

MINISTÈRE DE L'ÉDUCATION NATIONALE, DE
L'ENSEIGNEMENT SUPÉRIEUR ET DE LA RECHERCHE
SCIENTIFIQUE

ÉCOLE NATIONALE D'INGÉNIEURS
ABDERHAMANE BABA TOURE (ENI~ABT)

ENI  ABT

RÉPUBLIQUE DU MALI

Un peuple ~ Un but ~ Une foi



EXPOSE

Thème

ALGORITHME K-VOISINS LES PROCHES
K-NEAREST NEIGHBORS

Présenté et soutenu par : Souleymane Kodjo

Abdrahamane Idrissa Doumbia

Ouzairou Djiré

Responsable pédagogique : Dr A. SIDIBE

Sommaire

Chapitre 1 : Introduction	4
1.1 Cas d'usage et applications.....	4
Chapitre 2 : Algorithme K voisins les plus proches	5
2.1 Classification « Quel classe ? »	5
2.2 Régression « Combien ? »	6
Chapitre 3 : Principe de l'algorithme.....	7
3.1 Sélectionne le nombre k de voisin	7
3.2 On calcule les distances entre la donnée u et chaque donnée appartenant à E à l'aide de la fonction d	7
3.2.1 Distance Euclidienne.....	7
3.2.2 Distance Manthattan.....	7
3.3 On retient les k données du jeu de données E les plus proches de u	8
3.4 On attribue à u la classe qui est la plus fréquente parmi les k données les plus proches.....	8
3.5 Ecriture algorithmique	8
Chapitre 4 : Etude d'un Exemple	10
4.1 Jeux de données	10
4.2 Les logiciel utilisé	Erreur ! Signet non défini.
Anaconda navigator.....	10
Jupyter notebook	10
Chrome	10
4.3 Bibliothèques Python utilisées	10
4.3.1 Pandas.....	10
4.3.2 NumPy.....	11
4.3.3 Matplotlib	11
4.3.4 Sklearn.....	11
4.4 Importation des bibliothèques	11
4.5 Importation des données	11
4.6 Formatage	11
4.7 Première visualisation des données	12
4.8 Echantillon a prédire.....	13
4.9 Visualisation de Echantillon a prédire dans le graphe.....	13

4.10	Séparation du jeux de donnée en training set et testing set.....	14
4.11	Algorithmique KNN pour $K=2$	14
4.12	Optimisation du taux d'erreurs	15
•	Sous Forme Graphique.....	15
4.13	Création d'une fonction python résultat qui retourne la prédiction	16
4.14	Visualisation finale et prédiction	16
5	Conclusion.....	18

Chapitre 1 : Introduction

L'Apprentissage automatique ou machine Learning (en anglais) peut être définie comme une branche de l'intelligence artificiel englobant de nombreuses méthodes permettant de créer automatiquement les modèles à partir de données. Ces méthodes sont en fait des algorithmes.

Un programme informatique traditionnel effectue une tâche en suivant des instructions précises et donc systématiquement de la même façon. Au contraire, un système machine Learning ne suit pas d'instructions, mais à donner la capacité aux ordinateurs d'apprendre à partir de données, c'est-à-dire d'améliorer leur performance à résoudre les tâches sans être explicitement programmée pour chacun. En conséquence, ces performances s'améliorent au fil de son « Entraînement » à mesure que l'algorithmique est exposée à davantage de données.

1.1 Cas d'usage et applications

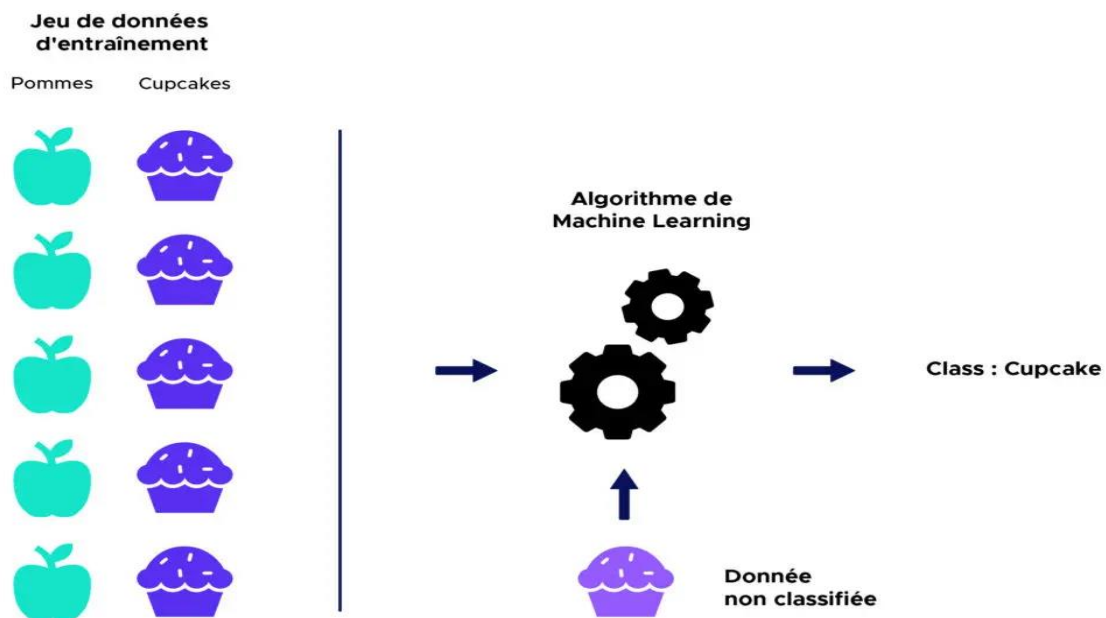
Le Machine Learning alimente de nombreux services modernes très populaires pour la recommandation de produits. On peut citer comme : Netflix, YouTube, Amazon ...

Il en va de même pour les moteurs de recherche web Google ou Baidu et les assistants vocaux tels Siri et Alexa.

Il est aussi utilisé dans la prédiction des prix, la détection des spam ...

Chapitre 2 : Algorithme K voisins les plus proches

L'algorithme des k plus proches voisins ou encore appelée KNN de l'anglais K-Nearest Neighbors, appartient à la famille des algorithmes d'apprentissage automatique (Machine Learning). L'idée d'apprentissage automatique ne date pas d'hier, puisque le terme de Machine Learning a été utilisé pour la première fois par l'informaticien américain Arthur Samuel en 1959. L'algorithme des k plus proches voisins est un algorithme de d'apprentissage supervisé. En apprentissage supervisé, un algorithme reçoit un ensemble de données qui est étiqueté avec des valeurs de sorties correspondantes sur lequel il va pouvoir s'entraîner et définir un modèle de prédiction. Cet algorithme pourra par la suite être utilisé sur de nouvelles données afin de prédire leurs valeurs de sorties correspondantes.



L'algorithme d'apprentissage supervisé permet de répondre de deux problèmes :

2.1 Classification « Quel classe ? »

L'algorithme de classification est d'abord alimenté par un nombre fini d'exemple catalogués qu'il utilise pour son apprentissage. Cet apprentissage lui permet de sélectionner une « hypothèse », ou règle censée commettre peu d'erreurs de classification sur ces exemples futures

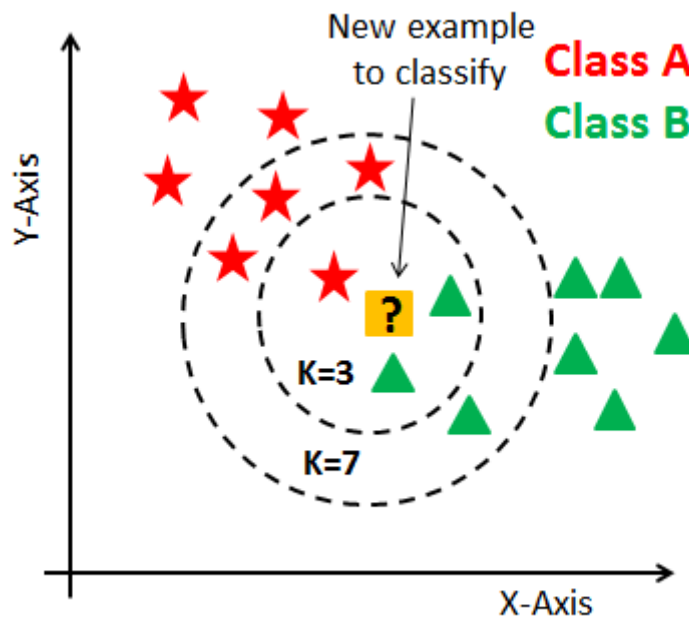
Le résultat est une classe d'appartenance. Un objet d'entrée est classifié selon le résultat majoritaire des statistiques de classes d'appartenance de ses k plus proches voisins.

2.2 Régression « Combien ? »

L'algorithme de régression est alimenté à partir de la variable cible ou de la variable explicative (Y), le modèle a pour but de faire une prédiction grâce à des variables dites explicatives (X) ou prédictives.

Le résultat est la valeur pour cet objet. Cette valeur est la moyenne des valeurs des k les plus proches.

KNN est une méthode non paramétrique dans laquelle le modèle mémorise les observations de l'ensemble d'apprentissage pour la classification des données de l'ensemble de test.



En effet, cet algorithme est qualifié comme paresseux, car il n'apprend rien pendant la phase d'entraînement. Pour prédire la classe d'une nouvelle donnée d'entrée, il va chercher ses K voisins les plus proches (en utilisant la distance euclidienne, ou autres) et choisira la classe des voisins majoritaires.

Chapitre 3 : Principe de l'algorithme

L'algorithme de k plus proches voisins ne nécessite pas de phase d'apprentissage à proprement parler, il faut juste stocker le jeu de données d'apprentissage. Soit un ensemble E contenant n données labellisées : $E = \{(y_i, \vec{x}_i)\}$ avec i compris entre 1 et n, où y_i correspond à la classe (le label) de la donnée i et où le vecteur \vec{x}_i de dimension p ($\vec{x}_i = (x_{1i}, x_{2i}, \dots, x_{pi})$) représente les variables prédicatrices de la donnée i. Soit une donnée u qui n'appartient pas à E et qui ne possède pas de label (u est uniquement caractérisée par un vecteur x_u de dimension p). Soit d'une fonction qui renvoie la distance entre la donnée u et une donnée quelconque appartenant à E. Soit un entier k inférieur ou égal à n. Voici le principe de l'algorithme de k plus proches voisins :

3.1 Sélectionne le nombre k de voisin

On fixe le nombre de k voisin, On détecte les k voisins les plus proches des nouvelles données l'entrée que l'on veut classe

3.2 On calcule les distances entre la donnée u et chaque donnée appartenant à E à l'aide de la fonction d

Il est possible d'utiliser différents types de distance : euclidienne, Manhattan, ...

3.2.1 Distance Euclidienne

- Sur une droite graduée :

Distance entre A et B est : $d(A,B) = |\text{abscisse de B} - \text{abscisse de A}|$

= la plus grande abscisse – la plus petite abscisse

On peut donc écrire une fonction distance, avec Python, qui prend en arguments d'entrée deux coordonnées x1 et x2 et qui renvoie la valeur de la distance entre les deux points de coordonnées x1 et x2

- Dans un plan :

On suppose que le plan est muni d'un repère orthonormal (O;I;J).

Soient A ($x_a ; y_a$) et B ($x_b ; y_b$) deux points du plan, la distance euclidienne entre A et B est :

$$d(A,B) = AB = \sqrt{(x_B - x_A)^2 + (y_B - y_A)^2}$$

3.2.2 Distance Manhattan

La distance de Manhattan a été utilisée dans une analyse de régression en 1757 par Roger

Joseph Boscovich. L'interprétation géométrique remonte à la fin du XIXe siècle et au développement de géométries non euclidiennes, notamment par Hermann Minkowski et son

inégalité de Minkowski, dont cette géométrie constitue un cas particulier, particulièrement utilisée dans la géométrie des nombres (Minkowski 1910).

La distance de Manhattan est appelée aussi taxi-distance, est la distance entre deux points Parcourue par un taxi lorsqu'il se déplace dans une ville où les rues sont agencées selon un réseau ou quadrillage. Un taxi-chemin est le trajet fait par un taxi lorsqu'il se déplace d'un nœud du réseau à un autre en utilisant les déplacements horizontaux et verticaux du réseau.

Entre deux points A et B, de coordonnées respectives (x_A, y_A) et (x_B, y_B)

La distance de Manhattan est définie par :

$$d(A, B) = |x_A - x_B| + |y_A - y_B|.$$

3.3 On retient les k données du jeu de données E les plus proches de u

3.4 On attribue à u la classe qui est la plus fréquente parmi les k données les plus proches

Suivant que l'on raisonne sur une, deux, trois dimensions, le calcul de la distance entre deux points est plus au moins simple.

Pour appliquer ce principe, il faudra :

- ✓ Évaluer la distance qui sépare le nouvel élément de chacun des autres points de l'ensemble E. Chaque point de l'ensemble est caractérisé par son indice i.
- ✓ Stocker ces valeurs de distance d dans une liste du type : $[[d, i], [\dots], \dots]$, où d est la distance qui sépare le nouvel élément du point d'indice i.
- ✓ Trier la liste selon les valeurs des distances d.
- ✓ Choisir les k premiers points de la liste triée qui sont donc les k-plus proches voisins.
- ✓ Assigner une classe au nouvel élément en fonction de la majorité des classes représentées parmi les k-plus proches voisins.

On aura donc besoin de trois fonctions :

- ✓ une fonction **distance** pour calculer la distance entre deux points de coordonnées connues.
- ✓ Une fonction **kvoisins** qui détermine les k-plus proches voisins d'un nouvel élément.
- ✓ Une fonction **predire_classe** qui détermine le résultat majoritaire des classes d'appartenance des k-plus proches voisins et assigne la classe du nouvel élément à cette classe majoritaire.

3.5 Ecriture algorithmique

On peut schématiser le fonctionnement de K-NN en l'écrivant en pseudo-code suivant :

Début Algorithme

Données en entrée :

- Un ensemble de données D .
- Une fonction de définition distance d .
- Un nombre entier K

Pour une nouvelle observation X dont on veut prédire sa variable de sortie y Faire :

1. Calculer toutes les distances de cette observation X avec les autres observations du jeu de données
2. Retenir les K observations du jeu de données D les proches de X en utilisation la fonction de calcul de distance d
3. Prendre les valeurs de y des K observations retenues :
 1. Si on effectue une régression, calculer la moyenne (ou la médiane) y de retenues
 2. Si on effectue une classification, calculer le mode de y retenues
4. Retourner la valeur calculée dans l'étape 3 comme étant la valeur qui a été prédite par K-NN pour l'observation X .

Fin Algorithme

Chapitre 4 : Etude d'un Exemple

4.1 Jeux de données

Nous avons choisi ici de nous baser sur le jeu de données « iris de Fisher ».

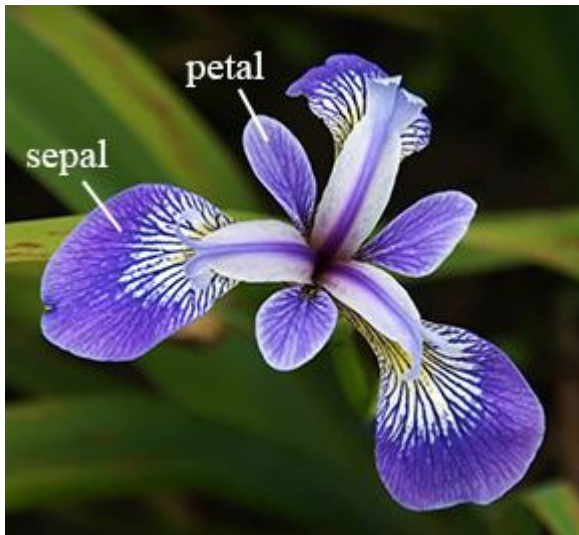


Photo 1 IRIS

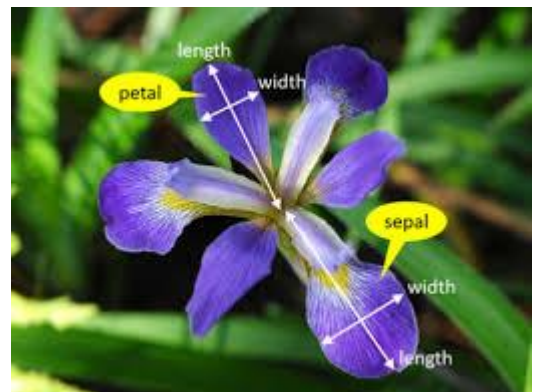


Photo 2 IRIS

Ce jeu de données est composé de 50 entrées, pour chaque entrée nous avons :

- ▷ la longueur des sépales (en cm)
- ▷ la largeur des sépales (en cm)
- ▷ la longueur des pétales (en cm)
- ▷ la largeur des pétales (en cm)
- ▷ l'espèce d'iris :
 - Iris setosa
 - Iris virginica
 - Iris versicolor

4.2 Les logiciels utilisés

Anaconda navigator
Jupyter notebook
Chrome

4.3 Bibliothèques Python utilisées

4.3.1 Pandas

Pandas est une excellente librairie pour importer vos tableaux Excel (et autres formats) dans Python dans le but de tirer des statistiques et de charger votre Dataset dans Sklearn.

4.3.2 NumPy

Numpy est la librairie qui permet de créer et manipuler des matrices ou tableau simplement et avec efficacité. Ainsi le calcul matriciel représente l'essentiel du Machine Learning. Il est important de le comprendre, mais les fonctions présentes dans Numpy font les calculs matriciels à notre place... Magique !

4.3.3 Matplotlib

Matplotlib est la librairie qui permet de visualiser nos Datasets, nos fonctions, nos résultats sous forme de graphes, courbes et nuages de points

4.3.4 Sklearn

Sklearn est la librairie qui contient toutes les fonctions de l'état de l'art du Machine Learning. On y trouve les algorithmes les plus importants ainsi que diverses fonctions de pre-processing.

4.4 Importation des bibliothèques

Importation des bibliothèques requies

```
Entrée [1]: 1 import pandas as pd
            2 import numpy as np
            3 import matplotlib.pyplot as plt
            4 from sklearn.neighbors import KNeighborsClassifier
```

4.5 Importation des données

Importation des données

```
Entrée [3]: 1 df = pd.read_csv(r"C:\Users\Kodjo\data_science\iris_min.csv")
```

```
Entrée [ ]: 1
```

4.6 Formatage

Ici on va remplacer

Iris-setosa par 0

Iris-virginica par 1

Iris-versicolor par 2

Formatage

```
Entrée [4]: 1 df['class'] = df['class'].replace("Iris-setosa", 0)
2 df['class'] = df['class'].replace("Iris-virginica", 1)
3 df['class'] = df['class'].replace("Iris-versicolor", 2)
```

```
Entrée [ ]: 1
```

4.7 Première visualisation des données

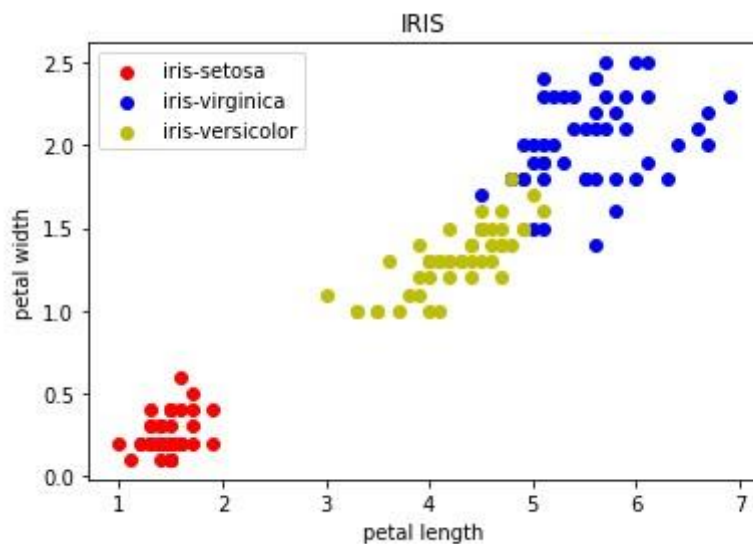
Une fois le fichier csv modifié, il est possible d'écrire un programme permettant de visualiser les données sous forme de graphique (abscisse : "petal_length", ordonnée : "petal_width")

Visualisation en nuage de point

```
Entrée [ ]: 1 x=df['petal_length']
2 y=df['petal_width']
3 label=df['class']
4 plt.scatter(x[label==0], y[label==0], c='r', label='iris-setosa')
5 plt.scatter(x[label==1], y[label==1], c='b', label='iris-virginica')
6 plt.scatter(x[label==2], y[label==2], c='y', label='iris-versicolor')
7 plt.xlabel("petal length")
8 plt.ylabel("petal width")
9 plt.title("IRIS")
10 plt.legend()
11 plt.show()
12
```

```
Entrée [ ]: 1
```

Graphe :



4.8 Echantillon a prédire

```
1 # Echantillon predire
```

Entrée [6]:

```
1 #Les attribut de echantillon
```

Entrée [7]:

```
1 #Valeur
2 longueur = 2.5
3 largeur = 0.75
4 k=3
```

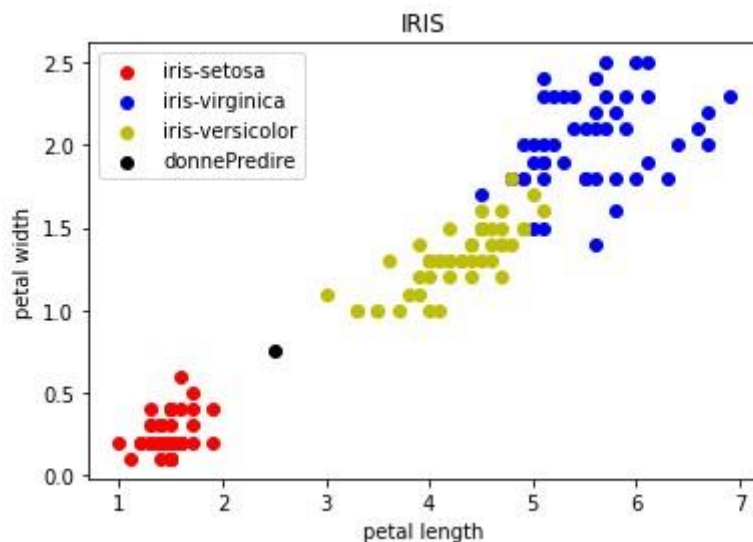
4.9 Visualisation de Echantillon a prédire dans le graphe

```
1 # Visualisation
```

Entrée [8]:

```
1 plt.scatter(x[label==0], y[label==0], c='r', label='iris-setosa')
2 plt.scatter(x[label==1], y[label==1], c='b', label='iris-virginica')
3 plt.scatter(x[label==2], y[label==2], c='y', label='iris-versicolor')
4 #On a joute au graphe
5 plt.scatter(longueur, largeur, color='k')
6 plt.xlabel("petal length")
7 plt.ylabel("petal width")
8 plt.title("IRIS")
9 plt.legend()
10 plt.show()
```

Graphe



4.10 Séparation du jeux de donnée en training set et testing set

Séparation du jeux de donnée en training set et testing set

```
Entrée [25]: 1 #feature
              2 feature_columns = ['petal_length', 'petal_width']
              3 X = df[feature_columns].values

Entrée [26]: 1 #target
              2 label = df['class']

Entrée [27]: 1 from sklearn.model_selection import train_test_split
              2 X_train, X_test, y_train, y_test = train_test_split(X,label, train_size=0.8, random_state=0)

Entrée [ ]: 1
```

4.11 Algorithmique KNN pour K=2

Utilisation d'algorithme KNN pour K=2 et le taux d'erreur de prédiction

Algorithme KNN

```
Entrée [28]: 1 #algo KNN
              2 k=2
              3 #Creation de model de classification
              4 model = KNeighborsClassifier(n_neighbors=k)
              5 #(data, target)
              6 model.fit(X_train,y_train)
              7 model.score(X_test, y_test)
```

Out[28]: 0.9333333333333333

4.12 Optimisation du taux d'erreurs

Optimisation du score sur les données test

```
Entrée [30]: 1 errors = {}  
2 valeur_k = list(range(2,15))
```

```
Entrée [31]: 1 for k in valeur_k:  
2     model = KNeighborsClassifier(n_neighbors=k)  
3     model.fit(X_train, y_train)  
4     score = 100*(1-model.score(X_test, y_test))  
5     error = {k:score}  
6     errors.update(error)
```

```
Entrée [32]: 1 errors
```

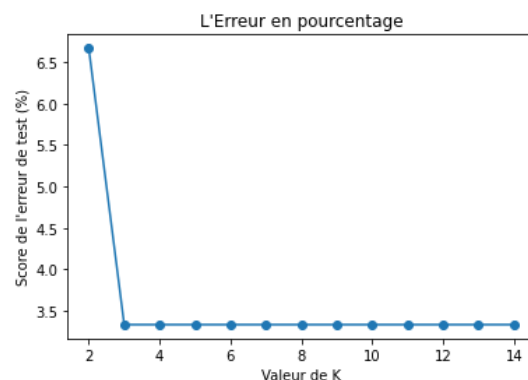
```
Out[32]: {2: 6.666666666666665,  
3: 3.333333333333326,  
4: 3.333333333333326,  
5: 3.333333333333326,  
6: 3.333333333333326,  
7: 3.333333333333326,  
8: 3.333333333333326,  
9: 3.333333333333326,  
10: 3.333333333333326,  
11: 3.333333333333326,  
12: 3.333333333333326,  
13: 3.333333333333326,  
14: 3.333333333333326}
```

- Sous Forme Graphique

Visualiser du Score d'erreur sous forme graphique

```
Entrée [34]: 1 plt.plot( valeur_k, list(errors.values()), 'o-' )  
2 plt.title("L'Erreur en pourcentage")  
3 plt.xlabel("Valeur de K")  
4 plt.ylabel("Score de l'erreur de test (%)")
```

```
Out[34]: Text(0, 0.5, "Score de l'erreur de test (%)")
```



4.13 Création d'une fonction python résultat qui retourne la prédiction

C'est la classe dans lequel échantillon a appartient.

```
Entrée [ ]: 1 #Fonction resltat
2 def resultat():
3     if prediction == 0 :
4         return "Resultat : setosa"
5     if prediction == 1 :
6         return "Resultat : virginica"
7     if prediction == 2 :
8         return "Resultat : verginicolor"
```

4.14 Visualisation finale et prédiction

```
Entrée [ ]: 1 #Donnée d'entr"
2 longueur = 2.5
3 largeur = 0.75

Entrée [36]: 1 #L'algorithmique est deja importée
2 d = list(zip(x,y))
3 model = KNeighborsClassifier(n_neighbors=3)
4 model.fit(d, label)
5 prediction = model.predict([[longueur, largeur]])

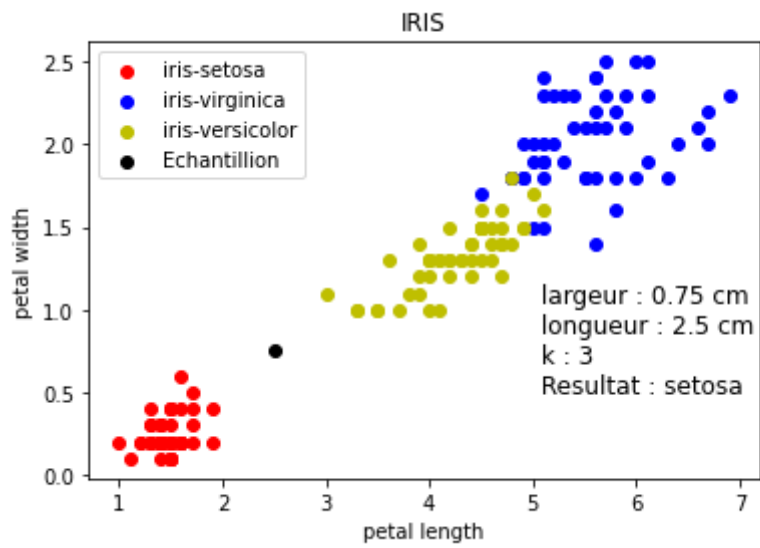
Entrée [40]: 1 #Fonction resltat
2 def resultat():
3     if prediction == 0 :
4         return "Resultat : setosa"
5     if prediction == 1 :
6         return "Resultat : virginica"
7     if prediction == 2 :
8         return "Resultat : verginicolor"
9 resultat()

Out[40]: 'Resultat : setosa'
```

Visualisation dans

```
Entrée [14]: 1 plt.scatter(x[label==0], y[label==0], c='r', label='iris-setosa')
2 plt.scatter(x[label==1], y[label==1], c='b', label='iris-virginica')
3 plt.scatter(x[label==2], y[label==2], c='y', label='iris-versicolor')
4 plt.scatter(longueur, largeur, color='k', label="Echantillon")
5 res = resultat()
6 plt.text(5,0.5, f" largeur : {largeur} cm\n longueur : {longueur} cm\n k : {k}\n {res}", fontsize=12)
7 plt.xlabel("petal length")
8 plt.ylabel("petal width")
9 plt.title("IRIS")
10 plt.legend()
11 plt.show()
12
```


Graphe : Le résultat de prédiction sur le graphe



5 Conclusion

Dans ce rapport, vous avez découvert l'algorithme k-NN qui est l'un des algorithmes de l'apprentissage automatique supervisé il est simple et facile à mettre en œuvre. il n'a aucune hypothèse sur les données (linéaires, affines...) en plus de cela il est polyvalent. Il peut être utilisé pour la classification, la régression.

Vous avez appris également que :

- K-NN stocke tout le jeu de données pour effectuer une prédiction,
- K-NN ne calcule aucun modèle prédictif et il rentre dans le cadre du Lazy Learning