



Mini Projets 2019 (Info 232)

Isabelle Guyon

info232@chalearn.org

1. Visualisation

We are visualizing the famous [iris dataset \(https://archive.ics.uci.edu/ml/datasets/iris\)](https://archive.ics.uci.edu/ml/datasets/iris).

Instructions

Ce TP vaut 5 points. Répondez à toutes les questions.

- Pour créer une nouvelle cellule, allez dans le menu "Insert".
- Pour transformer une cellule en commentaire texte, allez dans Cell + Cell Type + Markdown.
- Pour exécuter une cellule: SHIFT+RETURN

Les cellules doivent être exécutées dans l'ordre. Pour plus d'information [CONSULTEZ la DOCUMENTATION.](#)

Question 0: "Markdown" cells

Créez une nouvelle cellule de type Markdown en dessous de celle-ci. Recopiez dedans le paragraphe sur k-fold cross validation que vous trouverez sur Wikipedia. Essayez de colorer la cellule en VERT!

Question 1: "Code" cells

Maintenant vous allez exécuter la cellule ci-dessous après avoir remplacé la réponse par 1.

In [1]:

```
1 code_dir = 'code/'
2 from sys import path; path.append(code_dir)
3 %matplotlib inline
4 %load_ext autoreload
5 %autoreload 2
6 from checker import check
7 import warnings
8 warnings.simplefilter(action='ignore', category=FutureWarning)
9 question = 1
10 answer = 1 # Replace by 1
11 score = 0
12 score += check(answer, question)
```

1
1

CORRECT

:-)

Question 2: AutoML format

Les donnees que vous aurez a analyser dans votre projet seront au [format AutoML](https://github.com/codalab/chalab/wiki/Help:-Wizard-%E2%80%90-Challenge-%E2%80%90-Data) (<https://github.com/codalab/chalab/wiki/Help:-Wizard-%E2%80%90-Challenge-%E2%80%90-Data>): . En utilisant les moyens que vous voulez, remplacez les valeurs des variables de dimension des donnees par leur valeurs correctes dans la deuxieme cellule avant de l'executer.

In [2]:

```
1 data_dir = 'data'
2 data_name = 'iris'
3 !ls $data_dir*
4 !cat $data_dir/iris_public.info
```

```
iris_feat.name      iris_public.info    iris_train.data    iris_vali
d.solution
iris_label.name     iris_test.data      iris_train.solution
iris_private.info   iris_test.solution  iris_valid.data
usage = 'Sample dataset Iris data'
name = 'iris'
task = 'multiclass.classification'
target_type = 'Numerical'
feat_type = 'Numerical'
metric = 'bac_metric'
time_budget = 1200
feat_num = 4
target_num = 3
label_num = 3
train_num = 35
valid_num = 35
test_num = 35
has_categorical = 0
has_missing = 0
is_sparse = 0
```

In [3]:

```
1 feature_number = 4
2 training_sample_number = 35
3 validation_sample_number = 35
4 test_sample_number = 35
5 question = 2
6 reponse = feature_number*(training_sample_number+validation_sample_number+test_s
7 score += check(reponse, question)
```

420
420

CORRECT
:-)

Question 3: Pandas

En Anglais pour changer :-)

This time we are going to do simple "exploratory data analysis". To simplify, we lump all the data together in one big data structure called a "pandas" data frame. This will allow us to use the rich libraries "pandas" and "seaborn" to explore the data.

In the next cell, replace the "head" function, which just shows the first few rows of the dataset, by a pandas function providing descriptive statistics. To that end, you may want to check the [Pandas DataFrame reference page \(https://pandas.pydata.org/pandas-docs/stable/reference/frame.html\)](https://pandas.pydata.org/pandas-docs/stable/reference/frame.html). Then, in the following cell, replace the variables with their correct values and execute it.

In [4]:

```
1 from data_io import read_as_df
2 data = read_as_df(data_dir + '/' + data_name) # The data are located in data_dir
3 data.head()
4 data.describe()
```

Reading data/iris_train from AutoML format

Number of examples = 35

Number of features = 4

Class

0 setosa

1 versicolor

2 virginica

Number of classes = 3

Out[4]:

	sepal_length	sepal_width	petal_length	petal_width
count	35.000000	35.000000	35.000000	35.000000
mean	5.625714	3.005714	3.402857	1.054286
std	0.892565	0.421442	1.962514	0.839588
min	4.300000	2.000000	1.000000	0.100000
25%	4.900000	2.800000	1.450000	0.200000
50%	5.500000	3.000000	3.700000	1.000000
75%	6.300000	3.250000	5.100000	1.800000
max	7.700000	4.000000	6.700000	2.500000

In [5]:

```

1 std_sepal_length = 0.892565 # Standard deviation of the sepal length
2 mean_sepal_width = 3.005714 # Mean of the sepal width
3 min_petal_length = 1 # Minimum value of the petal length
4 max_petal_width = 2.5 # Maximum value of the petal width
5 question = 3
6 reponse = std_sepal_length+mean_sepal_width+min_petal_length+max_petal_width
7 score += check(reponse, question)

```

7.3982790000000005

7.3982790000000005

CORRECT

:-)

Question 4: Histograms

Un truc sympa de Pandas c'est que ça permet aussi de faire des graphes pour visualiser les données.

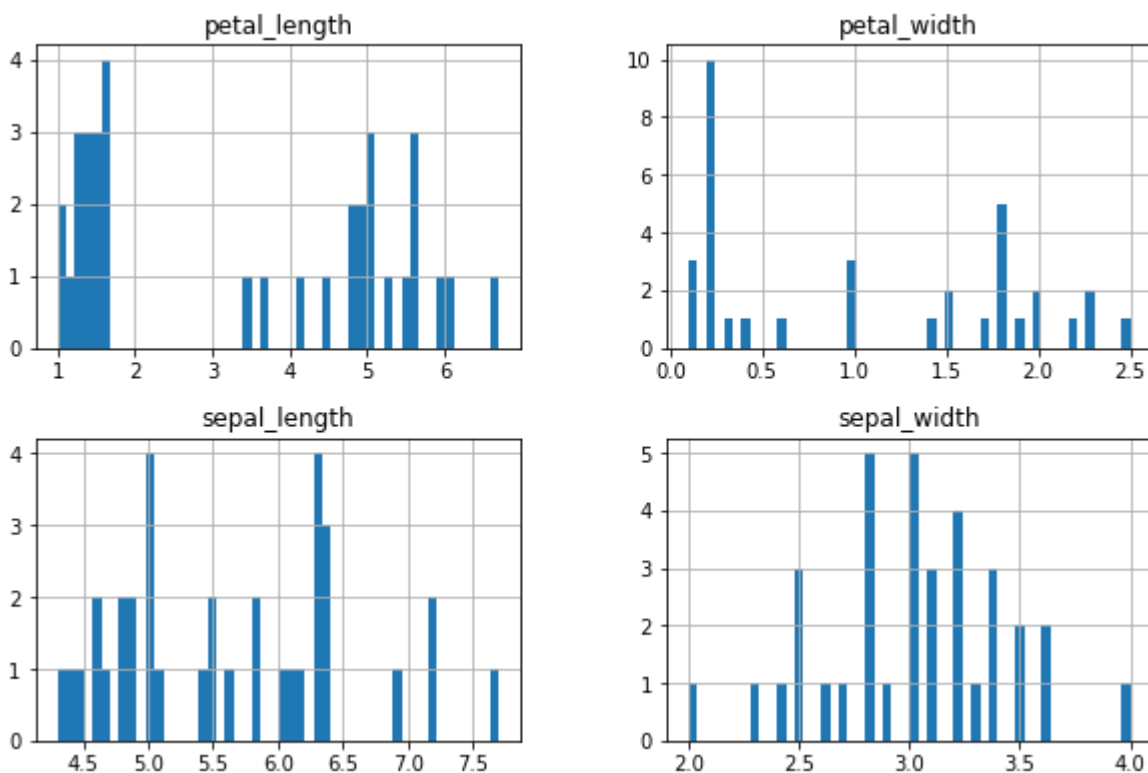
Remplacez le nombre de bins par un plus petit nombre. Répondez en changeant la variable answer dans la deuxième cellule: answer=1 si la hauteur de la plus grande barre diminue et answer=0 sinon. Comprenez-vous pourquoi?

In [6]:

```

1 data.hist(figsize=(10, 10), bins=50, layout=(3, 2));

```



In [7]:

```

1 question = 4
2 answer = 1          # one if the maximum bar height increases when the bin num
3 score += check(answer, question)

```

1
1

CORRECT

:-)

Have fun with the [Pandas DataFrame reference page \(https://pandas.pydata.org/pandas-docs/stable/reference/frame.html\)](https://pandas.pydata.org/pandas-docs/stable/reference/frame.html). There are lots of other plotting functions. Try some of them in the cell below.

In [8]:

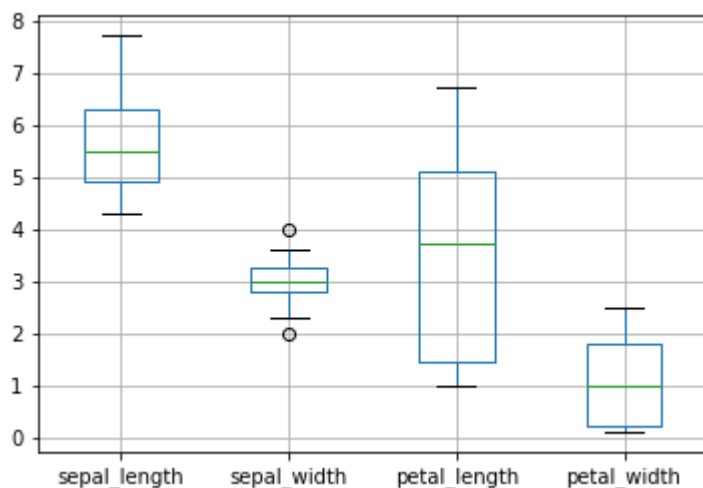
```

1 # What's this one for instance? Change it to something else.
2 data.boxplot()

```

Out[8]:

<matplotlib.axes._subplots.AxesSubplot at 0x1a1d303780>



Question 5: Pair plots

Seaborn (sns) is a package of data visualization functions: <https://seaborn.pydata.org/> (<https://seaborn.pydata.org/>). Quite useful! It is convenient to visualize data in 2 dimensions. One way of doing that is to plot a variable (feature) against another one, one point representing a sample (a flower). The pairplot function shows all the possibilities (off-diagonal graphs).

On the diagonal, what do you see? Compare with the histograms of the previous question. Then add another argument to the pairplot function

```
hue="target"
```

(if you do not understand, consult the documentation

<https://seaborn.pydata.org/generated/seaborn.pairplot.html>).

(<https://seaborn.pydata.org/generated/seaborn.pairplot.html>). After executing the next cell again, in the following cell, answer the questions:

What is the color of the class, which is best separated from all others?

color_best_separated = 1 if blue; 2 if orange; 3 if green.

Which iris type does this correspond to?

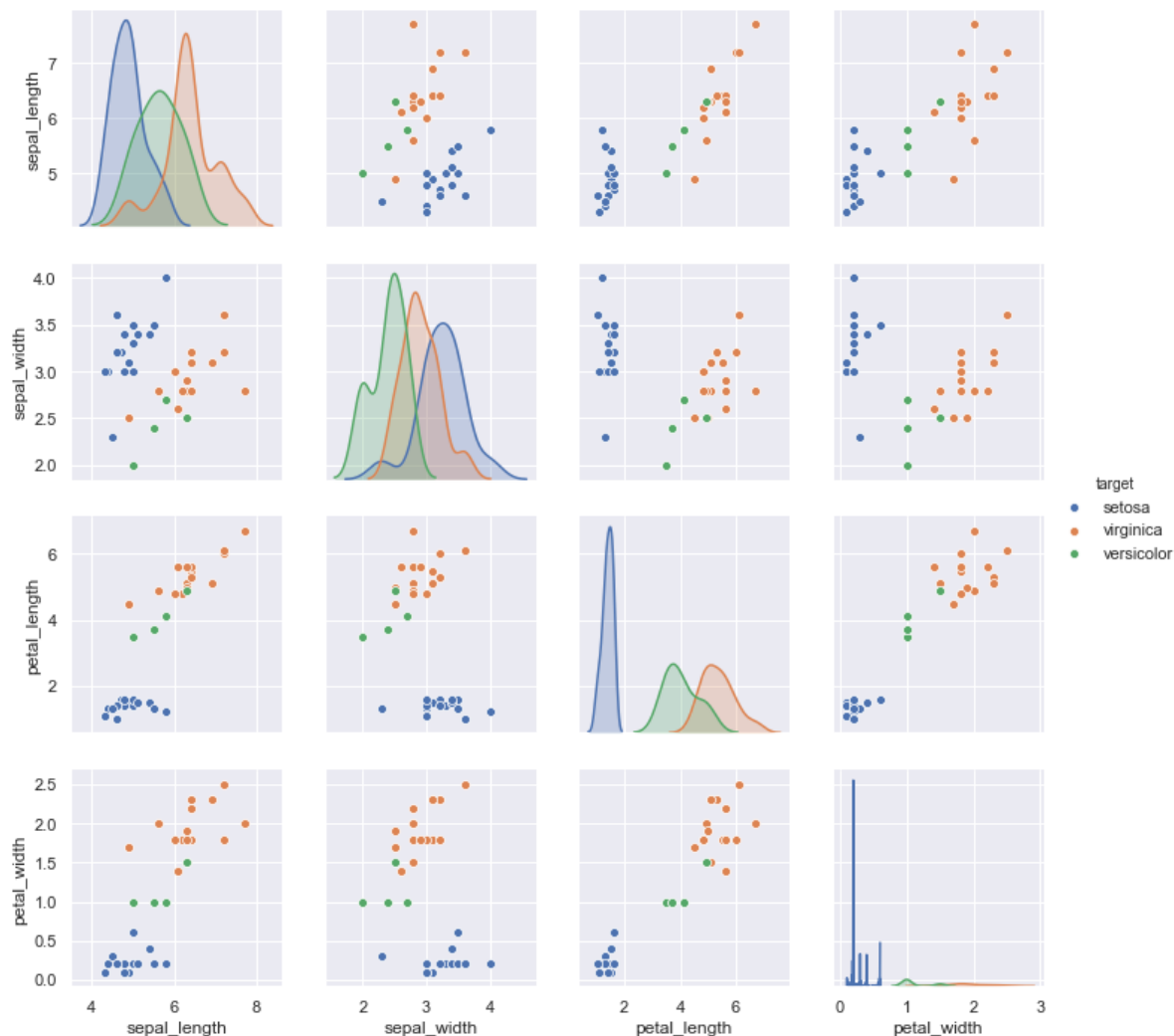
iris_best_separated = 1 if virginica; 2 if versicolor; 3 if setosa.

In [9]:

```
1 import seaborn as sns; sns.set()
2 sns.pairplot(data, hue="target"),
```

Out[9]:

(<seaborn.axisgrid.PairGrid at 0x1a1d7054a8>,,)



In [10]:

```
1 question = 5
2 color_best_separated = 1
3 iris_best_separated = 3
4 score += check(color_best_separated*iris_best_separated, question)
```

3
3

CORRECT

:~)

Question 6: Feature correlation

Les variables (features) peuvent être redondantes (c'est à dire capturer des informations similaires). Le coefficient de corrélation de Pearson (voir [page Wikipedia \(https://en.wikipedia.org/wiki/Pearson_correlation_coefficient\)](https://en.wikipedia.org/wiki/Pearson_correlation_coefficient)) permet de détecter la corrélation (c'est à dire la similarité au sens d'une dépendance linéaire).

En regardant les "pair plots" ci-dessus, à votre avis, quelle paire de variable est la plus corrélée? Repérez la paire par son numéro de ligne et de colonne pour répondre. Vérifiez votre intuition en exécutant la cellule suivante qui représente graphiquement la matrice de corrélation. Changez la méthode 'pearson' pour d'autres coefficients de corrélation (voir [documentation \(https://pandas.pydata.org/pandas-docs/stable/reference/api/pandas.DataFrame.corr.html\)](https://pandas.pydata.org/pandas-docs/stable/reference/api/pandas.DataFrame.corr.html)). Est-ce que ça change quelle paire est la plus corrélée? Regardez les définitions de [Kendall rank correlation coefficient \(Kendall tau\)](https://en.wikipedia.org/wiki/Kendall_rank_correlation_coefficient) et [Spearman correlation coefficient \(https://en.wikipedia.org/wiki/Spearman%27s_rank_correlation_coefficient\)](https://en.wikipedia.org/wiki/Spearman%27s_rank_correlation_coefficient). Essayez de comprendre la différence entre [corrélation et dépendance \(https://en.wikipedia.org/wiki/Correlation_and_dependence\)](https://en.wikipedia.org/wiki/Correlation_and_dependence).

In [11]:

```
1 question = 6
2 numero_ligne = 3          # Lignes numérotée de 0 à 3
3 numero_colonne = 2        # Colonnes numérotée de 0 à 3
4 score += check(numero_ligne*numero_colonne, question)
```

6
6

CORRECT

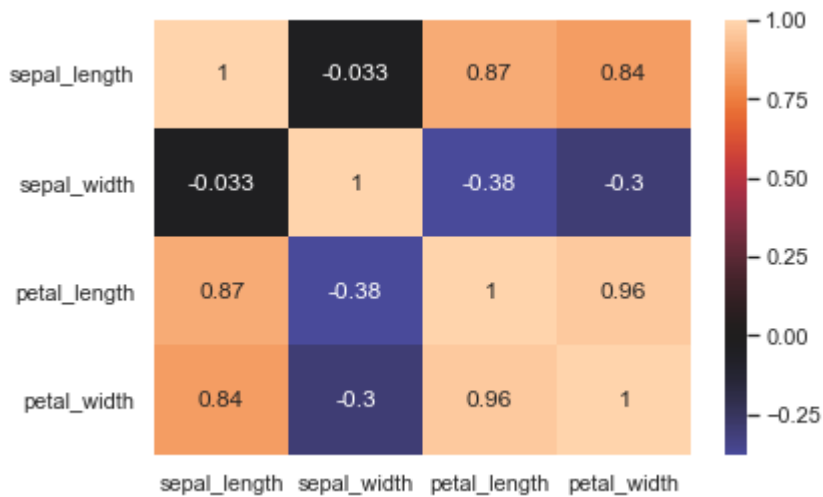
:~)

In [12]:

```
1 corr_mat = data.corr(method='pearson')
2 sns.heatmap(corr_mat, annot=True, center=0)
```

Out[12]:

<matplotlib.axes._subplots.AxesSubplot at 0x1a1d14f710>



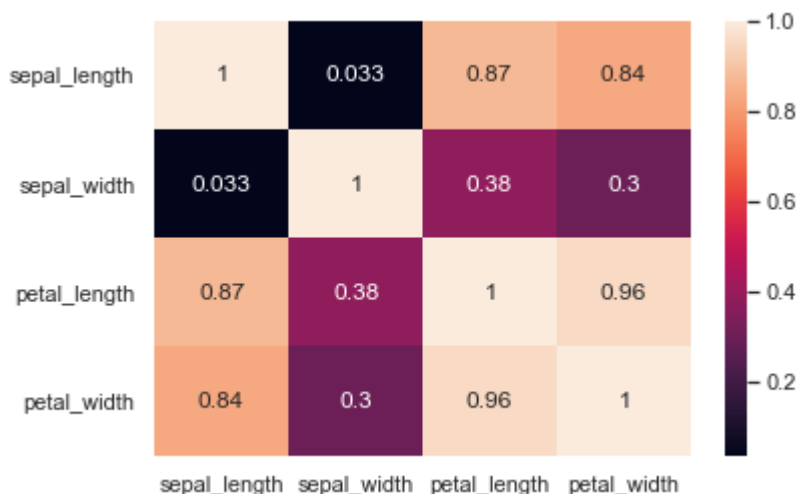
Remarquez que les variables peuvent être corrélées ou anti-corrélées. Vous voulez donc peut-être vous intéresser à la valeur absolue du coefficient de corrélation. Cela change-t-il votre réponse? Peut-être pas, mais ça pourrait!

In [13]:

```
1 sns.heatmap(abs(corr_mat), annot=True)
```

Out[13]:

<matplotlib.axes._subplots.AxesSubplot at 0x1a1d0d5be0>



Notez que la matrice est symétrique. Si vous voulez vous amuser, effacez les valeurs au dessus de la diagonale ([voir en bas de cette page \(https://seaborn.pydata.org/generated/seaborn.heatmap.html\)](https://seaborn.pydata.org/generated/seaborn.heatmap.html)).

Question 7: Feature selection

Representing a matrix of coefficients with colors seems to be pretty convenient for visualization purposes. We would like to do that also for the data matrix itself. Note that, since the last column (target) contains strings ("categorical variables"), we first need to convert them to numbers.

Observing the heatmap, which column is most correlated with the target? Insert another cell in which you plot the correlation matrix of data_new (inspiring yourself from the previous question), then confirm your intuition and answer the question.

In [14]:

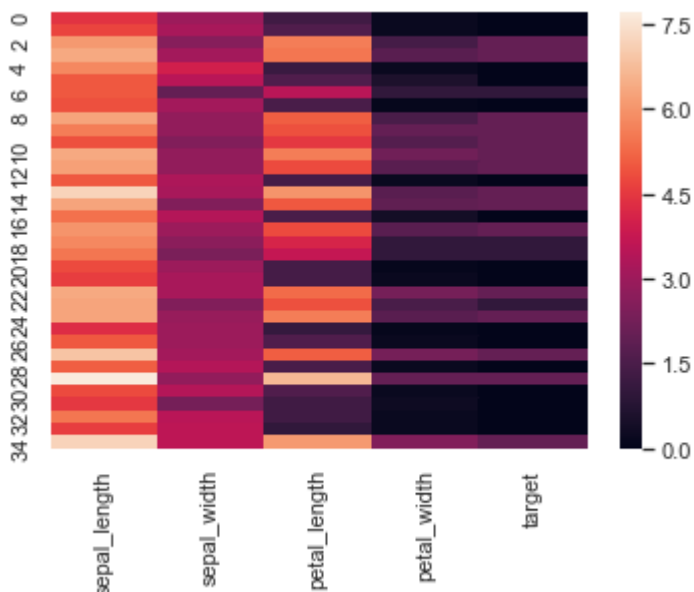
```
1 print(data.head())
2 data_num = data.copy() # If you don't use "copy", any change in data_num will affect data
3 data_num['target'] = data_num['target'].astype('category')
4 data_num['target'] = data_num['target'].cat.codes
5 print(data_num.head())
6 sns.heatmap(data_num)
```

	sepal_length	sepal_width	petal_length	petal_width	target
0	4.4	3.0	1.3	0.2	setosa
1	4.7	3.2	1.6	0.2	setosa
2	6.1	2.6	5.6	1.4	virginica
3	6.4	3.1	5.5	1.8	virginica
4	5.8	4.0	1.2	0.2	setosa

	sepal_length	sepal_width	petal_length	petal_width	target
0	4.4	3.0	1.3	0.2	0
1	4.7	3.2	1.6	0.2	0
2	6.1	2.6	5.6	1.4	2
3	6.4	3.1	5.5	1.8	2
4	5.8	4.0	1.2	0.2	0

Out[14]:

<matplotlib.axes._subplots.AxesSubplot at 0x1a1e8966d8>



In [15]:

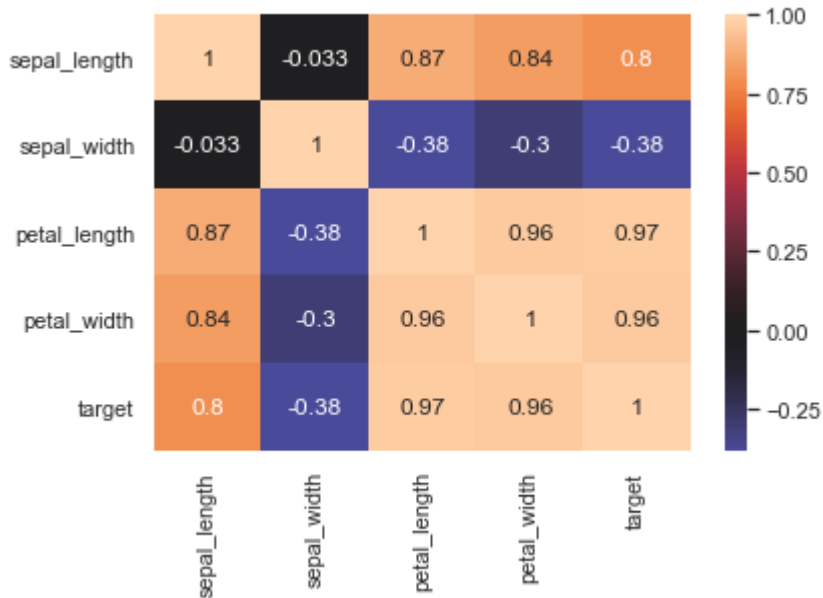
```

1 # Mettez ici votre code montrant la matrice de corrélation de data_num
2 corr_mat = data_num.corr(method='pearson')
3 sns.heatmap(corr_mat, annot=True, center=0)

```

Out[15]:

<matplotlib.axes._subplots.AxesSubplot at 0x1a1e896278>



La question 7 est ici:

Quelle est la variable (feature) la plus corrélée avec la colonne "target"? Quelle est la valeur du coefficient de corrélation de Pearson correspondant?

In [16]:

```

1 question = 7
2 numero_variable = 2 # Variables numérotées de 0 à 3
3 pearson_correlation = 0.97
4 score += check(numero_variable+pearson_correlation, question)

```

2.9699999999999998

2.97

CORRECT

:-)

In []:

1

Question 8: One-Rule classifier

In this section, we show how we can create a very simple classifier based on just ONE rule to separate the 3 types of flowers. That rule classifies irises on the basis of their petal length only. Check the code to see whether you understand it.

This classifier respects the structure of [scikit-learn \(https://scikit-learn.org/stable/\)](https://scikit-learn.org/stable/) learning machines. To make it compatible with other scikit-learn tools, we derive it from the base class BaseEstimator and overload 2 methods: "fit" and "predict". Then we use the Iris data to train and test a model (in this case we lumped all the data into a single matrix and use it as training data). After that we compute the training error.

Scikit-learn allows you to compute [a lot of other metrics \(https://scikit-learn.org/stable/modules/classes.html#sklearn-metrics-metrics\)](https://scikit-learn.org/stable/modules/classes.html#sklearn-metrics-metrics). To answer this question, you will have to compute the [BAC \(https://scikit-learn.org/stable/modules/generated/sklearn.metrics.balanced_accuracy_score.html#sklearn.metrics.balanced_accuracy_score\)](https://scikit-learn.org/stable/modules/generated/sklearn.metrics.balanced_accuracy_score.html#sklearn.metrics.balanced_accuracy_score) that is the Balanced ACcuracy score.

In [17]:

```

1  import numpy as np
2  from sklearn.base import BaseEstimator
3
4  class oneR(BaseEstimator):
5      ''' One Rule classifier '''
6      def __init__(self):
7          ''' The "constructor" initializes the parameters '''
8          self.selected_feat = 0 # The chosen variable/feature
9          self.theta1 = 0       # The first threshold
10         self.theta2 = 0        # The second threshold
11
12     def fit(self, X, Y, F=[]):
13         ''' The method "fit" trains a super-simple classifier '''
14         if not F: F=[str(item) for item in range(X.shape[1])]
15         # First it selects the feature most correlated to the target
16         correlations = np.corrcoef(X, Y, rowvar=0)
17         self.selected_feat = np.argmax(correlations[0:-1, -1])
18         best_feat = X[:, self.selected_feat]
19         print('Feature selected = ' + F[self.selected_feat])
20         # Then it computes the average values of the 3 classes
21         mu0 = np.median(best_feat[Y==0])
22         mu1 = np.median(best_feat[Y==1])
23         mu2 = np.median(best_feat[Y==2])
24         # Finally it sets two decision thresholds
25         self.theta1 = (mu0+mu1)/2.
26         self.theta2 = (mu1+mu2)/2.
27
28     def predict(self, X):
29         ''' The method "predict" classifies new test examples '''
30         # Select the values of the correct feature
31         best_feat = X[:, self.selected_feat]
32         # Initialize an array to hold the predicted values
33         Yhat = np.copy(best_feat) # By copying best_fit we get an array
34         # then classify using the selected feature according to the cutoff thresholds
35         Yhat[best_feat<self.theta1] = 0
36         Yhat[np.all([self.theta1<=best_feat, best_feat<=self.theta2], 0)] = 1
37         Yhat[best_feat>self.theta2] = 2
38         return Yhat

```

In [18]:

```

1  XY=data_num.as_matrix() # On transforme le Pandas dataframe en un NumPy array
2  X = XY[:,0:4]            # On recupere X (num_exemples x num_features)
3  Y = XY[:,4]             # et Y (num_exemples x 1) ==> les valeurs de la classe
4  feature_names = list(data_num)[-1] # On recupere aussi les noms des features
5  my_model = oneR()
6  my_model.fit(X, Y, feature_names) # Le nom des features est "optionnel", il peut être utile
7  Yhat = my_model.predict(X)

```

Feature selected = petal_length

In [19]:

```

1 def error_rate(solution, prediction):
2     return np.mean(solution!=prediction)
3
4 errate = error_rate(Y, Yhat)
5 print('Training error = %5.2f' % errate)
6 Yperm = np.random.permutation(Y)
7 print('Random permutation error (for comparison)= %5.2f' % error_rate(Y, Yperm))
8 print('Ideal error rate (for comparison)= %5.2f' % error_rate(Y, Y))

```

```

Training error = 0.06
Random permutation error (for comparison)= 0.63
Ideal error rate (for comparison)= 0.00

```

La question 8 est ici:

Remplacez xxx par la bonne fonction qui calcule le [BAC \(https://scikit-learn.org/stable/modules/generated/sklearn.metrics.balanced_accuracy_score.html#sklearn.metrics.balanced_accuracy_score\)](https://scikit-learn.org/stable/modules/generated/sklearn.metrics.balanced_accuracy_score.html#sklearn.metrics.balanced_accuracy_score), i.e. le Balanced ACcuracy score, et enlevez les # de commentaire. Observez que error_rate est different de Balanced Error Rate BER=1-BAC.

In [20]:

```

1 BER = errate
2 #from sklearn.metrics import xxx
3 #BAC = xxx(Y, Yhat)
4 #print(' BAC = %5.2f' % BAC)
5 #BER = 1 - BAC
6 from sklearn.metrics import balanced_accuracy_score
7 BAC = balanced_accuracy_score(Y, Yhat)
8 print('BAC = %5.2f' % BAC)
9 BER = 1 - BAC
10 print('BER = %5.2f' % BER)
11 question = 8
12 score += check(BER, question)

```

```

BAC = 0.89
BER = 0.11
0.10555555555555551
0.10555555555555551

```

CORRECT

:-)

Question 9: Viewing the classification results

Pour se faire une idee de comment marche les classifieurs il est utile de visualiser les regions de decision. Ce petit programme permet de le faire en deux dimensions.

In [21]:

```

1  import numpy as np
2  import matplotlib.pyplot as plt
3  from matplotlib.colors import LinearSegmentedColormap
4  colors = [(1, 0, 0), (0, 1, 0), (0, 0, 1)] # Red, lime, blue
5  cm = LinearSegmentedColormap.from_list('rgb', colors, N=3)
6
7  def ClfScatter(clf, X, Y, F, dim1=0, dim2=1, title=''):
8      '''clf_scatter(clf, X, Y, F, dim1=0, dim2=1)
9      Display decision function and training examples.
10     clf: a classifier with at least a fit and a predict method
11     like a scikit-learn classifier.
12     X: a 2 dimensional data matrix, samples in line and features in columns
13     Y: a target vectors of class values 0, 1, 2, etc.
14     F: feature names
15     dim1 and dim2: chosen features.
16     title: Figure title.
17     Returns: Predictions on training examples.
18     '''
19     # Fit model in chosen dimensions
20     X2 = X[:, (dim1, dim2)]
21     try:
22         clf.fit(X2, Y, F=[F[dim1], F[dim2]])
23     except:
24         clf.fit(X2, Y)
25     # Define a mesh
26     x_min, x_max = X2[:, 0].min() - 1, X2[:, 0].max() + 1
27     y_min, y_max = X2[:, 1].min() - 1, X2[:, 1].max() + 1
28     h = 0.1 # step
29     xx, yy = np.meshgrid(np.arange(x_min, x_max, h),
30                          np.arange(y_min, y_max, h))
31     Xtest = np.c_[xx.ravel(), yy.ravel()]
32     # Compute the training error
33     Yhat = clf.predict(X2)
34     training_error = error_rate(Y, Yhat)
35     # Make your predictions on all mesh grid points (test points)
36     Yhat = clf.predict(Xtest)
37     # Make contour plot for all points in mesh
38     Yhat = Yhat.reshape(xx.shape)
39     plt.contourf(xx, yy, Yhat, cmap=plt.cm.Paired)
40     # Overlay scatter plot of training examples
41     plt.scatter(X2[:, 0], X2[:, 1], c=Y, cmap=cm)
42     plt.title('{}: training error = {:.2f}'.format(title, training_error))
43     plt.xlabel(F[dim1])
44     plt.ylabel(F[dim2])
45     plt.show()
46     return clf.predict(X2)

```

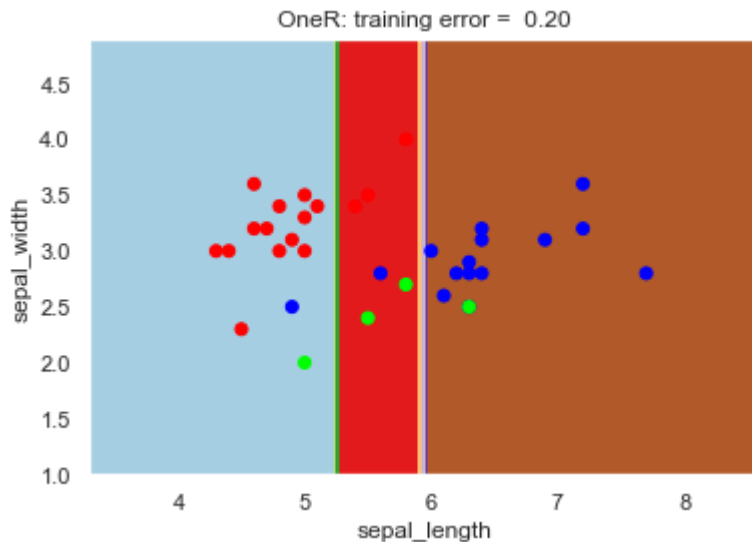
In [22]:

```

1 Yhat = ClfScatter(my_model, X, Y, feature_names, dim1=0, dim2=1, title='OneR')
2 errate = error_rate(Y, Yhat)
3 print('Training error = %5.2f' % errate)

```

Feature selected = sepal_length



Training error = 0.20

La question 9 est ici:

Remplacez `dim1=2, dim2=3` par `dim1=0, dim2=1` dans la cellule du dessus et re-exécutez la. Vous devriez obtenir de moins bonnes performances. Pensez vous qu'un classifieur qui utiliserait plusieurs règles pourrait obtenir 0 erreurs?

In [23]:

```

1 question = 9
2 reponse = 1 # Oui = 1, Non = 0
3 score += check(errate+reponse, question)

```

1.2

1.2

CORRECT

:-)

Question 10: More classifiers

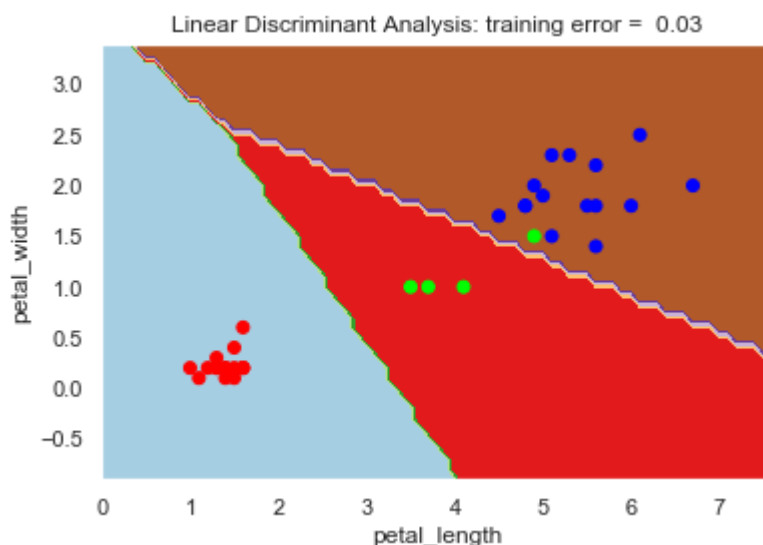
En utilisant le code précédent, on peut comparer plein de classifieurs. Certains obtiennent les mêmes résultats si on les fait tourner plusieurs fois, d'autres donnent toujours la même chose. Combien d'entre eux (parmi les exemples ci-dessous) donnent toujours la même chose? Vérifiez bien dans la doc de [scikit-learn](https://scikit-learn.org/stable/) (<https://scikit-learn.org/stable/>).

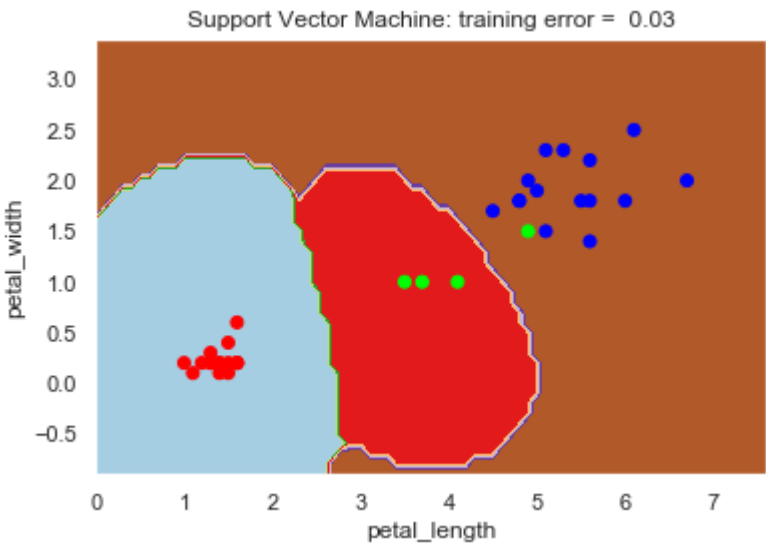
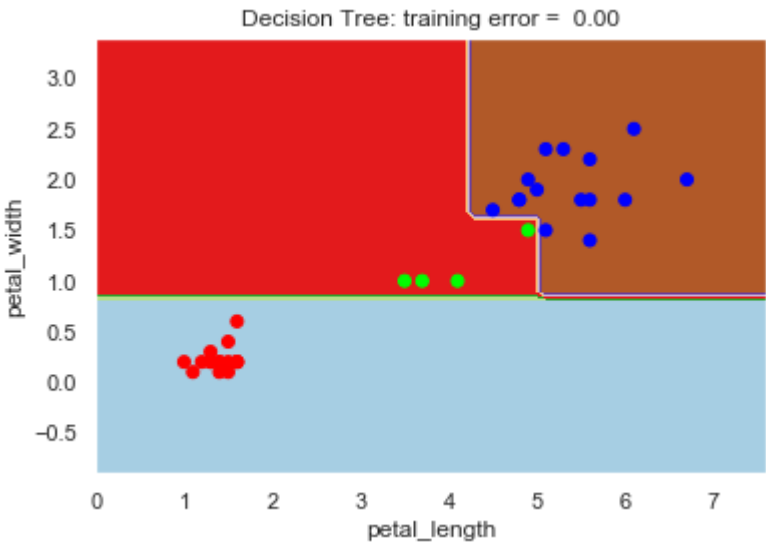
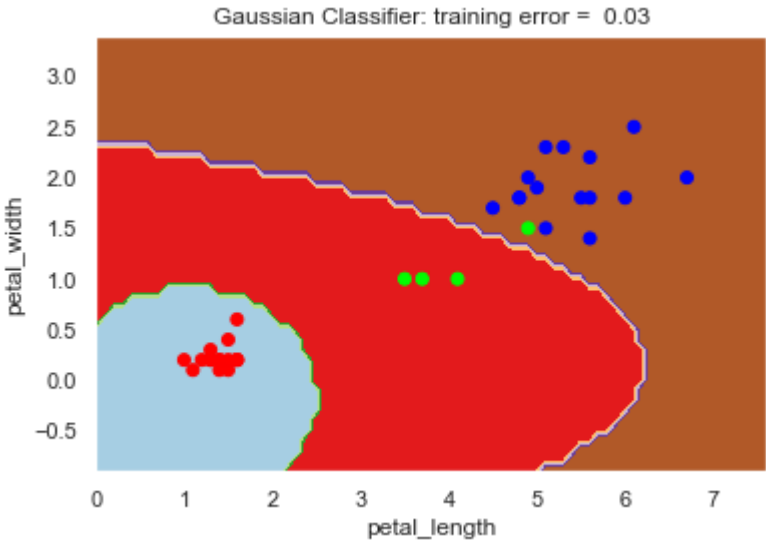
In [24]:

```

1 import random
2 from sklearn.linear_model import Perceptron
3 from sklearn.discriminant_analysis import LinearDiscriminantAnalysis
4 from sklearn.naive_bayes import GaussianNB
5 from sklearn.tree import DecisionTreeClassifier
6 from sklearn.svm import SVC
7 classifiers = [
8     Perceptron(random_state=random.randint(1,101)),
9     LinearDiscriminantAnalysis(),
10    GaussianNB(),
11    DecisionTreeClassifier(random_state=random.randint(1,101)),
12    SVC()]
13 names = ['Perceptron',
14          'Linear Discriminant Analysis',
15          'Gaussian Classifier',
16          'Decision Tree',
17          'Support Vector Machine']
18 for clf, name in zip(classifiers, names):
19     # This does a two dimensional fit in dim1 and dim2
20     ClfScatter(clf, X, Y, feature_names, dim1=2, dim2=3, title=name)

```





In [25]:

```
1 question = 10
2 reponse = 3 # Nombre de classifieurs qui retournent toujours la même separation
3 score += check(reponse, question)
```

3

3

CORRECT

:-)

In [26]:

```
1 print('Your final score is %d / 10, congratulations!' % score)
```

Your final score is 10 / 10, congratulations!

In []:

1