



## Introduction to Machine Learning with Scikit-learn

### Part II: Simple classification models

For this second part of an introduction to the `scikit-learn` module, we will focus on the second type of problem in Machine Learning: the **classification** problem.

The objective of this introduction is to:

- Introduce the classification problem.
- Learn to use the `scikit-learn` module to build a classification model, also known as a "classifier".
- Introduce metrics adapted to the evaluation of classification models.

### Introduction to classification

#### Objective of classification

In supervised learning, the objective is to predict the value of a target variable from explanatory variables.

- In a **regression** problem, the target variable takes **continuous values**. These values are numerical: price of a house, quantity of oxygen in the air of a city, etc. The target variable can therefore take an **infinity of values**.
- In a **classification** problem, the target variable takes **discrete values**. These values can be numeric or literal, but in both cases the target variable takes a **finite** number of values.

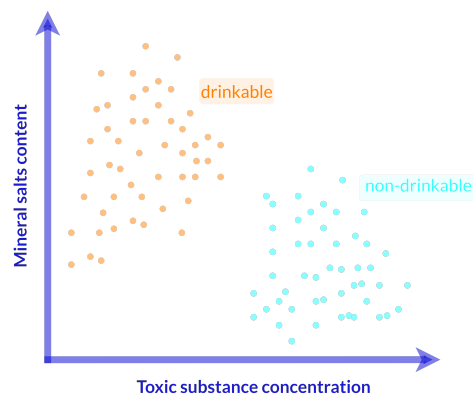
The different values taken by the target variable are called **classes**.

The objective of classification therefore consists in predicting the class of an observation from its explanatory variables.

#### An example of classification

We will look at a problem of a **binary** classification, i.e. where there are **two** classes. We are trying to determine whether the water in a stream is drinkable or not depending on its concentration of toxic substances and its mineral salts content.

The two classes are therefore 'drinkable' and 'non-drinkable'.



In the figure above, each point represents a stream whose position on the map is defined by its values for the concentration of toxic substances and the content of mineral salts.

The objective will be to build a **model capable of assigning one of the two classes** ('drinkable' / 'non-drinkable') to a stream of which only these two variables are known.

The figure above suggests the existence of two zones allowing easy classification of streams:

- An area where the streams are drinkable (top left).
- An area where the streams are not drinkable (bottom right).

We would like to create a model capable of **separating the dataset into two parts** corresponding to these areas.

A simple technique would be to separate the two areas **using a line**.

- (a) Run the next cell to display the interactive figure.

- The **orange** dots are the **drinkable** streams and the **blue** dots are the **non-drinkable** streams.
- The **red arrow** corresponds to a **vector** defined by  $w = (w_1, w_2)$ . The red line corresponds to the orthogonal (i.e. perpendicular) plane to  $w$ . You can change the coordinates of the vector  $w$  in two ways:

- By moving the sliders of  $w_1$  and  $w_2$ .

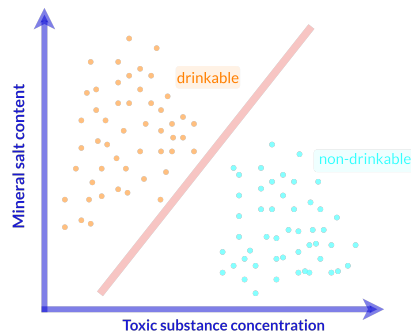
In [ ]:

```
from classification_widgets import linear_classification
linear_classification()
```

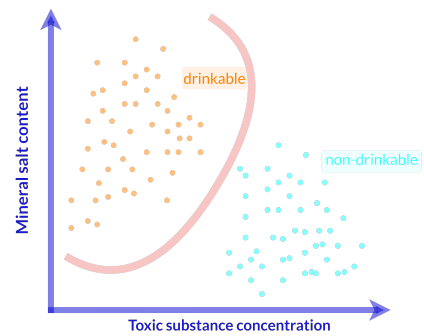
The classification we just performed is of **linear** type, that is to say that we used a flat linear plane to separate our classes.

This plane was parametrized by the vector  $w$ . Thus, **the objective of linear classification models is to find the vector  $w$  allowing the best possible separation of the different classes**. Each model of linear type has its own technique to find this vector.

There are also **non-linear** classification models, which we will see later.



Linear Classification



Non-Linear Classification

## 1. Using scikit-learn for classification

We will now introduce the main tools of the `scikit-learn` module for solving a classification problem.

In this exercise we will use the [Congressional Voting Records](https://archive.ics.uci.edu/ml/datasets/congressional+voting+records) (<https://archive.ics.uci.edu/ml/datasets/congressional+voting+records>) dataset, containing a number of votes cast by members of Congress of the United States House of Representatives.

The objective of our classification problem will be to **predict the political party** ("Democrat" or "Republican") of the members of the House of Representatives according to their votes on subjects like the education, health, budget, etc.

The explanatory variables will therefore be the votes on various subjects and the target variable will be the "democrat" or "republican" political party.

To solve this problem we will use:

- A non-linear classification model: **K-Nearest Neighbors**.
- A linear classification model: **Logistic Regression**.

### Data preparation

- (a) Run the following cell to import the `pandas` and `numpy` modules needed for the exercise.

In [ ]:

```
import pandas as pd
import numpy as np
%matplotlib inline
```

- (b) Load the data contained in the file `'votes.csv'` into a `DataFrame` named `votes`.

In [ ]:

```
# Insert your code here
```

Hide solution

In [ ]:

```
votes = pd.read_csv('votes.csv')
```

In order to briefly visualize our data:

- (c) Display the number of rows and columns of `votes`.
- (d) Show a preview of the first 20 rows of `votes`.

In [ ]:

```
# Insert your code here
```

Hide solution

In [ ]:

```
# Shape of the DataFrame
print('The DataFrame has', votes.shape[0], 'rows and', votes.shape[1], 'columns.')

# Display the first 20 rows
votes.head(20)
```

- The first column **"party"** contains the name of the **political party** to which each member of the Congress of the House of Representatives belongs. This is the target variable.
- The following **16** columns contain the votes of each member of Congress on legislative proposals:
  - 'y' indicates that the elected member voted **for** the bill.
  - 'n' indicates that the elected member voted **against** the bill.

In order to use the data in a classification model, it is necessary to transform these columns into binary **numeric** values, i.e. either 0 or 1.

- (e) For each of the columns 1 to 16 (column 0 being our target variable), replace the values 'y' by 1 and 'n' by 0. To do so, we can use the **replace** method from the **DataFrame** class.
- (f) Display the first 10 rows of the modified **DataFrame**.

In [ ]:

```
# Insert your code here
```

Hide solution

In [ ]:

```
# Replacing the values
votes = votes.replace(['y', 'n'], (1, 0))

# Display the first 10 rows of the DataFrame
votes.head(10)
```

- (g) In a **DataFrame** named **X**, store the **explanatory** variables of the dataset (all columns except 'party'). For this, you can use the **drop** method of a **DataFrame**.
- (h) In a **DataFrame** named **y**, store the **target variable** ('party').

In [ ]:

```
# Insert your code here
```

Hide solution

In [ ]:

```
# Separation of the variables
X = votes.drop(['party'], axis = 1)
y = votes['party']
```

As for the regression problem, we will have to split the data set into 2 sets: a **training set** and a **test set**. As a reminder:

- The training set is used to **train the classification** model, that is to say find the parameters of the model which best separate the classes.
- The test set is used to **evaluate** the model on data that it has never seen. This evaluation will allow us to judge the **generalizability** of the model.

- (i) Import the **train\_test\_split** function from the **sklearn.model\_selection** submodule. Remember that this function is used as follows:

```
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.2)
```

- (j) Split the data into a training set (**X\_train**, **y\_train**) and a test set (**X\_test**, **y\_test**) keeping 20% of the data for the test set.

To **eliminate the randomness** of the **train\_test\_split** function, you can use the **random\_state** parameter with an integer value (for example **random\_state = 2**). This will make it so every time you use the function with the argument **random\_state = 2**, the datasets produced will be the same.

In [ ]:

```
# Insert your code here
```

Hide solution

In [ ]:

```
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.2, random_state = 2)
```

## Non-linear classification: K-Nearest Neighbors model

In order to assign a class to an observation, the K-Nearest Neighbors algorithm considers, as its name suggests, the **K nearest neighbors** of the observation and determines the **most represented class** among these neighbors.

Concretely, the algorithm is as follows:

- Suppose that **K = 5**.
- For an observation that we want to classify, we will look at the 5 points of the training set that are closest to our observation. The distance metric used is often the **euclidian norm**.
- If among the 5 neighbors, the majority is "democrat", then the observation will be classified "democrat".

To train a K-Nearest Neighbors model for our problem, we'll use the **KNeighborsClassifier** class from the `neighbors` submodule of `scikit-learn`.

The number **K** of neighbors to consider is entered using the parameter **n\_neighbors** of the **KNeighborsClassifier** constructor.

- (k) Run the following cell to import the **KNeighborsClassifier** class.

In [ ]:

```
from sklearn.neighbors import KNeighborsClassifier
```

- (l) Instantiate a **KNeighborsClassifier** model named **knn** which will consider the **6 nearest neighbors** for classification.
- (m) Using the **fit** method, train the model **knn** on the training dataset.
- (n) Using the **predict** method, perform a prediction on the **test** dataset. Store these predictions in **y\_pred\_test\_knn** and display the first 10 predictions.

In [ ]:

*# Insert your code here*

Hide solution

In [ ]:

```
# Instanciación del modelo
knn = KNeighborsClassifier(n_neighbors = 6)

# Entrenamiento del modelo sobre el juego de entrenamiento
knn.fit(X_train, y_train)

# Predicción sobre los datos de test
y_pred_test_knn = knn.predict(X_test)

# Afichaje de las 5 primeras predicciones
print(y_pred_test_knn[:10])
```

In the last part of this exercise, we will see how to evaluate the performance of our model using these predictions.

## Linear Classification: Logistic Regression

The logistic regression model is closely related to the **linear regression** model seen in the previous notebook.

They should not be confused since they do not solve the same types of problems:

- **Logistic Regression** is used for **classification** (predict classes).
- **Linear regression** is used for **regression** (predict a quantitative variable).

The linear regression model was defined with the following formula:

$$y \approx \beta_0 + \sum_{j=1}^p \beta_j x_j$$

Logistic regression no longer estimates  $y$  directly but the **probability** that  $y$  is equal to 0 or 1. Thus, the model is defined by the formula:

$$P(y = 1) = f(\beta_0 + \sum_{j=1}^p \beta_j x_j)$$

Where

$$f(x) = \frac{1}{1 + e^{-x}}$$

The  $f$  function, often called **sigmoid** or **logistic function**, transforms the linear combination  $\beta_0 + \sum_{j=1}^p \beta_j x_j$  into a value between 0 and 1 that can be interpreted as a **probability**:

- If  $\beta_0 + \sum_{j=1}^p \beta_j x_j$  is **positive**, then  $P(y = 1) > 0.5$ , so the predicted class of the observation will be **1**.
- If  $\beta_0 + \sum_{j=1}^p \beta_j x_j$  is **negative**, then  $P(y = 1) < 0.5$ , i.e.  $P(y = 0) > 0.5$ , so the predicted class of the observation will be **0**.

- (o) Import the **LogisticRegression** class from the `linear_model` submodule of `scikit-learn`.
- (p) Instantiate a **LogisticRegression** model named **logreg** without specifying constructor arguments.
- (q) Train the model on the training dataset.

- (r) Make a prediction on the **test** dataset. Store these predictions in **y\_pred\_test\_logreg** and display the first 10 predictions.

In [ ]:

```
# Insert your code here
```

Hide solution

In [ ]:

```
# Importation de la classe LogisticRegression sous-module linear_model de sklearn
from sklearn.linear_model import LogisticRegression

# Instanciation du modèle
logreg = LogisticRegression()

# Entraînement du modèle sur le jeu d'entraînement
logreg.fit(X_train, y_train)

# Prédiction sur les données de test
y_pred_test_logreg = logreg.predict(X_test)

# Affichage des 5 premières prédictions
print(y_pred_test_logreg[:10])
```

## 2. Evaluate the performance of a classification model

There are different metrics to evaluate the performance of classification models such as:

- The **accuracy**.
- The **precision and the recall**.
- The **F1-score**.

Each metric assesses the performance of the model with a different approach.

In order to explain these concepts, we will introduce 4 very important terms.

Arbitrarily, we will choose that the class '**republican**' will be the **positive class** (1) and '**democrat**' will be the **negative class** (0).

Thus, we will call:

- **True Positive (TP)** an observation classified as **positive** ('republican') by the model which is indeed **positive** ('republican').
- **False Positive (FP)** an observation classified as **positive** ('republican') by the model which was actually **negative** ('democrat').
- **True Negative (TN)** an observation classified as **negative** ('democrat') by the model and which is indeed **negative** ('democrat').
- **False Negative (FN)** an observation classified as **negative** ('democrat') by the model which was actually **positive** ('republican').

## Predictions

0

1

In [ ]:

```
# Insert your code here
```

Hide solution

In [ ]:

```
from sklearn.metrics import confusion_matrix

# Computation and display of the confusion matrix
conf_matrix = confusion_matrix(y_test, y_pred_test_knn)
print("Confusion Matrix:\n", conf_matrix)

print("\nThe knn model made", conf_matrix[0,1], "False Positives.")

# Computation of the accuracy, precision and recall
(TN, FP), (FN, TP) = confusion_matrix(y_test, y_pred_test_knn)
n = len(y_test)

print("\nKNN Accuracy:", (TP + TN) / n)
print("\nKNN Precision:", TP / (TP + FP))
print("\nKNN Recall: ". TP / (TP + FN))
```

The display of the confusion matrix can also be done with the **pd.crosstab** function as we had done in a previous notebook:

```
pd.crosstab(y_test, y_pred_test_knn, rownames = ['Reality'], colnames = ['Prediction'])
```

Which in our case will produce the following DataFrame :

	Prediction	democrat	republican
Reality			
democrat		48	5
republican		2	32

For this dataset, the KNN model performs quite well. When the classes are **balanced**, i.e. there are about as many positives as there are negatives in the dataset, accuracy is a good enough metric to assess the performance.

However, as you will see later, when a **class is dominant**, precision and recall are much more relevant metrics.

If you think you cannot remember the formulas for the metrics of accuracy, precision and recall, do not worry! The **sklearn.metrics** submodule contains functions to calculate them quickly:

```
accuracy_score(y_test, y_pred_test_knn)
>>> 0.9195402298850575
```

- (e) Import the **accuracy\_score**, **precision\_score** and **recall\_score** functions from the **sklearn.metrics** submodule.
- (f) Display the confusion matrix of the predictions made by the **logreg** model using **pd.crosstab**.
- (g) Calculate the accuracy, precision and recall of model predictions **logreg**. To use the **precision\_score** and **recall\_score** metrics, you will need to fill in the argument **pos\_label = 'republican'** in order to specify that the 'republican' class is the positive class.

In [ ]:

```
# Insert your code here
```

Hide solution

In [ ]:

```
from sklearn.metrics import accuracy_score, precision_score, recall_score

# Computation and display of the confusion matrix
pd.crosstab(y_test, y_pred_test_logreg, rownames=['Reality'], colnames=['Prediction'])

# Computation of the accuracy, precision and recall
print("\nLogReg Accuracy:", accuracy_score(y_test, y_pred_test_logreg))

print("\nLogReg Precision:", precision_score(y_test, y_pred_test_logreg, pos_label = 'republican'))

print("\nLogReg Recall:", recall_score(y_test, y_pred_test_logreg, pos_label = 'republican'))
```

The `classification_report` function of the `sklearn.metrics` submodule displays all these metrics for each class.

- (h) Import the `classification_report` function from the `sklearn.metrics` submodule.
- (i) Display using the `print` and `classification_report` functions the classification reports of the models **Logreg** and **knn** on the test set.

In [ ]:

*# Insert your code here*

Hide solution

In [ ]:

```
from sklearn.metrics import classification_report

print("LogReg report:\n", classification_report(y_test, y_pred_test_logreg))

print("\n\n")

print("KNN report:\n", classification_report(y_test, y_pred_test_knn))
```

The classification report is a little more complete than what we have done so far. It contains an additional metric: the **F1-Score**.

The F1-Score is a sort of average between precision and recall. The F1-Score adapts very well to classification problems with balanced or unbalanced classes.

For most classification problems, the model with the highest F1-Score will be considered the model whose recall and precision **performances** are the most **balanced**, and is therefore **preferable** to others.

- (j) Import the `f1_score` function from submodule `sklearn.metrics`.
- (k) Compare the F1-Scores of the models **knn** and **logreg** on the test set. Which model has the best performance? As for the recall and the precision, it will be necessary to fill in the argument `pos_label = 'republican'`.

In [ ]:

*# Insert your code here*

Hide solution

In [ ]:

```
from sklearn.metrics import f1_score

print("F1 KNN:", f1_score(y_test, y_pred_test_knn, pos_label = 'republican'))

print("F1 LogReg:", f1_score(y_test, y_pred_test_logreg, pos_label = 'republican'))
```

## Conclusion and recap

Scikit-learn offers many classification models that can be grouped into two families:

- Linear **models** like **LogisticRegression** .
- Non-linear **models** like **KNeighborsClassifier** .

The implementation of these models is done in the same way for **all** models of scikit-learn:

- **Instantiation** of the model.
- **Training** of the model: **model.fit(X\_train, y\_train)** .
- **Prediction**: **model.predict(X\_test)** .



Validate