eyword arguments to axis() will raise a TypeError in 3.3.
```

```
ax.axis(xmin=xx.min(), xmax=xx.max(), y_min=yy.min(), y_max=yy.m
ax())
```

```
/usr/local/lib/python3.7/dist-packages/mlxtend/plotting/decision_r
egions.py:313: MatplotlibDeprecationWarning: Using a string of sin
gle character colors as a color sequence is deprecated. Use an exp
licit list instead.
```

```
**scatter_highlight_kwargs)
```

```
/usr/local/lib/python3.7/dist-packages/mlxtend/plotting/decision_r
egions.py:244: MatplotlibDeprecationWarning: Passing unsupported k
eyword arguments to axis() will raise a TypeError in 3.3.
```

```
ax.axis(xmin=xx.min(), xmax=xx.max(), y_min=yy.min(), y_max=yy.m
ax())
```

```
/usr/local/lib/python3.7/dist-packages/mlxtend/plotting/decision_r
egions.py:313: MatplotlibDeprecationWarning: Using a string of sin
gle character colors as a color sequence is deprecated. Use an exp
licit list instead.
```

```
**scatter_highlight_kwargs)
```

