

Introduction à l'algorithme KNN et son utilisation avec SQL

1.1 1. Principe de l'algorithme KNN (K-Nearest Neighbors)

L'algorithme KNN (K plus proches voisins) est une méthode de classification ou de recherche basée sur la **similarité**. Le principe est simple :

Pour classer ou analyser un objet, on cherche les K objets les plus proches dans notre base de données, puis on utilise leurs caractéristiques pour prendre une décision.

1.1.1 Exemple concret : Classification d'exoplanètes

Imaginons que nous voulons déterminer si une nouvelle planète découverte est plutôt de type "terrestre" (comme la Terre) ou "gazeuse" (comme Jupiter).

Démarche KNN : 1. On calcule la distance entre cette nouvelle planète et toutes les planètes connues 2. On sélectionne les K planètes les plus proches (par exemple K=5) 3. On regarde à quel type appartiennent ces 5 voisins (si les objets appartiennent à un certain type) 4. On attribue le type majoritaire à notre nouvelle planète

1.2 2. Calcul de la distance

La distance utilisée est généralement la **distance euclidienne** dans un espace à plusieurs dimensions.

Pour une exoplanète caractérisée par sa **masse** (mass) et son **rayon** (radius), la distance à Jupiter est :

$$distance = \sqrt{(mass - mass_Jupiter)^2 + (radius - radius_Jupiter)^2}$$

1.2.1 Visualisation

Sur un graphique avec la masse en abscisse et le rayon en ordonnée : - Chaque planète est un point - Jupiter est notre point de référence - La distance KNN correspond au rayon d'un cercle centré sur Jupiter - Les planètes à l'intérieur du cercle sont les "voisins proches"

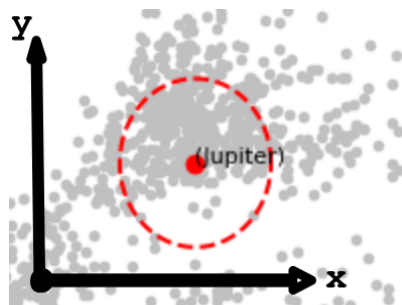


FIGURE 1 – Jupiter dans un diagramme en nuage de points masse-rayon

1.3 3. Implémentation en SQL

1.3.1 A. Création de la fonction de distance

SQLite ne possède pas nativement la fonction `power()` ni toutes les fonctions mathématiques. On doit donc les créer en Python :

```

1 import sqlite3
2 import math
3
4 # Connexion à la base de données
5 conn = sqlite3.connect('exoplanetes.db')
6
7 # Création de la fonction power
8 conn.create_function("power", 2, lambda x, y: x**y if x != None and y
    != None else None)
9
10 # Création de la fonction sqrt
11 conn.create_function("sqrt", 1, lambda x: math.sqrt(x) if x != None and
    x >= 0 else None)

```

1.3.2 B. Calcul de la distance dans une requête SQL

Une fois les fonctions créées, on peut calculer directement la distance dans nos requêtes :

```

1 SELECT
2     name,
3     mass,
4     radius,
5     sqrt(power(mass - 1.0, 2) + power(radius - 1.0, 2)) as
        distance_jupiter
6 FROM exoplanets
7 WHERE mass IS NOT NULL
8     AND radius IS NOT NULL
9 ORDER BY distance_jupiter
10 LIMIT 5;

```

QUESTION 1 : Que retourne la requête sql ci-dessus ?

1.3.3 C. Utilisation d'une fonction personnalisée

Pour simplifier et rendre le code plus lisible, on peut créer une fonction SQL dédiée :

```

1 def distance_jupiter_sql(mass, radius):
2     """Calcule la distance euclidienne à Jupiter"""
3     if mass is None or radius is None:
4         return None
5
6     mass_J = 1.0    # Masse de Jupiter (en masses joviennes)
7     radius_J = 1.0  # Rayon de Jupiter (en rayons joviens)
8
9     import math
10    distance = math.sqrt((mass - mass_J)**2 + (radius - radius_J)**2)
11
12    return distance
13
14 # Enregistrement dans SQLite

```

```
15 conn.create_function("distance_jupiter", 2, distance_jupiter_sql)
```

Utilisation simplifiée :

```
1 SELECT
2     name,
3     mass,
4     radius,
5     distance_jupiter(mass, radius) as distance_jupiter
6 FROM exoplanets
7 WHERE distance_jupiter(mass, radius) IS NOT NULL
8 ORDER BY distance_jupiter
9 LIMIT 5;
```

1.3.4 D. Mise à jour de la base avec les distances

On peut aussi ajouter une colonne dans la table pour stocker les distances :

```
1 # Ajout de la colonne
2 cursor.execute("ALTER TABLE exoplanets ADD COLUMN distance_jupiter REAL
3 ")
4 # Calcul et mise à jour
5 cursor.execute("""
6     UPDATE exoplanets
7     SET distance_jupiter = distance_jupiter(mass, radius)
8 """)
9
10 conn.commit()
```

1.4 4. Exercice pratique

QUESTION 2 : Comment trouver les 10 exoplanètes les plus similaires à la Terre ? Compléter la série d'instruction ci-dessous

Données de référence

Important : Dans notre base de données, les masses et rayons sont exprimés **relativement à Jupiter**.

- Masse de la Terre : **0.00315** masses joviennes (la Terre est ~318 fois moins massive que Jupiter)
- Rayon de la Terre : **0.0892** rayons joviens (le rayon de la Terre est ~11.2 fois plus petit que celui de Jupiter)

Étape 1 : Créer la fonction `distance_terre`

En vous inspirant de la fonction `distance_jupiter_sql` fournie ci-dessus, complétez le code suivant :

```
1 def distance_terre_sql(mass, radius):
2     """Calcule la distance euclidienne à la Terre"""
3     if mass is None or radius is None:
4         return None
5
```

```

6     mass_T = _____ # À compléter : masse de la Terre en masses
    joviennes
7     radius_T = _____ # À compléter : rayon de la Terre en rayons
    joviens
8
9     import math
10    distance = math.sqrt((mass - _____)**2 + (radius - _____)**2)
11
12    return distance
13
14 # Enregistrement dans SQLite
15 conn.create_function("_____", 2, distance_terre_sql)

```

Étape 2 : Requête SQL

Complétez la requête SQL pour obtenir les 10 planètes les plus proches de la Terre :

```

1 SELECT
2     name,
3     mass,
4     radius,
5     _____(mass, radius) as distance_terre
6 FROM exoplanets
7 WHERE _____(mass, radius) IS NOT NULL
8 ORDER BY _____
9 LIMIT _____;

```

QUESTION 3 : Compléter le script pour mettre à jour de la base.

Ajoutez une colonne distance_terre dans la table et calculez les distances :

```

1 # Ajout de la colonne
2 cursor.execute("ALTER TABLE exoplanets ADD COLUMN _____ REAL")
3
4 # Calcul et mise à jour
5 cursor.execute("""
6     UPDATE exoplanets
7     SET _____ = _____(mass, radius)
8 """)
9
10 conn.commit()

```

1.5 5. Visualisation

La fonction `coordonnees_cercle_xlog_y` retourne les listes `x,y` de coordonnées des points du cercle centré sur `X,Y`, de rayon `r`. Cette fonction a pour paramètres `X, Y, r`.

Compléter le code de visualisation pour afficher la Terre et son cercle de voisinage (rayon = 0.3). Le point représentant la Terre sera de couleur bleue :

```

1 # Données Terre (en unités joviennes)
2 mass_T = 0.00315
3 radius_T = 0.0892

```

```

4
5 # Position de la Terre sur le graphique
6 plt.plot(mass_T, radius_T, marker='o', color='_____' )
7 axes.text(mass_T, radius_T, f"(Terre)", fontsize=8)
8
9 # Cercle autour de la Terre (à compléter)
10 x_circle_T, y_circle_T = coordonnees_cercle_xlog_y(_____, _____,
11 _____)
12 plt.plot(x_circle_T, y_circle_T, color='gold', linewidth=1.5)

```

QUESTION 4 :

1. Pourquoi est-il important d'utiliser les mêmes unités (masses et rayons joviens) pour toutes les planètes ?
2. Observez les valeurs de masse et rayon de la Terre par rapport à Jupiter. Que remarquez-vous ? Ces valeurs sont-elles du même ordre de grandeur ?
3. Les 10 planètes les plus proches de la Terre sont-elles toutes de type "terrestre" ?

Note : Dans le TP à venir, nous verrons comment adapter cette distance lorsque les données ont des échelles très différentes (par exemple avec des échelles logarithmiques).

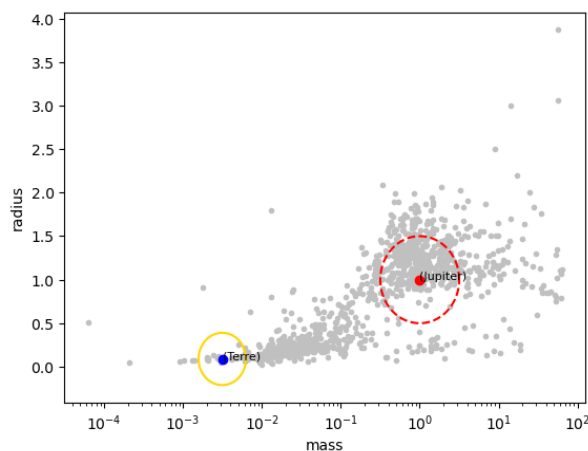


FIGURE 2 – représentation graphique avec axe des abscisses logarithmique