

## Scikit-Learn

Réalisé par :

Ayoub Oukchiren

Encadré par:

Mr. Mohammed Reda





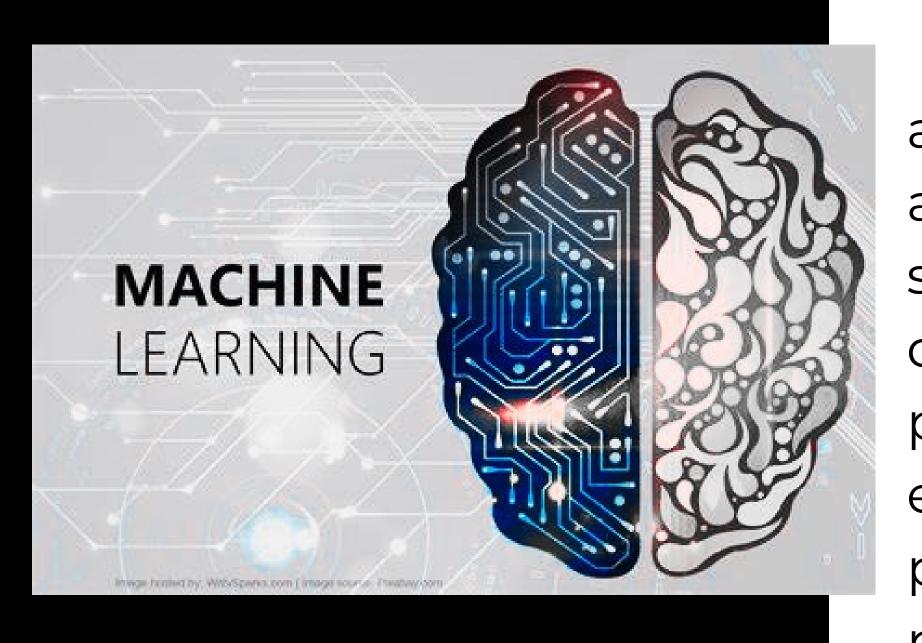
## Plan

Introduction Définition Installation Fonctionnalités Application Conclusion

## Introduction







Maching learning est une application de l'intelligence artificielle (IA) qui fournit aux systèmes la capacité d'apprendre et de s'améliorer automatiquement à partir de l'expérience sans être explicitement programmés. parmi les bibliothèques ML les plus populaires on trouve scikit-learn

## Définition





### c'est quoi?

scikit learn(sklearn) est la bibliothèque la plus puissante pour le machine learning en Python, elle contient de nombreux outils efficaces pour l'apprentissage automatique et la modélisation statistique, notamment la classification, la régression et le clustering.



## Installation



## prérequis:



Installing Scikit-Learn



python3 --version

La sortie devrait ressembler à :

Python 3.8.2



#### on tape la commande:

```
pip install scikit-learn
```

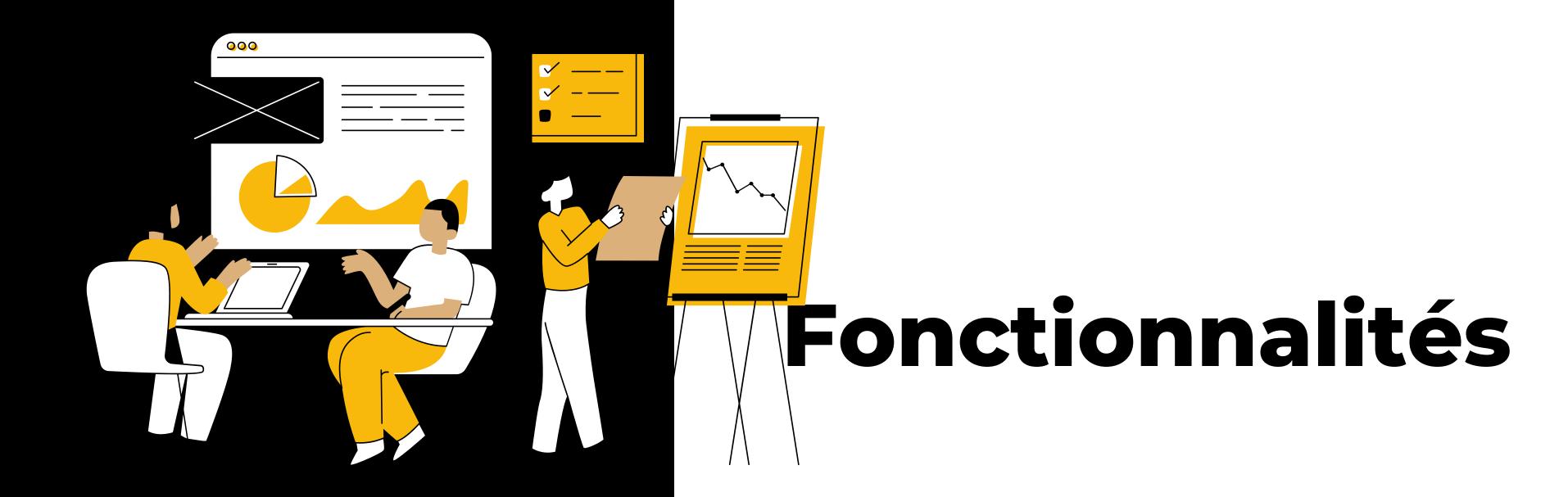
#### 2eme methode:

```
pip install -U scikit-learn
```

#### pour verifier l'installation

```
::\>python3 -m pip show scikit-learn
lame: scikit-learn
/ersion: 0.24.0
Summary: A set of python modules for machine learning and data mining
lome-page: http://scikit-learn.org
Author: None
Author-email: None
.icense: new BSD
.ocation: c:\python38\lib\site-packages
Requires: numpy, scipy, threadpoolctl, joblib
Required-by:
```

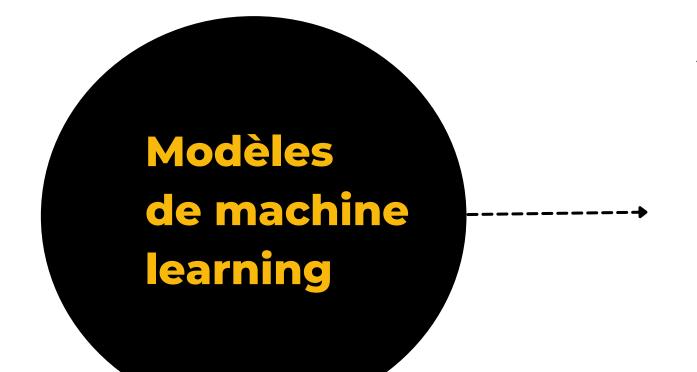
- NumPy 1.13.3+
- SciPy 0.19.1+
- Joblib 0.11+
- threadpoolctl 2.0.0+



Scikit-learn propose de nombreuses fonctionnalités pour aider les utilisateurs à résoudre des problèmes de machine learning



propose des outils sklearn normaliser, standardiser, encoder, et nettoyer les données pour les préparer à l'utilisation avec des algorithmes de machine learning.



sklearn propose une variété de modèles de machine learning tels que les régressions linéaires et logistiques, les arbres de décision, les k plus proches voisins, les réseaux de neurones, etc.



sklearn est compatible avec de nombreuses autres librairies de machine learning et de traitement de données, comme NumPy, pandas, et TensorFlow.



# les types de modèles de machine learning disponibles dans sklearn

1- Régressions: y compris la régression linéaire, la régression Ridge, la régression Lasso, la régression polynomiale, etc.



- 2- Arbres de décision: y compris l'arbre de décision, le random forest et l'extra tree.
- 3- k-moyennes: pour l'agrégation de données en groupes (clustering)
- 4- k-plus proches voisins (k-NN)
- 5- Réseaux de neurones : MLP, RBM etc.

• • • •

```
10
11
12
13
14
15
16
         port = "80"
 20
 21
 22
           path = "/"
 23
 24
 25
 26
 27
```

```
name = "ecs-application-load-balancer"
# target group with health_check
resource "aws_lb_target_group" "ecs-target-group-lb" {
  name = "ecs-target-group"
  protocol = "HTTP"
  vpc_id = aws_vpc.ecs_test.id
  depends_on = [aws_lb.ecs-load-balancer]
  health_check {
     protocol = "HTTP"
     matcher = "200"
     interval = 15
     timeout = 10
      healthy_threshold = 2
```

## Application

### la problématique

comment prédire les valeurs de la colonne "salary" dans le fichier (SalaireDatal.csv) en se basant sur le dataset de fichier (Salary\_Data.csv) en utilisant sklearn?

	Α	В
1	YearsExperie	Salary
2	1.1	39343.00
3	1.3	46205.00
4	1.5	37731.00
5	2.0	43525.00
6	2.2	39891.00
7	2.9	56642.00
8	3.0	60150.00
9	3.2	54445.00
10	3.2	64445.00
11	3.7	57189.00
12	3.9	63218.00
13	4.0	55794.00

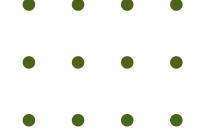
Α	В
YearsExperie	Salary
3	
5	
7	
9	
11	
13	
15	
17	
19	
21	
23	
25	

Salairy\_Data.csv

SalaireData1.csv

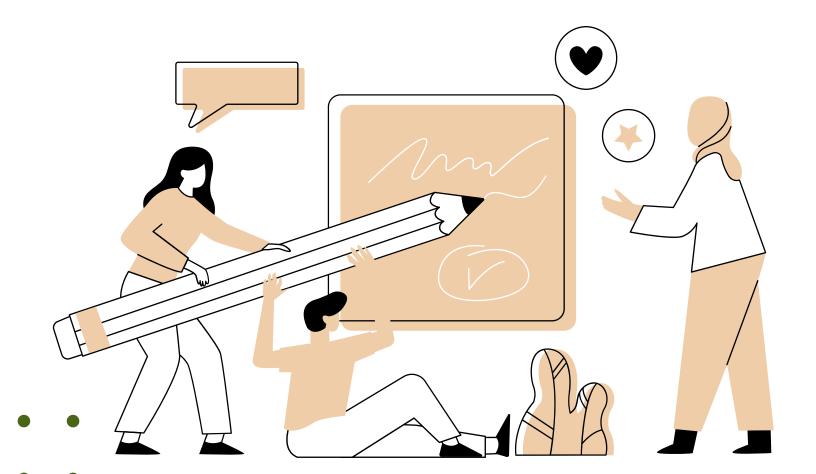
• • •

• • •



### l'objectif:

est de construire un bon modèle d'apprentissage automatique capable de prédire les salaires en fonction nombre d'années d'experiences

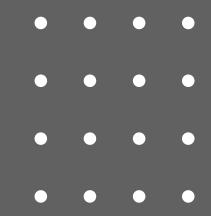


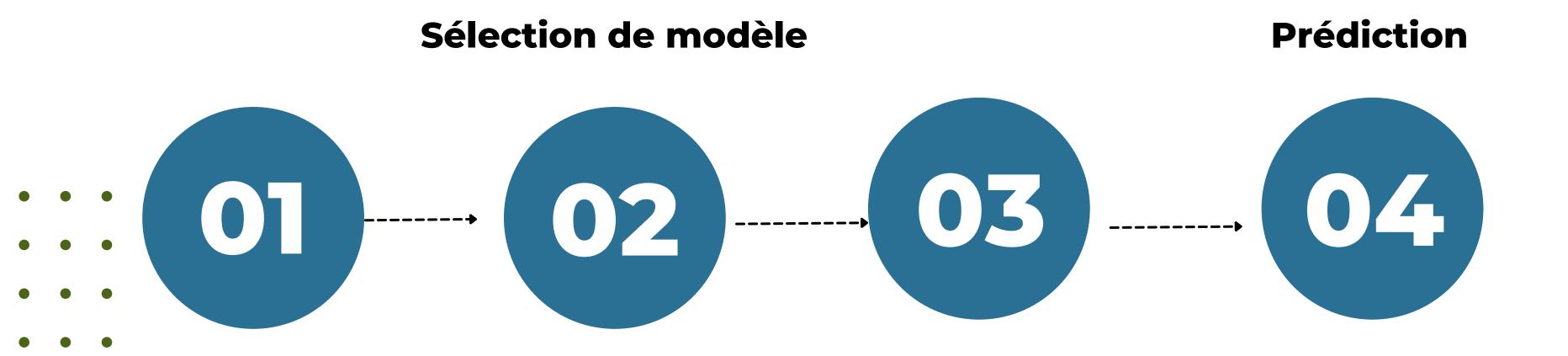
Α	В
YearsExperie	Salary
3	
5	
7	
9	
11	
13	
15	
17	
19	
21	
23	
25	

SalaireData1.csv

#### LES ETAPES SUIVIES:

Prétraitement des données





Entraînement et évaluation

### Importer les paquets

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
from sklearn.model selection import train test split
from sklearn.linear_model import LinearRegression
```

## Importer le dataset

```
data = pd.read_csv("Salary_Data.csv")
```

#### جِــامعة الحسن الثــاني يالدار البيـضاء - Records : 20 0020 Los + 1 المدد + 1 المدد - ENIVERSITE HASSAN TI DE CASABLANCA



### afficher le dataset (Salary\_data.csv)

```
data.head()
```

YearsExperience		Salary
0	1.1	39343.0
1	1.3	46205.0
2	1.5	37731.0
3	2.0	43525.0
4	2.2	39891.0

## ENIVERSITÉ HASSAN II DE CASABLANCA

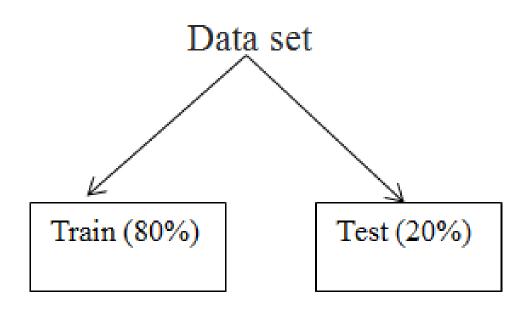
### Prétraitement des données

```
x = data.iloc[:,0].values
y = data.iloc[:,1].values
```

x stock toutes les lignes de la première colonne de data (YearsExperience) y stock toutes les lignes de la deuxième colonne de data (Salary)

```
x_train,x_test,y_train,y_test = train_test_split(x,y,test_size = 0.2,random_state=3)
```

train\_test\_split permet de diviser un ensemble de données en deux parties : un ensemble d'entraînement et un ensemble de test. tel que 20% des données seront utilisées pour tester le modèle (x\_test,y\_test) et les 80% pour l'entraînement (x\_train,y\_train)





## Définir le modèle de régression linéaire

ml = LinearRegression()

La variable **ml** est un objet de la classe de régression linéaire que on va utiliser pour entraîner un modèle sur les données d'entraînement, puis effectuer des prédictions sur les données de test.

## Entraîner le modèle sur les données d'entraînement

ml.fit(x\_train, y\_train)

LinearRegression()

fit() est une méthode de l'objet de régression linéaire ml, qui permet d'entraîner le modèle sur des données d'entraînement. Les arguments passés à la méthode fit sont les données d'entraînement x\_train et les étiquettes correspondantes y\_train

• • • •

• • • •

• • • •



## Faire des prédictions sur les données du test

y\_pred = ml.predict(x\_test)

La méthode predict prend en entrée les données pour lesquelles on souhaite effectuer des prédictions et renvoie les prédictions sous forme d'un tableau de valeurs.

y\_pred est un tableau qui contient les prédictions du modèle pour les données de test







## Comparer les résultats de la prédiction

- le nuage de points rouge correspond aux données réelles
- la ligne bleue correspond aux valeurs prédites
- on constate que les nuages de points sont proches de la droite blue donc on peut dire que notre modèle généralise assez bien

```
plt.scatter(x_test,y_test,color= 'red')
 plt.plot(x_test,y_pred,color='blue')
 plt.xlabel("Years of Experience")
 plt.ylabel("Salary")
 Text(0, 0.5, 'Salary')
120000
110000
100000
 90000
 80000
 70000
 60000
                              Years of Experience
```

## calculer la performance du modèle avec score()

ml.score(x\_test,y\_test)

0.9695039421049821

score(): renvoie un score de performance qui varie entre 0 et 1, on a trouvé un score proche de 1 signifie que notre modèle est performant donc nous pouvons utiliser ce modèle pour prédire les valeurs manquantes dans la colonne Salary du fichier salaire Data1.csv

• • • •

• • • •

• • • •

#### منامعة الحسن الثنائي بالدار البيضاء USOO OSI R GE-099EE بدعة HELLAND UNIVERSITÉ ILUSSAN II DE CASASILANC



### afficher le dataset (SalaireData1.csv)

```
data2 = pd.read_csv("SalaireData1.csv")
```

data2.head()

YearsExperience		Salary
0	3	NaN
1	5	NaN
2	7	NaN
3	9	NaN
4	11	NaN

• • • •

• • • •

## utiliser le modèle pour prédire les salaires



```
x_predict = data2[["YearsExperience"]].values

y_predict= ml.predict(x_predict)
```

x\_predict contient les années d'expérience pour lesquelles on souhaite effectuer des prédictions de salaires

y\_predict contient les prédictions de salaires

### résultat

```
y_predict.shape
df=pd.DataFrame(data2[["YearsExperience"]])
df['Salary']=y_predict
```

df

	YearsExperience	Salary
0	3	54453.467948
1	5	73290.601896
2	7	92127.735843
3	9	110964.869791
4	11	129802.003738
5	13	148639.137685
6	15	167476.271633
7	17	186313.405580
8	19	205150.539527
9	21	223987.673475
10	23	242824.807422
11	25	261661.941370
12	27	280499.075317

• • • •

• • • •

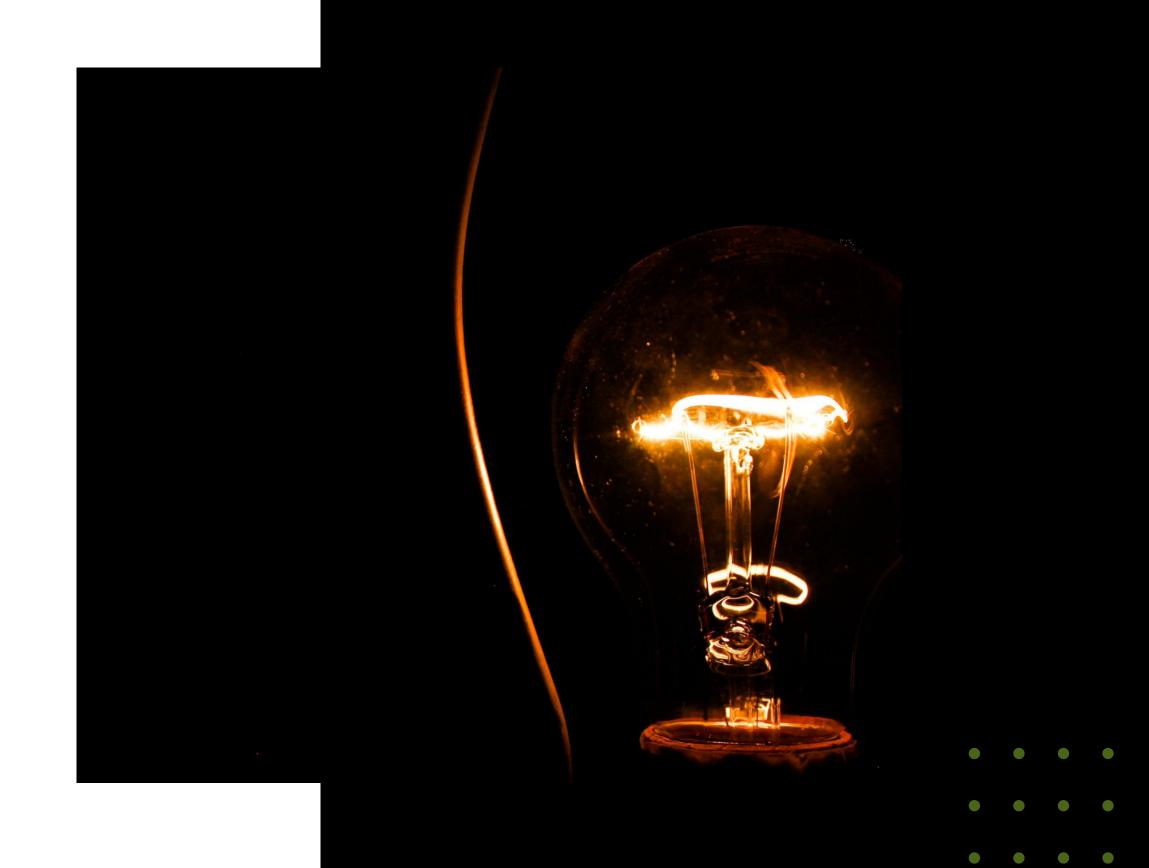
• • • •

جـــامعة الحسن الثــاني بالتدار البيضاء +LOLUSE : الحال USOO OSI E SELGHOSE ENTYERSITE HASSAN II DE CASABLANCA



## Conclusion

scikit-learn rest un outil incontournable pour tous les chercheurs et les développeurs d'IA qui cherchent à résoudre des problèmes d'apprentissage automatique







## 

pour votre attention