

# 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.

In [1]:

```
1  #Sickit learn met régulièrement à jour des versions et
2  #indique des futurs warnings.
3  #ces deux lignes permettent de ne pas les afficher.
4  import warnings
5  warnings.filterwarnings("ignore", category=FutureWarning)
6  # la ligne ici est ajoutée principalement pour SVC dont des mises à jour
7  # sont annoncées mais jamais mise à jour :)
```

In [2]:

```
1  from sklearn import datasets
2  import numpy as np
3  import pandas as pd
4
5  from sklearn.model_selection import train_test_split
6  from sklearn.metrics import accuracy_score
7
8  from sklearn.linear_model import LogisticRegression
9  from sklearn.naive_bayes import GaussianNB
10 from sklearn.tree import DecisionTreeClassifier
11 from sklearn.svm import SVC
12 from sklearn.svm import LinearSVC
13
14 import matplotlib.pyplot as plt
15 from mlxtend.plotting import plot_decision_regions
16 import matplotlib.gridspec as gridspec
17 import itertools
18
```

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 [3]:

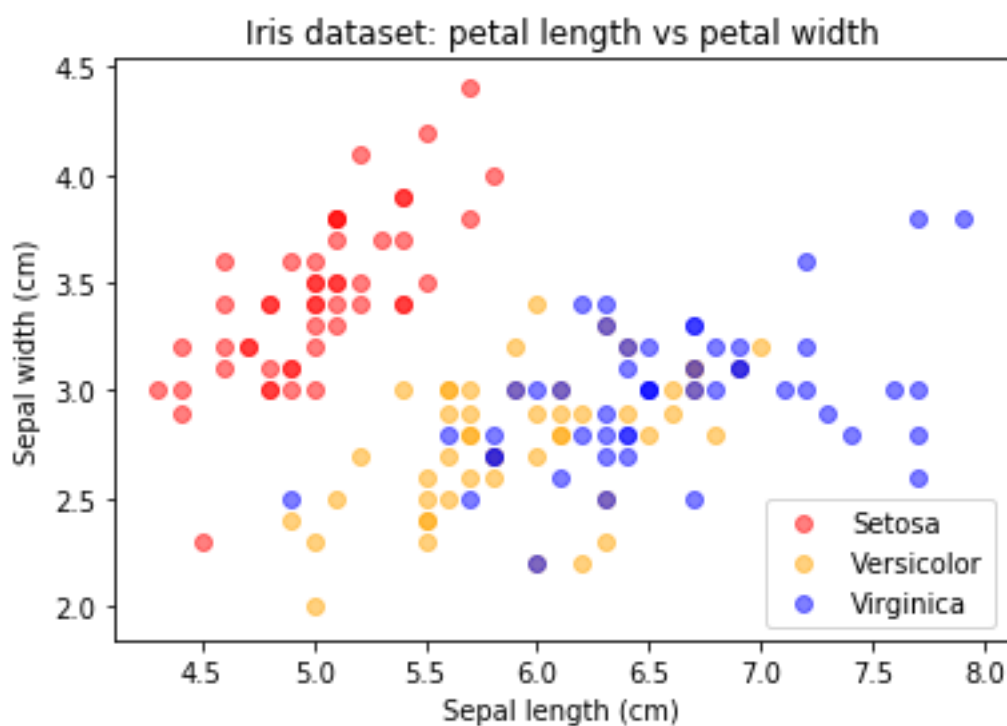
```
1
2  print ('Lecture du fichier iris\n')
3  iris = datasets.load_iris()
4  #sélection des deux attributs sepals
5  X = iris.data[:, :2]
6  y = iris.target
7
```

Lecture du fichier iris

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

In [4]:

```
1
2 #Passage par un dataframe par soucis de simplification
3 iris_df = pd.DataFrame(iris['data'], columns=iris['feature_names'])
4 iris_df['species'] = iris['target']
5
6 colours = ['red', 'orange', 'blue']
7 species = ['Setosa', 'Versicolor', 'Virginica']
8
9 for i in range(0, 3):
10     species_df = iris_df[iris_df['species'] == i]
11     plt.scatter(
12         species_df['sepal length (cm)'],
13         species_df['sepal width (cm)'],
14         cmap=plt.cm.coolwarm,
15         color=colours[i],
16         alpha=0.5,
17         label=species[i]
18     )
19
20 plt.xlabel('Sepal length (cm)')
21 plt.ylabel('Sepal width (cm)')
22 plt.title('Iris dataset: petal length vs petal width')
23 plt.legend(loc='lower right')
24
25 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 [5]:

```
1 validation_size=0.3 #30% du jeu de données pour le test
2
3 testsize= 1-validation_size
4 seed=30
5 ▼ X_train,X_test,y_train,y_test=train_test_split(X, y,
6                                                    train_size=validation_size,
7                                                    random_state=seed,
8                                                    test_size=testsize)
```

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

In [6]:

```
1 ▼ lr=LogisticRegression(random_state=1,
2                           solver='newton-cg',
3                           multi_class='multinomial')
4
5 gnb = GaussianNB()
6 deci= DecisionTreeClassifier(random_state=1)
7 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 clasifieur 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 [7]:

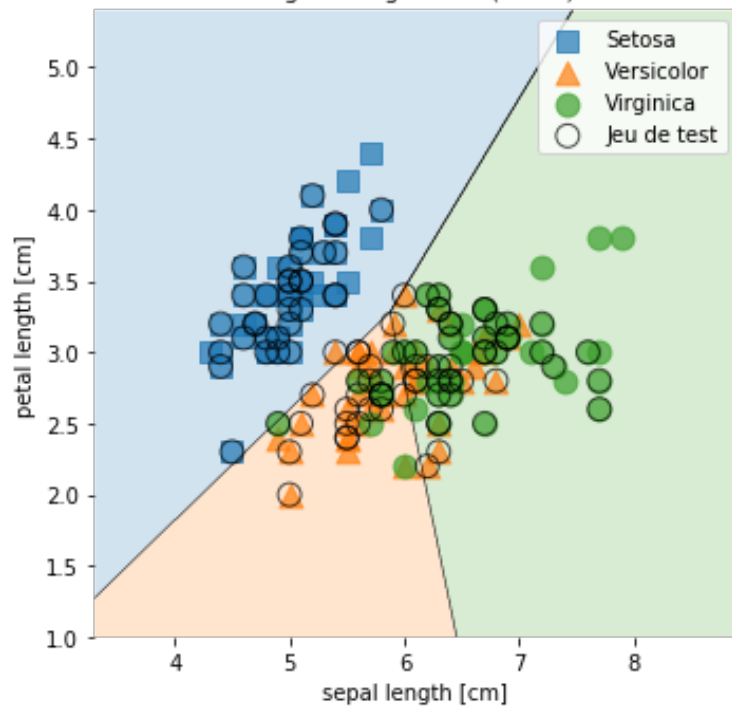
```
1  ▼ def call_decision_regions_iris (labels,list_clf,X_train,y_train,X_test,y_te
2      # pour afficher les points du jeu de test plus clair
3      scatter_kwargs = {'s': 120, 'edgecolor': None, 'alpha': 0.7}
4      contourf_kwargs = {'alpha': 0.2}
5      scatter_highlight_kwargs = {'s': 120, 'label': 'Jeu de test', 'alpha':
6
7      #pour afficher les 4 valeurs sur 2 colonnes et 2 lignes
8      gs = gridspec.GridSpec(2, 2)
9
10     fig = plt.figure(figsize=(12,12))
11
12  ▼ for clf, label, grd in zip(list_clf,
13                               labels,
14                               itertools.product([0, 1], repeat=2)):
15     #fit du modele
16     clf.fit(X_train, y_train)
17
18     #prediction sur le jeu de test
19     result=clf.predict(X_test)
20     acc=accuracy_score(result, y_test)
21
22     #affichage de la zone de décision
23     ax = plt.subplot(gs[grd[0], grd[1]])
24  ▼ fig=plot_decision_regions(X, y, clf=clf, legend=2,
25                             X_highlight=X_test,
26                             scatter_kwargs=scatter_kwargs,
27                             contourf_kwargs=contourf_kwargs,
28                             scatter_highlight_kwargs=scatter_highlight_kwargs)
29
30
31
32     L = plt.legend()
33     L.get_texts()[0].set_text('Setosa')
34     L.get_texts()[1].set_text('Versicolor')
35     L.get_texts()[2].set_text('Virginica')
36     accu='%0.3f'%acc
37     plt.xlabel('sepal length [cm]')
38     plt.ylabel('petal length [cm]')
39     label=label+ " (" +accu+')'
40     plt.title(label, size=11)
41 plt.show()
42
```

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

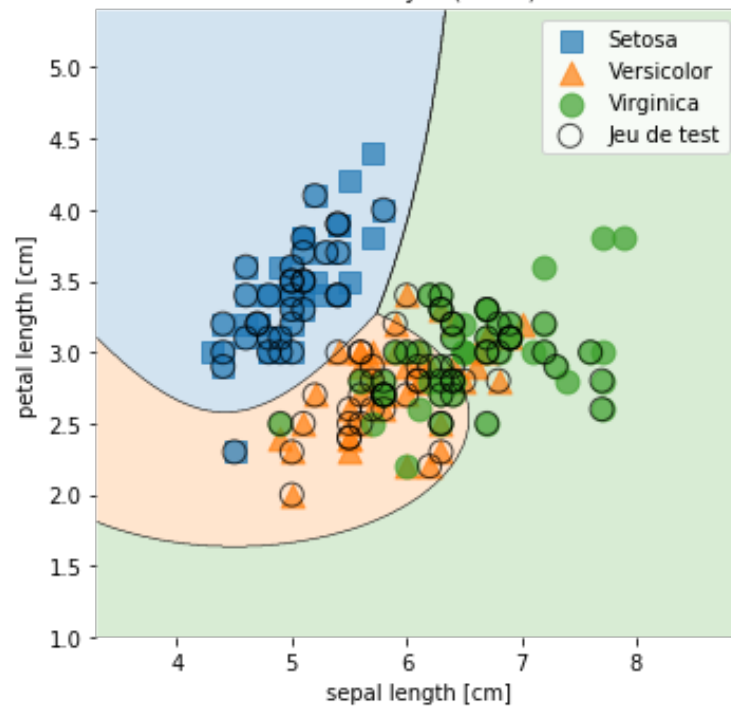
In [8]:

```
1 labels = ['Logistic Regression',  
2           'Naive Bayes',  
3           'Decision Tree',  
4           'SVM']  
5 list_clf=[lr,gnb,deci,svm]  
6  
7 call_decision_regions_iris (labels,list_clf,X_train,y_train,X_test,y_test)
```

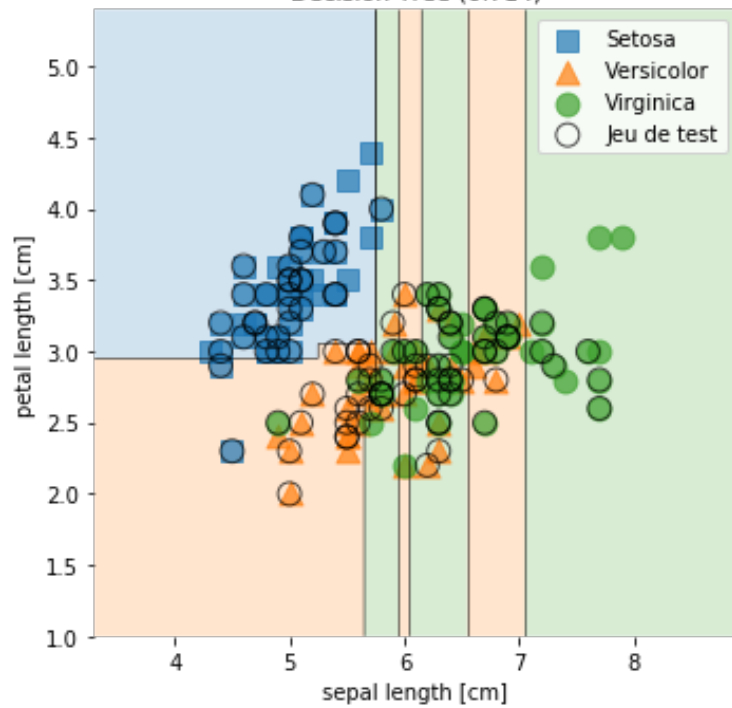
Logistic Regression (0.790)



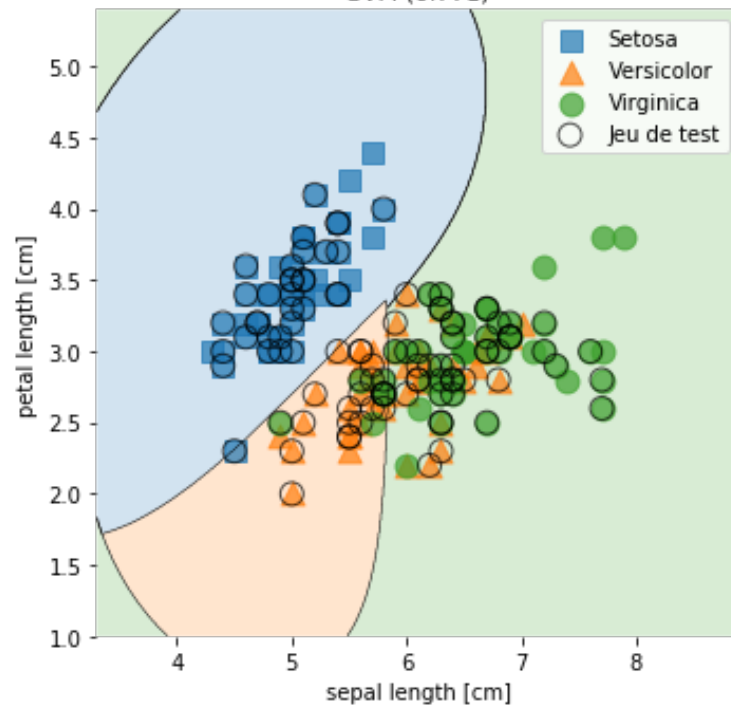
Naive Bayes (0.781)



Decision Tree (0.714)



SVM (0.771)



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 polynome de degré trois.

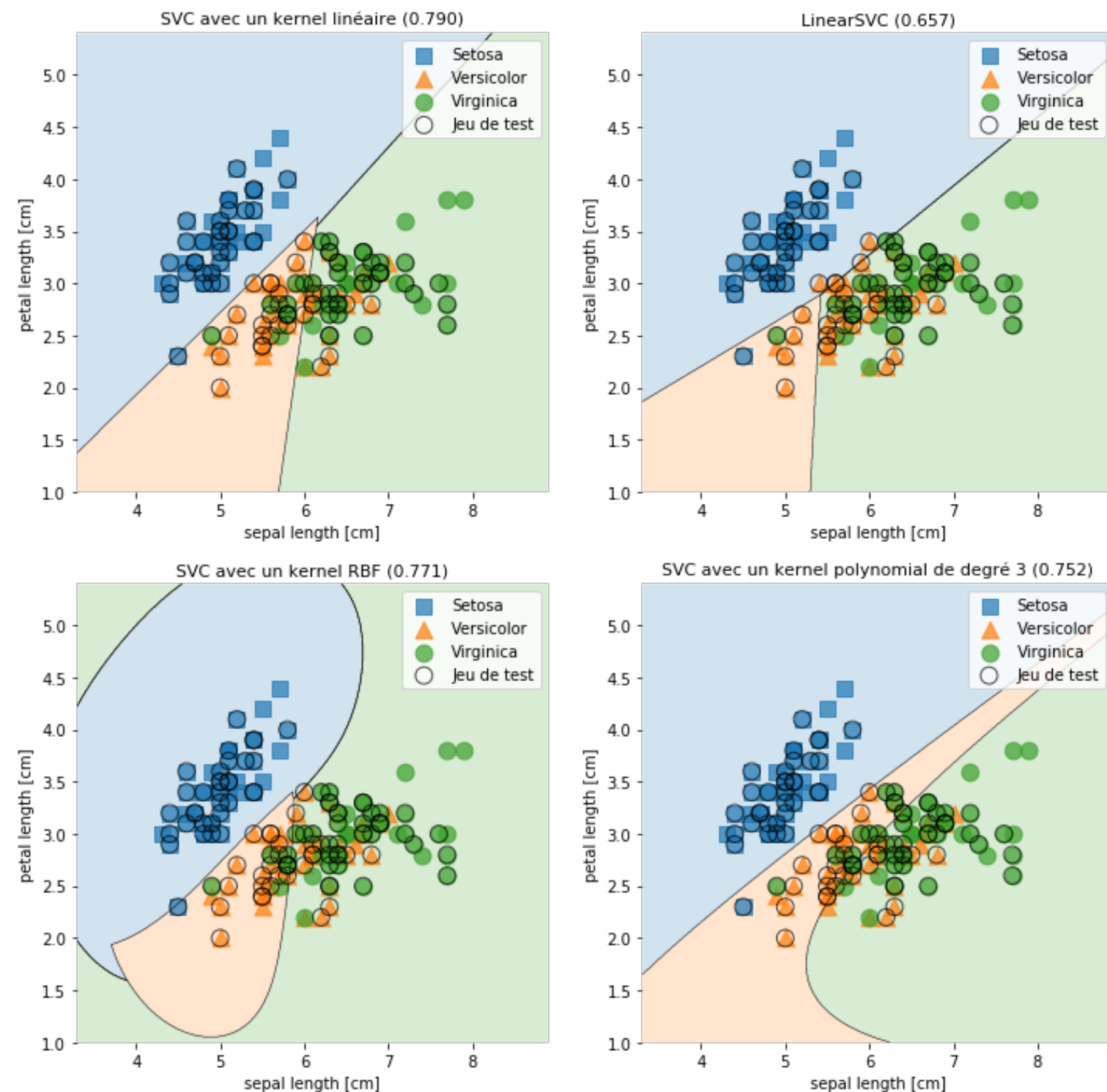
In [9]:

```
1 C = 1.0 # valeur de pénalité
2 svc = SVC(kernel='linear', C=C)
3 # LinearSVC (linear kernel)
4 lin_svc = LinearSVC(max_iter=2000, C=C)
5 # SVC avec noyau RBF
6 rbf_svc = SVC(kernel='rbf', gamma=0.7, C=C)
7 # SVC avec noyau polynomial de degre 3
8 poly_svc = SVC(kernel='poly', degree=3, C=C)
```



In [10]:

```
1 labels = ['SVC avec un kernel linéaire',  
2           'LinearSVC',  
3           'SVC avec un kernel RBF',  
4           'SVC avec un kernel polynomial de degré 3']  
5 list_clf=[svc,lin_svc,rbf_svc,poly_svc]  
6 call_decision_regions_iris (labels,list_clf,X_train,y_train,X_test,y_test)
```



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

In [11]:

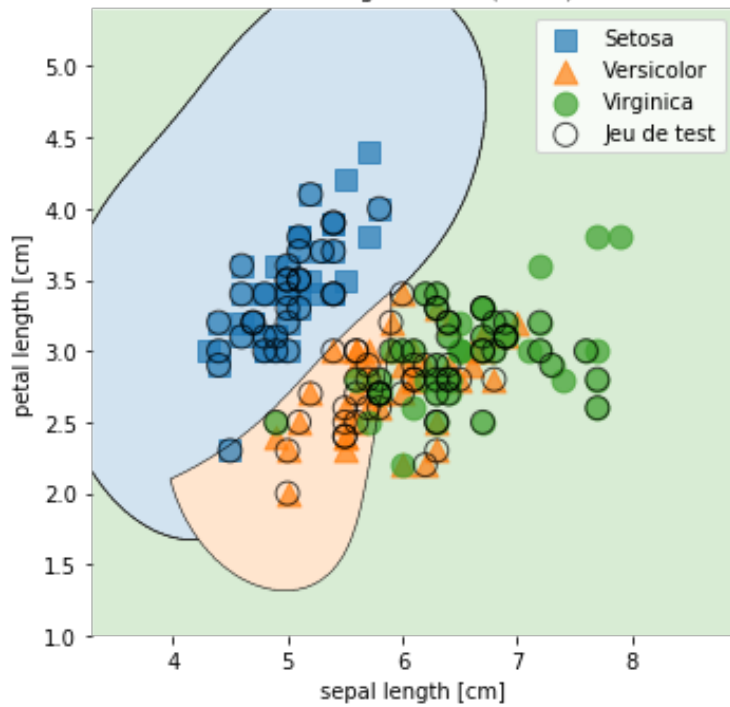
```
1 C = 1.0  
2  
3 # SVC avec des valeurs différentes de gamma  
4 rbf_svc1 = SVC(kernel='rbf', gamma=1, C=C)  
5 rbf_svc2 = SVC(kernel='rbf', gamma=10, C=C)  
6 rbf_svc3 = SVC(kernel='rbf', gamma=50, C=C)  
7 rbf_svc4 = SVC(kernel='rbf', gamma=100, C=C)
```



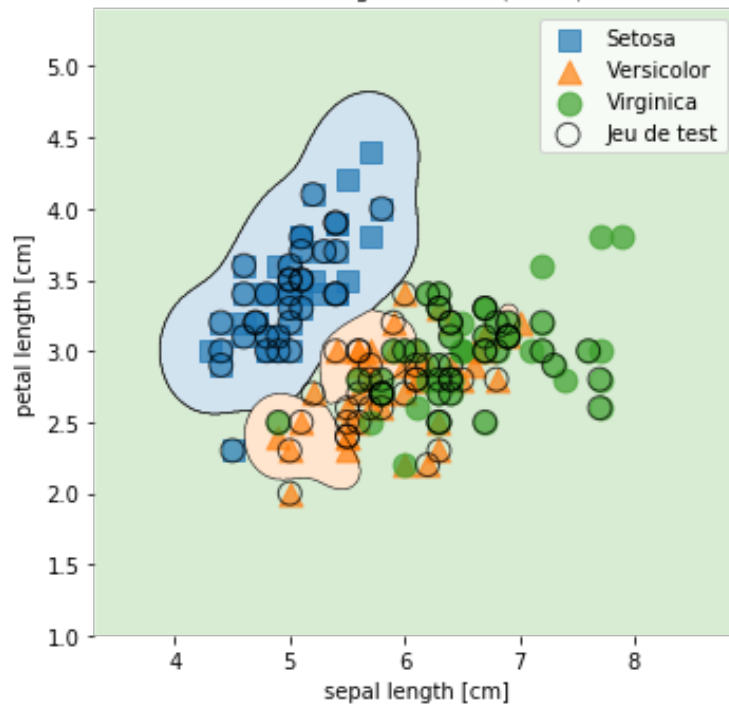
In [12]:

```
1 labels = ['SVC RBF gamma=1',  
2           'SVC RBF gamma=10',  
3           'SVC RBF gamma=50',  
4           'SVC RBF gamma=100']  
5 list_clf=[rbf_svc1,rbf_svc2,rbf_svc3,rbf_svc4]  
6 call_decision_regions_iris (labels,list_clf,X_train,y_train,X_test,y_test)
```

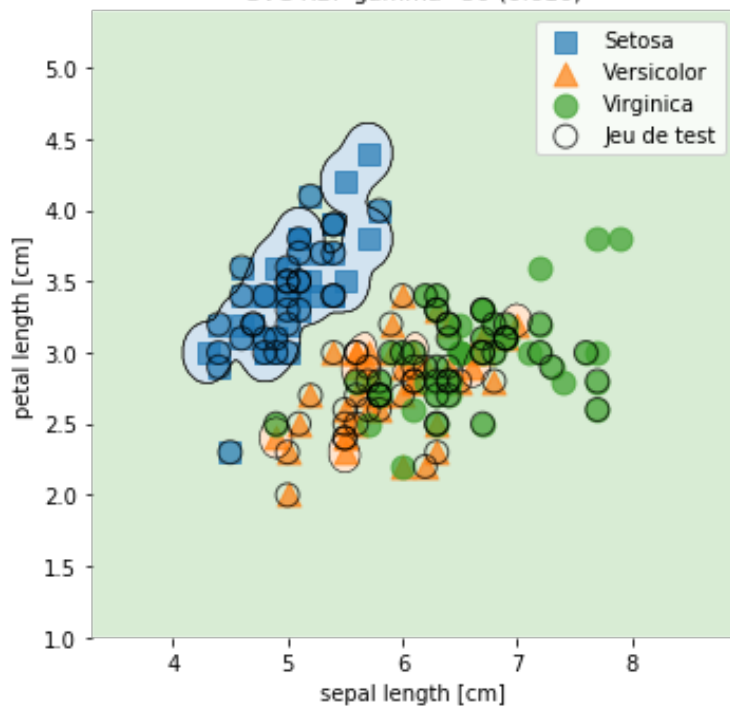
SVC RBF gamma=1 (0.771)



SVC RBF gamma=10 (0.676)



SVC RBF gamma=50 (0.610)



SVC RBF gamma=100 (0.562)

