



Royaume du Maroc
Université Mohammed Premier
École Supérieure De Technologie-Oujda
Département : Informatique
Filière : LP Informatique Décisionnelle

Analyse De Données

Analyse et Classification des Emails : Détection de diabétiques / non-diabétiques

Réalisé par : HNIOUA Abdessamad

Encadré par : M. Mounir Grari

Année Universitaire : 2023/2024

La classification supervisée : est souvent un choix approprié pour la détection de diabétiques en raison de la disponibilité des données étiquetées, de sa performance élevée, de son interprétabilité, de sa flexibilité et de son évolutivité. Cependant, il est important de choisir le bon algorithme et de bien prétraiter les données pour obtenir les meilleurs résultats possibles.

1. Dataset 'diabétiques':

Le jeu de données décrit est conçu pour l'analyse et la prédiction du diabète. Il contient plusieurs caractéristiques cliniques pertinentes des patients. Voici une brève description des caractéristiques incluses :

Grossesses : Indique le nombre de fois que la patiente a été enceinte.

Glucose : Niveau de glucose dans le sang après un jeûne de 2 heures, mesuré en mg/dL. Un indicateur crucial pour diagnostiquer le diabète.

Pression Artérielle : Mesure de la pression sanguine (mm Hg). La pression artérielle élevée peut être un indicateur de risque accru de diabète.

Épaisseur de la Peau : Épaisseur du pli cutané du triceps (mm), un indicateur de résistance à l'insuline.

Insuline : Niveau d'insuline dans le sang 2 heures après un repas, mesuré en UI/mL. Les niveaux bas peuvent indiquer le diabète.

IMC (Indice de Masse Corporelle) : Le poids en kilogrammes divisé par le carré de la taille en mètres. L'IMC élevé peut indiquer un risque de diabète.

Fonction de Prédiction du Diabète : Une fonction qui fournit une prédiction de la probabilité du diabète basée sur l'histoire familiale.

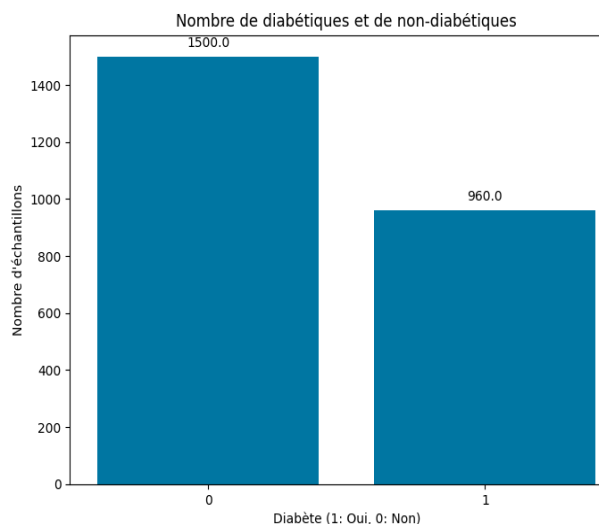
Âge : Âge du patient en années.

L'Issue (résultat) est la variable cible, indiquant si le patient a été diagnostiqué avec le diabète (1 pour oui, 0 pour non).

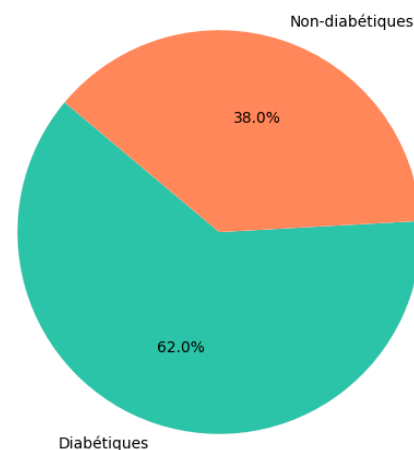


Cela crée un nouveau dataframe X qui contient uniquement les caractéristiques (facteurs).

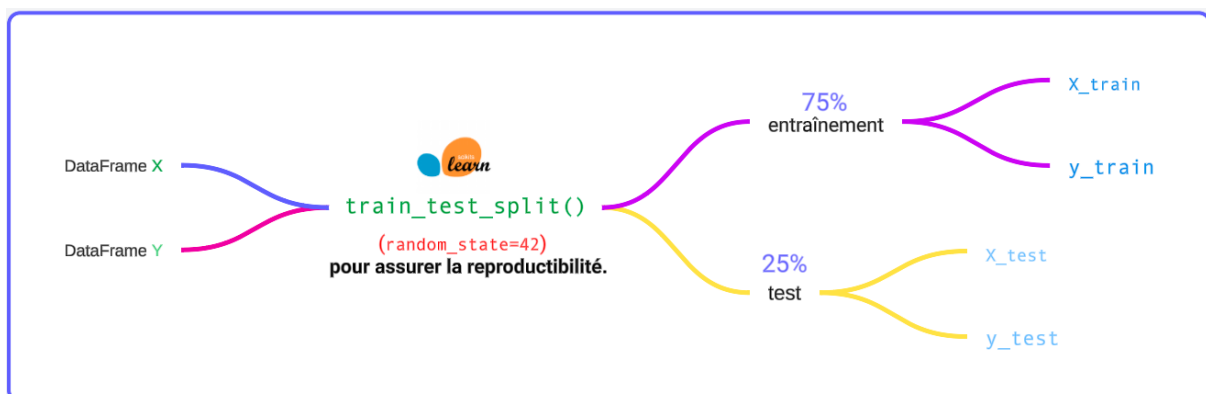
Nous attribuons la colonne 'Outcome' du dataframe data à la variable y, créant ainsi un dataframe y qui contient la cible.



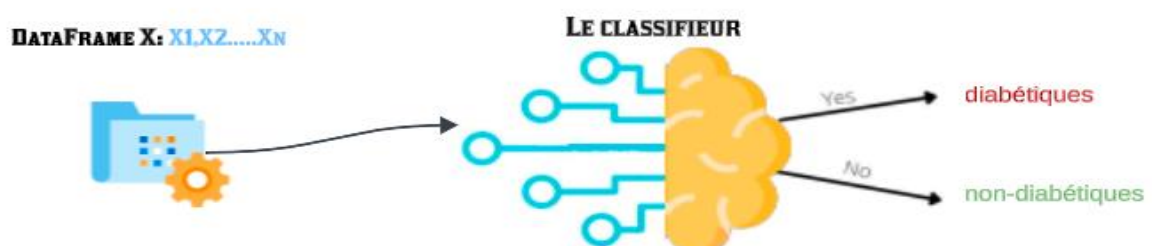
Pourcentage de diabétiques et de non-diabétiques "Train"



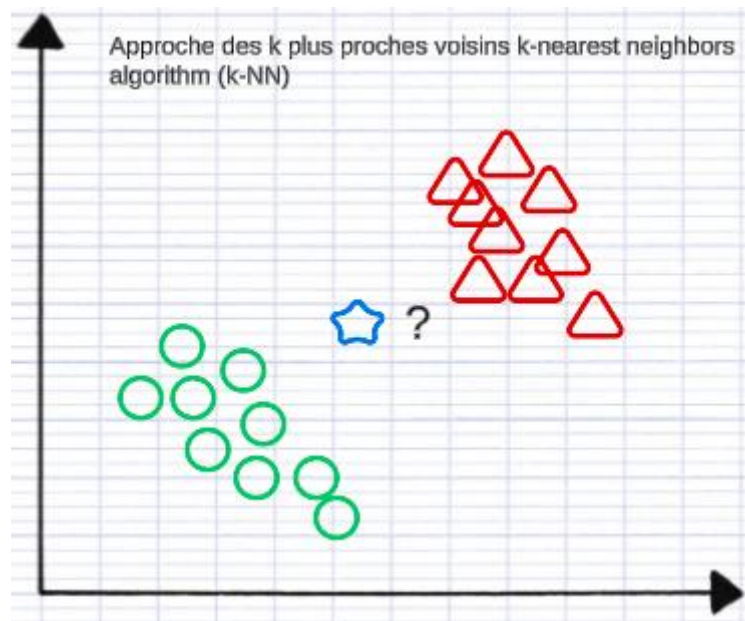
La visualisation représente un histogramme des données de diabétiques, distinguant les échantillons de courrier électronique en fonction de s'ils sont classés comme diabétiques (1) ou non-diabétiques (0). La colonne des courriels non-diabétiques est plus dominante, avec une valeur de 2788 échantillons, tandis que la colonne des courriels diabétiques présente une valeur de 1812 échantillons. Cela suggère que la majorité des échantillons de courrier électronique dans l'ensemble de données ne sont pas classés comme de diabétiques.



Le schéma illustre le processus de division des données en ensembles d'entraînement et de test à l'aide de la fonction `train_test_split` de scikit-learn. Dans cette configuration, 75% des données sont réservées pour l'ensemble d'entraînement et 25% pour l'ensemble de test. La graine aléatoire est fixée à 42 pour assurer que la division des données est reproductible.



I. Approche des k plus proches voisins k-nearest neighbors algorithm (k-NN)



L'algorithme des k plus proches voisins (k-nearest neighbors ou k-NN) est une méthode de classification simple mais puissante. Dans le contexte d'un projet de filtrage de spam, cette approche fonctionne en calculant la similarité entre les e-mails en se basant sur des caractéristiques sélectionnées, puis en attribuant une étiquette (diabétiques ou non-diabétiques) en fonction du vote majoritaire parmi les k voisins les plus proches. Bien que facile à comprendre et à mettre en œuvre, il est essentiel de choisir judicieusement le nombre k et **la mesure de distance** pour obtenir des performances optimales.

1. Sélectionne le modèle d'estimateur des k plus proches voisins avec k=3 en utilisant la classe KNeighborsClassifier de scikit-learn.

```
from sklearn.neighbors import KNeighborsClassifier
clf = KNeighborsClassifier(n_neighbors=3)
```

2. Entraîne le modèle des k plus proches voisins avec k=3 sur les données fournies dans les DataFrames X_train (caractéristiques) et y_train (étiquettes de classe).

```
clf.fit(X_train, y_train)
```

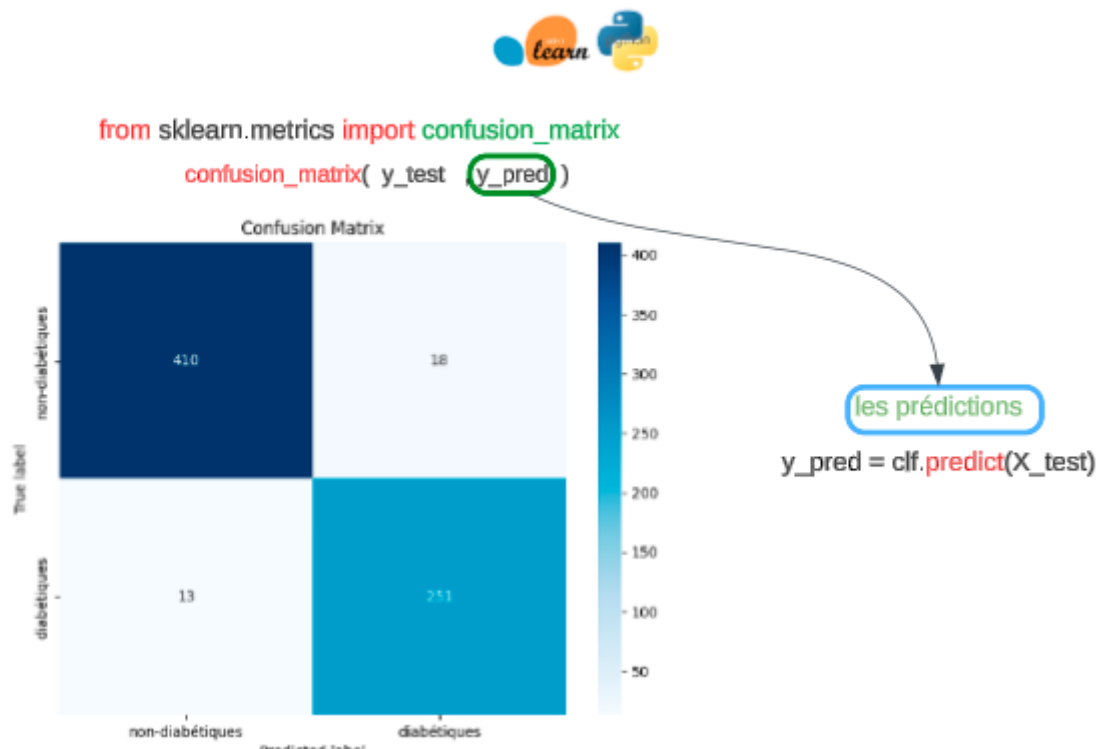
3. Le modèle a été évalué sur les données d'entraînement, obtenant un score de 90% de précision. Cela signifie que le modèle prédit correctement la classe des échantillons dans 90% des cas.

```
print(clf.score(X_train, y_train))
```

4. Utilise le modèle entraîné pour effectuer des prédictions sur les données de test X_test, stockant les prédictions dans la variable pred.

```
pred = clf.predict(X_test)
```

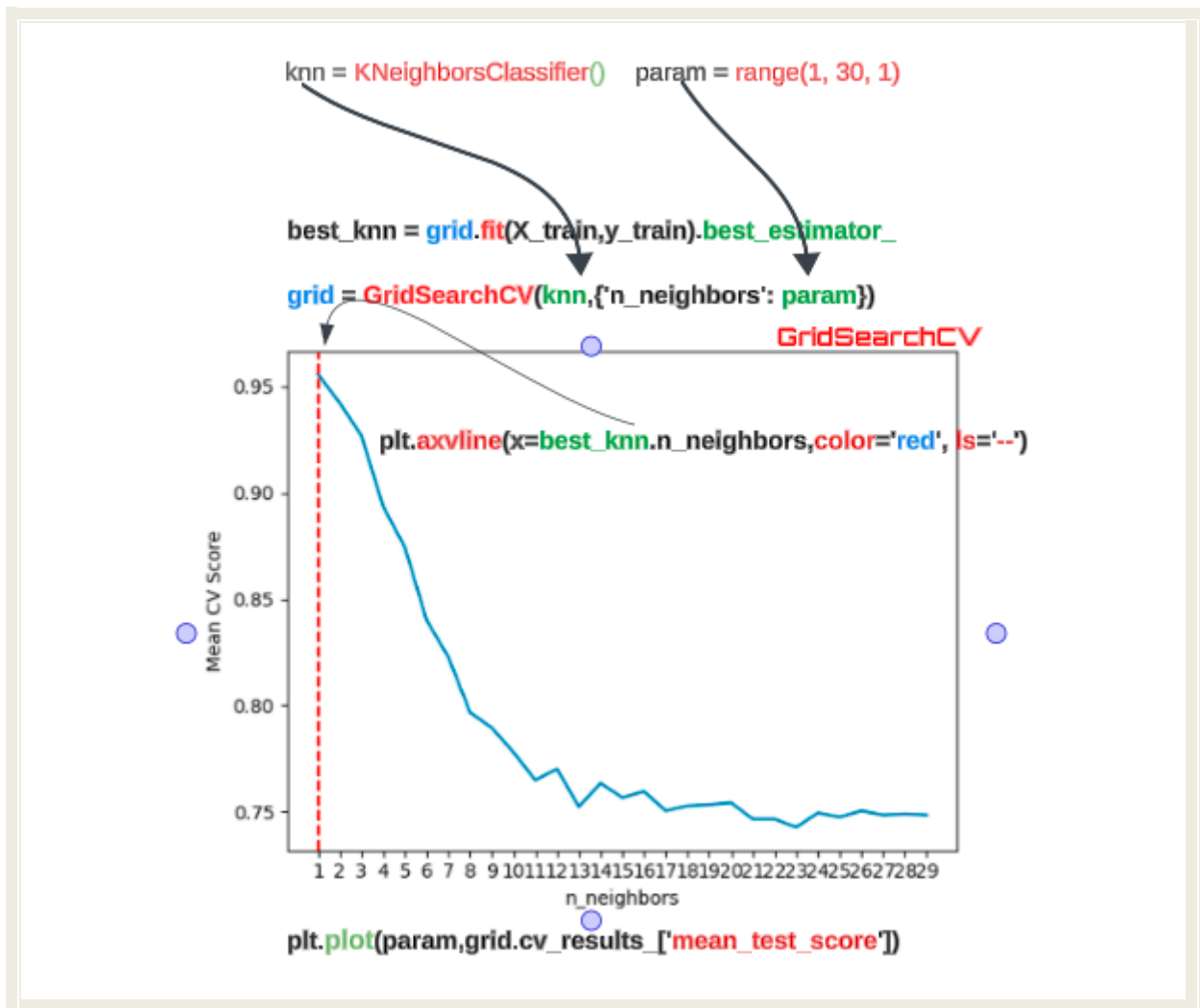
La matrice de confusion : résume la performance d'un modèle de classification en révélant ses prédictions correctes et incorrectes par rapport aux vraies étiquettes de classe, essentiellement en identifiant les vrais positifs, les vrais négatifs, les faux positifs et les faux négatifs.



L'accuracy évalue la proportion de prédictions correctes du modèle de classification sur l'ensemble des échantillons, ce qui est essentiel pour mesurer sa capacité à distinguer efficacement les diabétiques des non-diabétiques.



GridSearchCV de scikit-learn pour rechercher la meilleure valeur du paramètre `n_neighbors` (nombre de voisins) pour le modèle kNN. Il configure une grille de valeurs de `n_neighbors` à tester, puis utilise GridSearchCV pour évaluer les performances du modèle kNN avec chaque valeur de `n_neighbors` et sélectionner celle qui donne les meilleures performances sur les données d'entraînement à l'aide d'une validation croisée. Enfin, il renvoie les paramètres du modèle kNN optimal.



La partie d'optimisation par GridSearchCV a permis de déterminer que le meilleur nombre de voisins à utiliser dans l'algorithme kNN est 1, en se basant sur les performances du modèle évalué sur les données d'entraînement à l'aide d'une validation croisée. Cette démarche a abouti à la conclusion que le modèle kNN avec 1 voisins offre la meilleure performance parmi les valeurs de `n_neighbors` testées, ce qui permet d'optimiser la précision du modèle pour la classification des données.

```
y_pred = best_knn.predict(X_test)
from sklearn.metrics import accuracy_score
# Calcul de l'accuracy
accuracy = accuracy_score(y_pred, y_test)

# Affichage de l'accuracy
print("Accuracy:", accuracy*100,"%")
```

Après l'exécution de l'algorithme d'optimisation GridSearchCV, l'accuracy du modèle a légèrement augmenté, passant de 95% à 96%. Cette amélioration peut être attribuée à la recherche du meilleur paramètre `n_neighbors` effectuée par GridSearchCV, ce qui a permis de sélectionner une configuration plus adaptée du modèle kNN pour les données d'entraînement.