

```
In [ ]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
```

```
In [ ]: import sklearn as sk
```

## I- Régression

```
In [ ]: #un exemple simple en deux dimensions

rng = np.random.RandomState(42) #pour générer les mêmes données

#constituer un exemple de data
x = 10 * rng.rand(50)
print('la taille de notre échantillon est :',x.shape)

y=2*x-1 + rng.randn(50)
#afficher data y=f(x) [y en fonction de x] comme un nuage de points
plt.scatter(x, y);
```

### Y-a-t-il une relation entre $x$ et $y$ : trouver $f$ tel que $y=f(x)$ ? (Cours 1)

Pour répondre à cette question, nous allons supposer que  $f$  est une fonction de la forme  $f(x) = a * x + b$  avec  $a$  et  $b$  sont des réels à déterminer.

Input :  $(x_i, y_i)$ , pour  $i = 1 \dots 50$

on a  $y_i = a * x_i + b$ , pour  $i = 1 \dots 50$

qui forme un système linéaire facile à résoudre (plus de données que d'inconnus)

### Formulation (cours semaine 1)

```
In [ ]: # On peut résoudre ce problème de régression linéaire avec sklearn

# on choisit et charge le modèle
from sklearn.linear_model import LinearRegression

X = x[:, np.newaxis]
print('la tailles des entrées est :',X.shape)

models = LinearRegression(fit_intercept=True)
models.fit(X, y)
```

```
In [ ]: a=models.coef_
print('-'*5,'la solution','-'*5)
```

```
print('la valeur trouvée de a est : ', a[0])

b=models.intercept_
print('la valeur trouvée de b est : ', b)
```

Si maintenant on a un nouveau  $x_{new} = 2.5$  qui est différent de tous les  $x_i$  observés

On peut trouver son image  $y_{new}$  avec  $y_{new} = a * x_{new} + b$

```
In [ ]: xnew=np.array([2.50])
ynew = models.predict(xnew.reshape(-1, 1))
print(ynew)
```

On peut aussi appliquer la même méthode sur xnew comme tableau de valeurs au lieu d'un seul scalaire

```
In [ ]: xnew=np.linspace(-1,10,5)
#s'assurer d'avoir le bon format
xnew=xnew[:, np.newaxis]

ynew = models.predict(xnew)
print(ynew)
```

Vérification visuelle

```
In [ ]: plt.scatter(x, y,color='k');# données apprentissage en noir
plt.scatter(xnew, np.zeros(xnew.shape[0]),color='b');# x_i non observés en b
plt.scatter(xnew, ynew,color='r');# y_i prédit ave la régression linéaire (x
```

```
In [ ]: #on peut aussi afficher la fonction f
plt.scatter(x, y,color='k');
#plt.scatter(xnew, ynew);
plt.plot(xnew, ynew,'r');
#l'erreur est donnée par la somme cumulée des distances
#entre les points en noir et la droite en rouge

ypred=models.predict(X)
print(ypred.shape)
print('Erreur quadratique moyenne : ',np.mean((y-ypred)**2))
```

## Nous venons de faire notre premier exemple pour le cas simple $x \in \mathbb{R}$ et $y \in \mathbb{R}$

On peut généraliser ce résultat quelque soit la taille de  $x$  :  $x \in \mathbb{R}^d$ , pour toute dimension  $d$

Exemple :

```
In [ ]: from mpl_toolkits.mplot3d import Axes3D
#constituer un exmple de data
```

```

x = np.array(10 * rng.rand(100,2))
y=2*np.inner(np.array([-1,1]), x)+ 2*rng.randn(x.shape[0])

fig=plt.figure()
ax = fig.add_subplot(111, projection='3d')

ax.scatter(x[:,0], x[:,1],y,c='b', marker='o');
ax.set_xlabel('valeur de x[:,0]')
ax.set_ylabel('valeur de x[:,1]')
ax.set_zlabel('valeur de y ')

plt.show()

```

```

In [ ]: model = LinearRegression(fit_intercept=True)
        model.fit(x, y)

```

```

In [ ]: xnew = np.array(10 * rng.rand(1000,2))
        ynew = model.predict(xnew)

```

```

In [ ]: fig=plt.figure()
        ax = fig.add_subplot(111, projection='3d')

        ax.scatter(x[:,0], x[:,1],y,c='b', marker='o');
        ax.scatter(x[:,0], x[:,1],-y,c='b', marker='o');
        ax.set_xlabel('valeur de x[:,0]')
        ax.set_ylabel('valeur de x[:,1]')
        ax.set_zlabel('valeur de y ')

        ax.scatter(xnew[:,0], xnew[:,1],ynew,c='r', marker='*');

        plt.show()

```

```

In [ ]: # régression avec statsmodels avec d'autres modèles statistiques
        import statsmodels.api as sm

        x = [[0, 1], [5, 1], [15, 2], [25, 5], [35, 11], [45, 15], [55, 34], [60, 35]]
        y = [4, 5, 20, 14, 32, 22, 38, 43]
        x, y = np.array(x), np.array(y)

```

```

In [ ]: x = sm.add_constant(x)
        print(x, ' \n avec une taille : \n', x.shape)

```

```

In [ ]: #créer le modèle OLS (ordinary least squares) qui minimise l'erreur quadratique
        model = sm.OLS(y, x)

```

```

In [ ]: #avec plus de détails sur le modèle
        results = model.fit()
        #print(results.summary())

```

```

In [ ]: #utiliser le modèle pour la prédiction
        x_new = sm.add_constant(np.arange(10).reshape((-1, 2)))
        y_new = results.predict(x_new)
        print(y_new)

```

Bien que nous pouvons nous contenter d'un modèle linéaire dans un premier temps, il est bon à savoir que d'autres modèles existent.

## Règle d'or : on choisit selon l'application !!

On retient que :

- Le modèle dans sklearn est facile à comprendre et à implémenter
- Statsmodels offre des outils statistiques plus avancés

## Un test de la régression linéaire sur des données réelles

Dans cette partie nous allons découvrir comment on peut utiliser python, numpy, matplotlib et sklearn pour classer des fleurs.

Le dataset (jeu de données) est une base de données des caractéristiques de trois espèces de fleurs d'Iris (Setosa, Versicolour et Virginica).

Description :

- chaque ligne est une observation des caractéristiques d'une fleur d'Iris ;
- caractéristiques : longueur et largeur de (sépal, pétales) ;
- dans cet exemple, le dataset contient 150 observations (50 observations par espèce) ;
- plus d'information  
[https://fr.wikipedia.org/wiki/Iris\\_de\\_Fisher](https://fr.wikipedia.org/wiki/Iris_de_Fisher) .



```
In [ ]: #importer les bibliothèques

#pour l'affichage (si déjà fait pour np, plt)
%matplotlib inline

#des datasets dans sklearn pour les tests
from sklearn import datasets
```

```
In [ ]: #charger la base
iris = datasets.load_iris()
#vérifier le type de la variable iris
print(type(iris))
#vérifier quel est le type de données
print(type(iris.data))
```

```
#vérifier les dimensions
print(iris.data.shape)
```

## Le choix des variables

```
In [ ]: X = iris.data[:, :2] # Utiliser les deux premières colonnes afin d'avoir un

print(np.unique(iris.target))
#on va garder deux classes seulement pour un test simple
y = (iris.target != 0) * 1 # re-étiquetage des fleurs
print(X.shape)
print(np.unique(y))
```

## Pour mieux comprendre notre data, on va visualiser des exemples

```
In [ ]: #visualisation des données
plt.figure(figsize=(10, 6))
plt.scatter(X[y == 0][:, 0], X[y == 0][:, 1], color='b', label='classe 0')
plt.scatter(X[y == 1][:, 0], X[y == 1][:, 1], color='r', label='classe 1')
plt.legend();
```

Une petite vérification visuelle montre que les deux classes peuvent être séparées par une droite. On dira que elles peuvent être **linéairement** séparées.

## Pour la suite on va utiliser une régression logistique pour exploiter cette séparation.

```
In [ ]: #charger le modèle
from sklearn.linear_model import LogisticRegression

model = LogisticRegression(C=1e20) # Régression logistique
# Entraînement du modèle avec toutes les données
model.fit(X, y)
```

```
In [ ]: Xnew = np.array([
    [5.5, 2.5],
    [7, 3],
    [3, 2],
    [5, 3]
])
```

```
In [ ]: model.predict(Xnew)
```

```
In [ ]: type(Xnew)
```

Analyse des résultats :

- \* La première observation [5.5, 2.5] est de classe 1
- \* La deuxième observation [7, 3] est de classe 1
- \* La troisième observation [3,2] est de classe 0
- \* La quatrième observation [5,3] est de classe 0

```
In [ ]: #vérification visuelle ,,,

#visualisation des données
plt.figure(figsize=(10, 6))
plt.scatter(X[y == 0][:, 0], X[y == 0][:, 1], color='b', label='class 0')
plt.scatter(X[y == 1][:, 0], X[y == 1][:, 1], color='r', label='class 1')

s = np.random.rand(*Xnew[:, 0].shape) * 800 + 500
print(s.shape)
Color='kygm' #noir jaune vert magenta
for i in range(Xnew.shape[0]):
    plt.scatter(Xnew[i, 0], Xnew[i, 1],s[i], color=Color[i],marker=r'$\clubsuit$')
plt.legend();
```

## II- Classification

```
In [ ]: # classification d'images de chiffres
from sklearn.datasets import load_digits

digits = load_digits()
digits.images.shape
```

```
In [ ]: #visualiser les données
fig, axes = plt.subplots(10, 10, figsize=(8, 8),
                        subplot_kw={'xticks':[], 'yticks':[]},
                        gridspec_kw=dict(hspace=0.1, wspace=0.1))

for i, ax in enumerate(axes.flat):
    ax.imshow(digits.images[i], cmap='binary', interpolation='nearest')
    ax.text(0.05, 0.05, str(digits.target[i]),transform=ax.transAxes, color=
```

```
In [ ]: X = digits.data
X.shape
```

```
In [ ]: y = digits.target
y.shape
```

## Préparation pour l'apprentissage : train et test

```
In [ ]: #former les bases d'apprentissage et de test
from sklearn.model_selection import train_test_split
Xtrain, Xtest, ytrain, ytest = train_test_split(X, y,test_size=0.20, random_
```

```
In [ ]: print(Xtrain.shape)
print(Xtest.shape)
```

```
print('pourcentage: ', Xtrain.shape[0]/X.shape[0])  
  
print(ytrain.shape)
```

## Classification avec l'arbre de décision

```
In [ ]: from sklearn.tree import DecisionTreeClassifier  
  
Arbre_decision = DecisionTreeClassifier(random_state=0, max_depth=20)  
clf = Arbre_decision.fit(Xtrain, ytrain)
```

```
In [ ]: from sklearn.metrics import accuracy_score  
ypredit = clf.predict(Xtest)  
  
accuracy_score(ytest, ypredict)
```

```
In [ ]: from sklearn import metrics  
print(metrics.confusion_matrix(ytest, ypredict))
```

## Classification avec le plus proche voisin

```
In [ ]: from sklearn.neighbors import KNeighborsClassifier  
  
KNN = KNeighborsClassifier()  
clf = KNN.fit(Xtrain, ytrain)  
  
ypredit = clf.predict(Xtest)  
  
accuracy_score(ytest, ypredict)
```

```
In [ ]: print(metrics.confusion_matrix(ytest, ypredict))
```

## Classification avec SVM

```
In [ ]: from sklearn import svm  
  
clf = svm.SVC(gamma=0.001)  
clf.fit(Xtrain, ytrain)
```

```
In [ ]: from sklearn.metrics import accuracy_score  
ypredit = clf.predict(Xtest)  
accuracy_score(ytest, ypredict)
```

```
In [ ]: ypredict = clf.predict(Xtest)  
print(metrics.confusion_matrix(ytest, ypredict))
```

# \*\*\*\*\* Pratique avec le dataset Titanic

```
In [ ]: import pandas as pd  
  
df=pd.read_csv('./titanic/train.csv')
```

```
In [ ]: df.head(5)
```

```
In [ ]: df.shape
```

```
In [ ]: df.describe()
```

```
In [ ]: df.info()
```

## La signification des colonnes :

PassengerID : identifiant

Survived : 0 si ce passager n'a pas survécu, 1 sinon

Pclass : la classe (1, 2 ou 3)

Name : le nom du passager

Sex : femme ou homme (sexe)

Age : l'âge (en années)

SibSp : le nombre de frère, soeur et/ou épouse à bord

Parch : le nombre de parent et/ou d'enfant à bord

Ticket : numéro du ticket

Fare : prix du billet

Cabin : numéro de cabine

Embarked : port d'embarquement(C = Cherbourg, Q = Queenstown, S = Southampton)

## Exo 1 : dans un premier temps on explore pour comprendre les données

- Lire et explorer le dataset à partir des fichiers csv.



- Vérifier que le dataset contient des valeurs manquantes.
- Proposer une solution pour remplacer les valeurs manquantes.
- Est il possible de visualiser les données ?
- Créer une nouvelle colonne qui mentionne le sexe du passager.
- Combien de personnes ont survécu ?
- Est ce que la Pclass a plus d'importance que l'âge pour la survie ?
- Pour répondre à cette question, nous allons tester la classification : survived en fonction de l'âge (respectivement Pclass).
- On peut tirer une conclusion à partir de la performance mais relativement au modèle utilisé.

## Exo 2 : dans un deuxième temps on va transformer quelques colonnes

- Transformer le sexe en nombre.
- On teste si (âge et sexe) est mieux que (classe, sexe).
- On teste la performance de trois différentes méthodes pour la classification.
- Proposer comment transformer d'autres colonnes.
- Tester l'importance de chaque colonne à part, des combinaisons, etc.
- Conclure après plusieurs tests : modèles + features.