

Exercice 1 (débranché) : Classification k NN

Données pour l'apprentissage (volontairement réduite pour procéder à la main)

numero_departement	nom_departement	prefecture	nom_commune	codes_postaux	latitude(°)	longitude(°)
77	Seine-et-Marne	Melun	Fontainebleau	77300	48,4	2,7
77	Seine-et-Marne	Melun	Melun	77000	48,533333	2,666667
74	Haute-Savoie	Annecy	Annecy	74000	45,9	6,116667
91	Essonne	Evry	Evry	91000	48,633333	2,45

Les coordonnées GPS de Dammarie-les-Lys sont :

Latitude : 48,515088°

Longitude : 2,634702°

On rappelle que le rayon moyen de la terre mesure environ 6371 kilomètres et que la longueur d'un arc de cercle est égale au produit du rayon par l'angle au centre.

1) Distance Dammarie-Annecy

- Calculer l'écart en kilomètres lié à la différence de latitude entre Dammarie-les-Lys et Annecy.
- Calculer l'écart en kilomètres lié à la différence de longitude entre Dammarie-les-Lys et Annecy.
- En utilisant le théorème de Pythagore, en déduire une approximation de la distance entre Dammarie les lys et Annecy.

2) Calculer les distances entre Dammarie-les-Lys et chacune des villes apparaissant dans les données d'apprentissage.

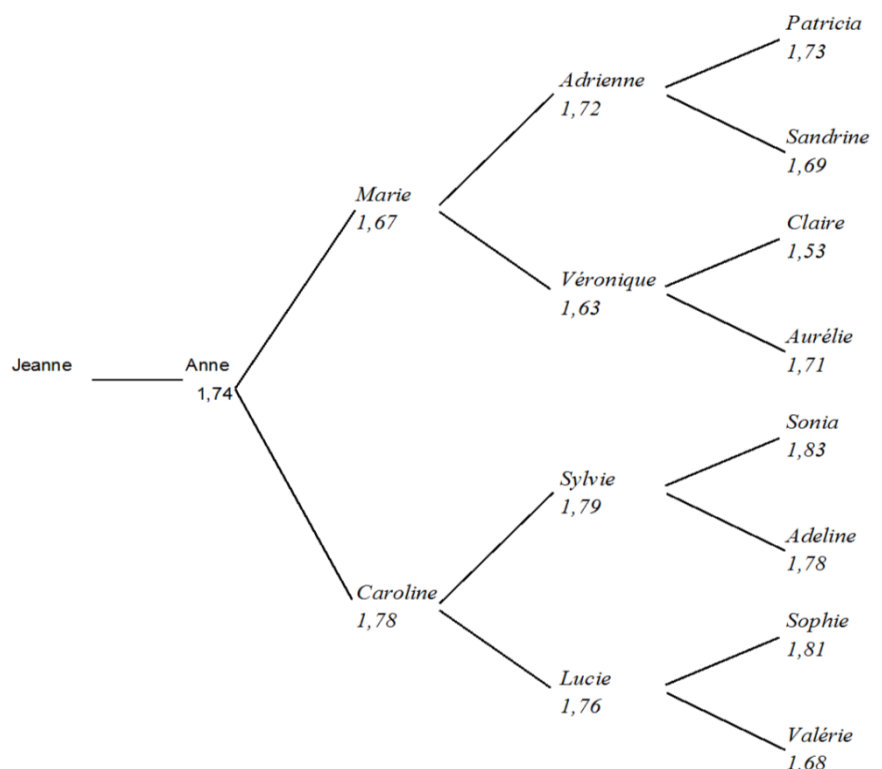
3) A quel département la ville de Dammarie-les-Lys appartient-elle selon un algorithme basé sur les 3 plus proches voisins s'appuyant sur les données d'apprentissage ci-dessus ? (détailler la démarche pour justifier)

Exercice 2 (débranché) : Régression k NN

L'arbre généalogique de Jeanne est donné ci-contre avec les tailles à l'âge adulte.

Selon le principe de l'exemple 2 du cours, quelle estimation de la taille de Jeanne à l'âge adulte est proposée par une régression à partir :

- d'un algorithme 1 NN ?
- d'un algorithme 3 NN ?
- d'un algorithme 7 NN ?
- d'un algorithme 5 NN ? (avec la convention qu'en cas d'égalité des distances, le premier voisin trouvé est conservé)



Exercice 3 (débranché) : formalisation de l'algorithme k NN

On dispose :

- de données d'apprentissages présentées sous la forme d'une liste de couples (prédicteurs, étiquette)
- d'une fonction distance qui à partir des prédicteurs de deux objets renvoie la distance entre les deux objets avec les propriétés mathématiques usuelles pour une distance

Ecrire en langage naturel, en pseudo-code ou en Python une fonction kNN qui :

- prend en entrées : un objet observé pour lequel on dispose des prédicteurs, une liste de données d'apprentissage et entier naturel k non nul et inférieur ou égal à la longueur de la liste des données d'apprentissage
- renvoie en sortie la liste des étiquettes des k plus proches voisins de l'objet observé parmi les données d'apprentissage.

En cas d'égalité de distance on prendra arbitrairement le premier rencontré.

Au cas où les prédicteurs d'un élément des données d'apprentissage correspondraient exactement à ceux de l'objet observé, on retournera une liste contenant k fois l'étiquette de cet élément.

Exercice 4 (TP informatique) : Classification k NN

Le fichier « `apprentissage_departement.csv` » a été constitué à partir de données extraites du site :

<https://www.data.gouv.fr/fr/datasets/listes-des-communes-geolocalisees-par-regions-departements-circonscriptions-nd/>

- 1) Ecrire une fonction en Python qui lit ce fichier et constitue une liste de données d'apprentissage exploitable par la fonction kNN de l'exercice 3, constitué de couples ((latitude, longitude), département) (deux prédicteurs et une étiquette). On pourra s'appuyer sur la bibliothèque « `csv` » et les dictionnaires.
- 2) En s'inspirant de l'exercice 1, écrire une fonction `distance` calculant la distance entre deux villes à partir des couples (latitude, longitude) de chacune de ces deux villes. On rappelle que la fonction « `sqrt` » de la bibliothèque « `math` » permet de calculer une racine carrée.
- 3) A l'aide des fonctions de l'activité préparatoire et de l'exercice 3, écrire une fonction prenant en entrée un entier non nul k et le couple (latitude, longitude) des coordonnées GPS d'une commune et effectue une prédiction sur le département auquel appartient cette commune à l'aide d'un algorithme basé sur les k plus proches voisins figurant dans les données du fichier « `apprentissage_departement.csv` ».
- 4) Tester cette fonction sur plusieurs communes de votre choix (les coordonnées GPS pourront être cherchées sur internet). On effectuera des tests avec des communes qui apparaissent dans le fichier « `apprentissage_departement.csv` » et d'autres qui n'y apparaissent pas. Les tests seront effectués avec $k = 7$ et $k = 6500$.

Exercice 5 (TP informatique) : Régression k NN

Le fichier « apprentissage_temperatures2016.csv » a été constitué à partir de données extraites du site : <https://www.data.gouv.fr/fr/datasets/temperature-quotidienne-regionale-depuis-janvier-2016/>

On souhaite estimer la température moyenne pour une région et un jour donné.

- 1) Ecrire une fonction en Python qui lit ce fichier et constitue une liste de données d'apprentissage exploitable par la fonction `kNN` de l'exercice 3, constitué de couples ((région, jour), température moyenne) où jour désigne le numéro du jour dans l'année (compris entre 1 et 366).
- 2) Définir une fonction `distance` entre des prédicteurs donnés sous la forme d'un couple (région, jour), qui vaut 400 (ou n'importe quelle valeur strictement supérieure à 366) si les régions sont différentes et qui vaut le nombre de jours d'écart entre les jours de l'année sinon (peu importe l'année, attention l'écart entre le 30 décembre et le 2 janvier vaut donc 4).
- 3) A l'aide des fonctions des exercices 0 (activité préparatoire) et 3, écrire une fonction prenant en entrée un entier non nul `k` et un couple (région, jour) qui renvoie une estimation de la température qu'il fera dans la région pour ce jour par un algorithme k NN.
- 4) Tester l'algorithme sur les jours de votre choix et pour différentes valeurs de `k`. Que peut-on penser des choix `k = 366` ? `k = 4000` ?
- 5) (bonus) Ecrire des fonctions permettant d'obtenir le numéro du jour dans l'année à partir de la date ou vice versa, puis les utiliser pour améliorer le programme précédent en travaillant à partir des dates.

1NSI	Algorithmique	TP - Algorithme des k plus proches voisins
	k plus proches voisins	

REPONSES AUX QUESTIONS DU TP :

Q1 :

- a) différence de latitude : ≈ 291 km b) différence de longitude : ≈ 387 km c) distance ≈ 484 km
- Dammarie-les-Lys \leftrightarrow Fontainebleau : 14,7 km ; Dammarie-les-Lys \leftrightarrow Melun : 4,09 km
Dammarie-les-Lys \leftrightarrow Evry : 24,4 km
- Les trois plus proches voisins sont Melun, Fontainebleau et Evry. Le département le plus fréquent sur ces 3 commune est « Seine et Marne », donc d'après un algorithme 3-NN, Dammarie-les-Lys est bien en Seine et Marne.

Q2 :

- D'après un algorithme 1-NN, Jeanne fera la taille de sa mère (1 plus proche voisine), donc **1,74 m**.
- D'après un algorithme 3-NN, Jeanne fera la taille de la moyenne de sa mère et ses deux grand-mères (3 plus proches voisines), donc $(1,74 + 1,67 + 1,78)/3 = \mathbf{1,73\ m}$.
- D'après un algorithme 7-NN, Jeanne fera la taille de la moyenne de sa mère, ses deux grand-mères et ses 4 arrières grand-mères (7 plus proches voisines), donc **1,727 m**.
- D'après un algorithme 5-NN, Jeanne fera la taille de la moyenne de sa mère, ses deux grand-mères et ses 2 arrières grand-mères maternelles (du côté de Marie), donc **1,708 m**.

Q3 :

```

FONCTION kNN(k, observation, apprentissage) :
k (entier) : nombre de plus proches voisins à prendre en compte
observation : prédicteurs de l'objet dont on cherche l'étiquette son type exact dépend du
              problème
apprentissage : une liste de couple au format (predicteurs, étiquette)
DEBUT
    voisins ← liste vide
    POUR i VARIANT DE 0 A k-1 FAIRE :
        ajouter (distance(observation, apprentissage[i][predicteurs]),
                apprentissage[i][étiquette]) à voisins
    FIN POUR
    trier(voisins)
    POUR i VARIANT DE k A longueur(apprentissage) FAIRE :
        dist ← distance(observation, apprentissage[i][predicteurs])
        SI dist < voisins[k-1][0]:
            POUR j VARIANT DE 0 A k :
                SI dist < voisins[j][0]:
                    insérer (dist, apprentissage[i][1]) à la position j de voisins
                    supprimer dernier élément de voisins
                SORTIR DE LA BOUCLE POUR
            FIN SI
        FIN POUR
    FIN SI
    retour ← liste vide
    POUR i VARIANT DE 0 A k-1 :
        ajouter voisin[i][1] à retour
    FIN POUR
    RENVOYER retour
FIN

```

Q4:

```

1 import csv #pour lire le fichier de données d'apprentissage
2 import math #pour les calculs de distances
3
4 ## variables globales et constantes
5 RAYON=6371
6
7 #####
8 # QUESTION 1 #
9 #####
10 def extraireDonnees(nom) :
11     """gène des données d'apprentissage sous la forme d'une liste de couple au
12     format objet=(prédicteurs, étiquette) les données sont extraite du fichier
13     .csv dont les champs sont numero_departement; nom_departement; prefecture;
14     nom_commune; codes_postaux; latitude; longitude"""
15     apprentissage=[]
16     with open(nom+'.csv', 'r', encoding='utf-8') as fichier:
17         for donnees in csv.DictReader(fichier,delimiter=';') :
18             d=dict(donnees)
19             apprentissage.append(((float(d['latitude']),float(d['longitude'])),
20                                 d['nom_departement'], d['nom_commune'] ) )
21     return apprentissage
22
23 #####
24 # QUESTION 2 #
25 #####
26 def distance(v1, v2):
27     d_latitude = (v1[0]-v2[0])*math.pi/180 * RAYON
28     d_longitude = (v1[1]-v2[1])*math.pi/180 * RAYON
29     return math.sqrt(d_latitude*d_latitude + d_longitude*d_longitude)
30
31 distance((46.2, 5.216667),(44.9, 2.666667))
32
33 #####
34 # QUESTION 3 #
35 #####
36 def valeurModale(etiquettes) :
37     '''retourne le mode des valeurs de la liste d'étiquettes fournie en entrée
38     On considère que les étiquettes sont de type "chaîne de caractères'''
39     assert etiquettes != []
40     valeursRecensees = {} # dictionnaire associant une valeur à son effectif
41     for v in etiquettes :
42         if not( v in valeursRecensees.keys() ) :
43             # la valeur n'a pas encore été recensée
44             valeursRecensees[v] = 1
45         else :
46             valeursRecensees[v] = valeursRecensees[v] + 1
47
48     #on cherche la valeur associée à l'effectif maximal (en cas d'égalité on conserve la première
49     rencontrée)
50     maxi = 0
51     for clef in valeursRecensees.keys() :
52         if maxi < valeursRecensees[clef]:
53             mode = clef
54             maxi = valeursRecensees[clef]
55
56     return mode
57
58 def kNN(k, observation, apprentissage):
59     """Fonction qui effectue une prédiction avec l'algorithme k NN
60     'k' est la valeur de k pour l'algorithme kNN (nombre de plus proches voisins à considérer)
61     'observation' est le(s) prédicteur(s) de l'objet dont on cherche l'étiquette
62     'apprentissage' est une liste de couple au format objet=(prédicteurs, étiquette)"""
63     # On détermine la liste des k plus proches voisins
64     # voisins est une liste de tuple (prédicteur, étiquette) on commence par
65     # la remplir avec les k premiers voisins des données d'apprentissages
66     voisins = [(distance(observation,e[0]),e[1]) for e in apprentissage[:k]]
67     voisins.sort() #Trie la liste par ordre croissant on s'assurera par la suite qu'elle reste triée
68     for i in range(k,len(apprentissage)):
69         # distance entre cette valeur du prédicteur et notre observation
70         dist = distance(observation,apprentissage[i][0])
71         if dist < voisins[k-1][0]: # Si c'est plus petit de le plus lointain voisin
72             for j in range(k): # on cherche où l'insérer dans la liste des voisins
73                 if dist < voisins[j][0]: # Si plus petit
74                     voisins.insert(j,(dist,apprentissage[i][1])) # On l'insère à cette position
75                     voisins.pop() # Et on éjecte le plus grand des voisins (dernier élément)
76                 break # on arrête la recherche
77     # La liste des voisins étant constituée, il ne reste plus qu'à renvoyer le mode des étiquettes
78     return valeurModale([v[1] for v in voisins])
79
80 #####

```

```

81 # QUESTION 4 #
82 #####
83 apprentissage = extraireDonnees('apprentissage_departements')
84
85 def test(k,commune) :
86     print("les coordonnées GPS ",commune,"sont vraisemblablement situées dans le département",
87           kNN(k, commune, apprentissage))
88
89 test(7,(49.316667,1.283333))
90 test(7,(45.95, 5.36))
91 test(7,(45.95, 5.35))
92
93 test(6500,(49.316667,1.283333))
94 test(6500,(45.95, 5.36))
95 test(6500,(45.95, 5.35))

```

Q5:

```

1 import csv #pour lire le fichier de données d'apprentissage
2
3 #####
4 # QUESTION 1 #
5 #####
6 def extraireDonnees(nom) :
7     """Génère des données d'apprentissage sous la forme d'une liste de couple au
8     format objet=(predicteurs, étiquette) les données sont extraites du fichier
9     .csv dont les champs sont jour;Date;Code INSEE région;Région;TMin (°C);TMax (°C);TMoy (°C)"""
10    apprentissage=[]
11    with open(nom+'.csv', "r", encoding='utf-8') as fichier:
12        for donnees in csv.DictReader(fichier,delimiter=',') :
13            d=dict(donnees)
14            apprentissage.append(((d['Région'],int(d['jour'])),float(d['TMoy (°C)'])))
15    return apprentissage
16
17 #####
18 # QUESTION 2 #
19 #####
20 def distance(p1, p2):
21     """Calcule la distance entre les deux prédicteurs p1 et p2 qui vaut 400 (ou n'importe quelle
22     valeur strictement supérieure à 366) si les régions sont différentes et qui vaut le nombre
23     de jours d'écart entre les jours de l'année sinon."""
24     if p1[0] != p2[0]: return 400
25     j1,j2 = min(p1[1], p2[1]), max(p1[1], p2[1])
26     return min(j2 - j1, j1 + 366-j2)
27
28 #####
29 # QUESTION 3 #
30 #####
31 def moyenne(etiquettes):
32     '''retourne la moyenne des valeurs numériques de la liste d'étiquettes fournie en entrée'''
33     assert etiquettes != []
34     total = 0
35     for valeur in etiquettes :
36         total = total + valeur
37     moyenne = total / len(etiquettes)
38     return moyenne
39
40 def kNN(k, observation, apprentissage):
41     """Fonction qui effectue une prédiction avec l'algorithme k NN
42     'k' est la valeur de k pour l'algorithme kNN (nombre de plus proches voisins à considérer)
43     'observation' est le(s) prédicteur(s) de l'objet dont on cherche l'étiquette
44     'apprentissage' est une liste de couple au format objet=(predicteurs, étiquette)"""
45     # On détermine la liste des k plus proches voisins
46     # voisins est une liste de tuple (prédicteur, étiquette) on commence par
47     # la remplir avec les k premiers voisins des données d'apprentissages
48     voisins = [(distance(observation,e[0]),e[1]) for e in apprentissage[:k]]
49     voisins.sort() #Trie la liste par ordre croissant on s'assurera par la suite qu'elle reste triée
50     for i in range(k,len(apprentissage)):
51         # distance entre cette valeur du prédicteur et notre observation
52         dist = distance(observation,apprentissage[i][0])
53         if dist < voisins[k-1][0]: # Si c'est plus petit de le plus lointain voisin
54             for j in range(k): # on cherche où l'insérer dans la liste des voisins
55                 if dist < voisins[j][0]: # Si plus petit
56                     voisins.insert(j,(dist,apprentissage[i][1])) # On l'insère à cette position
57                     voisins.pop() # Et on éjecte le plus grand des voisins (dernier élément)
58                 break # on arrête la recherche
59     # La liste des voisins étant constituée,il ne reste plus qu'à renvoyer la moyenne des étiquettes
60     return moyenne([v[1] for v in voisins])
61
62 #####
63 # QUESTION 5 #

```

```
64 #####
65 mois = [(31,"janvier"), (29,"février"), (31,"mars"), (30,"avril"), (31,"mai"), (30,"juin"),
66         (31,"juillet"), (31,"août"), (30,"septembre"), (31,"octobre"), (30,"novembre"), (31,"décembre")]
67 def jour_mois(no_jour):
68     assert no_jour > 0 and no_jour < 367
69     jour = no_jour
70     no_mois = 0
71     while jour > mois[no_mois][0]:
72         jour -= mois[no_mois][0]
73         no_mois += 1
74     return str(jour) + " " + mois[no_mois][1]
75
76 #####
77 # QUESTION 4 #
78 #####
79 apprentissage = extraireDonnees('apprentissage_temperatures2016')
80 def test(liste_k, liste_jours, region):
81     for jour in liste_jours:
82         print("La température estimée pour la région", region,"le", jour_mois(jour), "est de :")
83         for k in liste_k:
84             # Ecriture de la température avec 3 chiffres significatifs
85             print(str(float('%0.3g' % kNN(k, (region, jour), apprentissage))) + "°C (k=" + str(k) + ")")
86     print()
87
88 test([10, 20, 60, 366], [20, 80, 190], 'Ile-de-France')
89 test([10, 20, 60, 366], [20, 80, 190], 'Occitanie')
90 test([10, 20, 60, 366], [20, 80, 190], 'Bretagne')
```

Pour $k = 366$, on obtient la température moyenne annuelle (quelque soit la date choisie), tandis que pour $k = 4000$, on a la température moyenne en France quelques soient la date ou la région choisies.