

ZOUHRI YASSINE | DSE

1 - Créez une base de 15 textes à partir du web (10 textes de culture, 5 textes d'économie, et 5 textes de politique)

```
# import statements
import numpy
import pandas as pd
import re
```

```

from stop_words import get_stop_words
import math

In [2]: # Creation de BE
BE=[["La Francophonie est également présente dans le sport.", "sport"], ["Au début, le s
t = [b[0] for b in BE ]
tps = [b[1] for b in BE ]
print(tps)

['sport', 'sport', 'sport', 'sport', 'sport', 'économie', 'économie', 'économie', 'éco
nomie', 'économie', 'politique', 'politique', 'politique', 'politique', 'politique']

2 - Créez une fonction qui crée une liste L de tous les mots figurant dans cette base (sans répétition).

In [3]: # Creation de la liste L de tous les mots figurant dans cette base (sans répétition)
def creation(liste_s):
    words = []

```

```
for s in liste_s:
    w = extraction(s)
    words.extend(w)
```

[illegible]

0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 1. 0. 0. 0. 0. 1. 0. 0. 0. 0.]

Les sports aiment le basket-ball, le football, le football pourraient être considérés

[illegible]

Fais donc preuve de conscience politique.

```
[0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 1. 0. 0.
 0. 0. 0. 1. 0. 0. 0. 0. 0. 1. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0.]
```

```

0.0 0.0 0.0 1.0 0.0 0.0 0.0 1.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0
0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.]

Ce succès est le résultat d'une politique commune, et c'est une bonne chose que l'on s'
oit à nouveau parvenu à consacrer 839 millions d'euros aux Balkans dans le budget 200
1.
[0. 0. 0. 0. 1. 1. 0. 0. 0. 1. 1. 0. 0. 0. 0. 1. 1. 0. 1. 0. 0. 0. 1. 0. 0. 0. 0.
0. 0. 1. 0. 0. 0. 1. 0. 0. 0. 1. 0. 0. 0. 0. 0. 1. 1. 0. 0. 0. 0. 0. 0. 0. 0.
1. 0. 0. 0. 0. 0. 0. 2. 0. 0. 0. 0. 0. 0. 0. 0. 1. 0. 0. 0. 0. 0. 0. 0. 1. 0.
1. 0. 0. 0. 1. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 1. 0. 0. 1. 0. 0. 0. 1. 0. 0. 0. 1. 0.
0. 0. 0. 0. 1. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 1. 0. 2. 0. 0. 0. 0. 0. 0. 0.]

In [5]: df_t2v = pd.DataFrame(bg, columns=vocab)
df_t2v

Out[5]:
```

| | actuel | aiment | attende | autes | aux | balkans | basketball | beaux | bonne | budget | ... | un | une | vote |
|----|--------|--------|---------|-------|-----|---------|------------|-------|-------|--------|-----|-----|-----|------|
| 0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | .. | 0.0 | 0.0 | 0.0 |
| 1 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | .. | 1.0 | 0.0 | 0.0 |
| 2 | 0.0 | 1.0 | 0.0 | 0.0 | 0.0 | 0.0 | 1.0 | 0.0 | 0.0 | 0.0 | .. | 0.0 | 0.0 | 0.0 |
| 3 | 0.0 | 0.0 | 1.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | .. | 0.0 | 0.0 | 0.0 |
| 4 | 0.0 | 0.0 | 0.0 | 1.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | .. | 0.0 | 0.0 | 0.0 |
| 5 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 1.0 | 0.0 | 0.0 | .. | 0.0 | 0.0 | 0.0 |
| 6 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | .. | 0.0 | 0.0 | 0.0 |
| 7 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | .. | 0.0 | 0.0 | 0.0 |
| 8 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | .. | 1.0 | 0.0 | 0.0 |
| 9 | 1.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | .. | 0.0 | 0.0 | 1.0 |
| 10 | 0.0 | 0.0 | 0.0 | 0.0 | 1.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | .. | 0.0 | 0.0 | 0.0 |
| 11 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | .. | 0.0 | 1.0 | 0.0 |
| 12 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | .. | 0.0 | 0.0 | 0.0 |
| 13 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | .. | 0.0 | 0.0 | 0.0 |
| 14 | 0.0 | 0.0 | 0.0 | 0.0 | 1.0 | 1.0 | 0.0 | 0.0 | 1.0 | 1.0 | .. | 0.0 | 1.0 | 0.0 |

15 rows x 119 columns

4 - Créez une liste de mots vides MV (pronoms, prépositions, conjonctions) et utilisez la pour réduire la dimension de L. Quel est le taux de réduction obtenu ?

```

In [6]: def mv(voc):
        ignore = get_stop_words('french')
        v = vc
        for k in v:
            if k in ignore:
                v[k] = temp
        return v
mv_vocab = mv(vocab)
mv_len_vocab = len(mv_vocab)
print("Le taux de reduction est de : {} %\n".format(1-(mv_len_vocab/old_len_vocab)))
print("(0) words in this List for the '15 textes' \n\n (1) \n".format(len(mv_vocab),))

Le taux de reduction est de : 0.2016806722689075 %

95 words in this List for the '15 textes'
```

- 'actuel', 'aiment', 'attende', 'autres', 'balkans', 'basketball', 'beaux', 'bonne', 'budget', 'budgétaire', 'canadienne', 'cardio', 'cest', 'chose', 'commune', 'comprend', 'compétent', 'conférence', 'consacrer', 'conscience', 'conseillère', 'considerer', 'e', 'considérés', 'continuent', 'deux', 'devriezvous', 'dici', 'divisions', 'dune', 'débat', 'démocratie', 'effectif', 'européenne', 'fiscale', 'football', 'forment', 'francophonie', 'grandes', 'gérées', 'instruit', 'internationale', 'jeune', 'jours', 'la

esser, league, l'excellent, liberaliser, l'objectif, ion, l'économie, mard
es', 'mardi', 'mercredi', 'millions', 'moment', 'mondiales', 'notamment', 'parvenu',
'petit', 'pleine', 'politiciens', 'politique', 'porte', 'portoricaine', 'pourraient',
'preuve', 'professionnelles', 'prosperité', 'présente', 'puissances', 'renoue', 'réa

[illegible]

Les Maldives sont une jeune démocratie, en pleine transition politique.

[illegible]

```

    if checker(s) and checker(e) and checker(p):
        mots.append(c)
print("Les mots de la liste L, ayant figuré dans les 3 classes, sont :\n {}".format(mots))

```

```
print("Apres T2V : '(n')\n")  
elim = elim(df_t2v_tz)
```

Après T2V :

Les mots de la liste L, ayant figuré dans les 3 classes, sont :
[de', le', 'soint']

```
n [10]:  
  
elim_vocab = [temp for temp in vocab if temp not in elim ]  
elim_len_vocab = len(elim_vocab)  
  
print("\nLe taux de reduction est de : {}{}\n".format(1-(elim_len_vocab/old_len_vocab)))  
print("(f)() words in this List for the '15 textes' \\\n\n {}".format(len(elim_vocab),  
bg_mv_elim = t2v(t, elim_vocab))
```

Le taux de reduction est de : 0.02521008403613467
116 words in this List for the '15 textes'

[l'actuel', 'aient', 'attendue', 'autres', 'aux', 'baikans', 'basketball', 'beaux',
'bonnes', 'budget', 'budgétaires', 'canadienne', 'cardinal', 'dest', 'chose', 'comme', 'co-
munne', 'comprend', 'compétent', 'conférence', 'consacrer', 'conscience', 'conservati-
ce', 'considerer', 'considérés', 'continent', 'dans', 'des', 'deuros', 'devriezvous',
'dicit', 'divisions', 'donc', 'dune', 'débat', 'début', 'démocratie', 'effectif', 'en',
'est', 'et', 'européenne', 'fiscale', 'football', 'forment', 'francophonie', 'grande s'
'gerées', 'instruit', 'internationale', 'jeune', 'jours', 'la', 'laisser', 'leagu e'
, 'leur', 'l'excellent', 'libéraliser', 'objectif', 'lon', 'léconomie', 'maldives',
'mardi', 'mercredi', 'millions', 'moment', 'mondiales', 'notamment', 'nouveau', 'par',
'parvenu', 'petit', 'pleine', 'politiciens', 'politique', 'porte', 'portoricaïne', 'po-
'pourrait-être', 'renewe', 'professionnelles', 'prosperité', 'présente', 'puissance s'
, 'que', 'reposé', 'réduction', 'résultats', 'second', 'secteur', 'services', 'simula-
teur', 'soit', 'soldat', 'sport', 'sports', 'stimulerait', 'succès', 'suite', 'sur',
'tant', 'terminés', 'the', 'tourisme', 'transition', 'trois', 'un', 'une', 'vote',
'à', 'économie', 'économique', 'égaleme't', 'émissions', 'épaisser', 'être']

La Francophonie est également présente dans le sport.

```
[0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0.  
0. 0. 1. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 1. 0. 0. 0. 0.  
0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0.  
0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0.  
0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0.  
0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0.]
```

Au début, le football est considéré comme un petit moment économique.

```
[0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0.  
0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0.  
0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0.  
0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0.  
0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0.  
0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0.]
```

Les sports aiment le basket-ball, le football, le football pourraient être considérés cardio.

```
[0. 1. 0. 0. 0. 0. 1. 0. 0. 0. 0. 0. 1. 0. 0. 0. 0. 0. 0. 0. 0. 0.  
1. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0.  
0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0.  
0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0.  
0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0.  
0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0.]
```

FIFA 14 est la suite tant attendue de l'excellent simulateur de football.

```
[0. 0. 1. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0.  
0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0.  
0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0.  
0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0.  
0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0.  
0. 0. 1. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0.]
```

Les trois autres divisions professionnelles sont gérées par The Football League.

```
[0. 0. 0. 0. 1. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0.  
0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0.  
0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0.  
0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0.  
0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0.  
0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0.]
```

Les beaux jours de l'économie conservatrice sont terminés, pour le moment.

```
[0. 0. 0. 0. 0. 0. 0. 0. 1. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0.  
0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0.  
0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0.  
0. 0. 0. 1. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0.]
```

0. 0. 0. 0. 0. 1. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0.
0. 0. 0. 0. 0. 1. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0. 0.]

Les grandes puissances mondiales continuent de libéraliser leur économie

[illegible]

7 - Proposez un texte pour lequel : (en utilisant T2V)

```

n [12]: #définition de KNN
def knn(k,t,B):
    #créer la table des distances T-B
    B = [[b,t] for b,t in zip(B,tps)]
    tD=[]
    for i in range(len(B)):
        tD.append([dist(t,B[i][0]),B[i][1]])
    #trier la tD
    for i in range(len(tD)):
        for j in range(len(tD)):
            if tD[i][0]>tD[j][0]:
                temp=tD[i]
                tD[i]=tD[j]
                tD[j]=temp
    pp.pprint(tD)
    #calculer les k derniers de la tD
    #créer les tables tC et tO
    tC=[]
    tO=[]
    for i in range(len(tD)):
        if not B[i][1] in tC:
            tC.append(B[i][1])
            tO.append(0)

    #remplir la table tO
    for i in range(1,k+1):
        for j in range(len(tC)):
            if tD[-i][1]==tC[j]:
                tO[j]=tO[j]+1
    pp.pprint(tO)
    #retourner la classe prédite
    if tO.count(max(tO)) == 1:
        indmax=tO.index(max(tO))
        return tC[indmax]
    return tD[-1][1]

n [13]: texte_eco = "Les exportations de l'économie grecque sont traditionnellement faibles."
bag_vector_eco = vector(vocab, texte_eco)

predicted_eco = knn(4, bag_vector_eco, bg)
print("\nParmi les 15 textes de la base, la DIST minimale (en utilisant T2V) correspond à la classe : ÉCONOMIE"

[[[5.830951894845301, 'politique'],
[4.795851523312719, 'sport'],
[3.872983346207417, 'économie'],
[3.7416573867739413, 'économie'],
[3.605551275463989, 'sport'],
[3.4641016151377544, 'politique'],
[3.3166247903554, 'politique'],
[3.3166247903554, 'sport'],
[3.3166247903554, 'sport'],
[3.162276601683795, 'politique'],
[3.162276601683795, 'sport'],
[3.162276601683795, 'économie'],
[3.0, 'économie'],
[2.6457513110645907, 'économie'],
[2.4494849742783178, 'politique']]
[0, 3, 1]

n [14]: Parmi les 15 textes de la base, la DIST minimale (en utilisant T2V) correspond à la classe : ÉCONOMIE

8 - Proposez un texte pour lequel : (en utilisant T2V + MV + ELIM)

• Parmi les 15 textes de la base, la DIST minimale correspond à la classe sport.</strong>

n [14]: texte_sport = "Je suis très heureux de la performance du défenseur et de celle de l'éc
bag_vector_sport = vector(mv_vocab, texte_sport)

```

```
predicted_sport = knn(4, bag_vector_sport, bg_mv)
print("\nParmi les 15 textes de la base, la DIST minimale (en utilisant T2V + MV +
```

```
[14.0, 'politique'],
[3.7416573867739413, 'économie'],
[13.622776601683795, 'économie'],
[2.8284271247461903, 'économie'],
[2.6457513110645907, 'sport'],
[2.6457513110645907, 'sport'],
[2.6457513110645907, 'économie'],
[2.6457513110645907, 'politique'],
[2.6457513110645907, 'politique'],
[2.6457513110645907, 'économie'],
[2.4494897427832178, 'politique'],
[2.23606797749979, 'sport'],
[2.23606797749979, 'sport'],
[2.0, 'politique'],
[2.0, 'sport']]
[3, 0, 1]

Parmi les 15 textes de la base, la DIST minimale (en utilisant T2V + MV + ELIM) corres-
pond à la classe : SPORT

9 - Proposez un texte pour lequel :

• La classe de la DIST minimale (en utilisant T2V) est différente de la classe de la DIST minimale
(en utilisant T2V + MV + ELIM) </strong>

n [15]:
texte_9 = "Au début, le football est comme une jeune démocratie qui repose sur la poli-
bag_vector_t2v = vector(vocab, texte_9)

predicted_t2v = knn(4, bag_vector_t2v, bg)
print("\nParmi les 15 textes de la base, la DIST minimale (en utilisant T2V) correspo-

[[5.744562646538029, 'politique'],
[5.477225575051661, 'économie'],
[4.795831523312719, 'économie'],
[4.69041575982343, 'sport'],
[4.472135954993958, 'sport'],
[4.472135954993958, 'économie'],
[4.472135954993958, 'économie'],
[4.358898943540674, 'économie'],
[4.242640687119285, 'sport'],
[4.0, 'politique'],
[3.872983346207417, 'politique'],
[3.872983346207417, 'politique'],
[3.872983346207417, 'sport'],
[3.605551275463989, 'politique'],
[3.4641016151377544, 'sport']]
[2, 0, 2]

Parmi les 15 textes de la base, la DIST minimale (en utilisant T2V) correspond à la c-
lasse : SPORT

n [16]:
bag_vector_t2v_mv_elim = vector(elim_vocab, texte_9)

predicted_mv_t2v_elim = knn(4, bag_vector_t2v_mv_elim, bg_mv_elim)
print("\nParmi les 15 textes de la base, la DIST minimale (en utilisant T2V + MV + EL-

[[5.656854249492381, 'politique'],
[5.0, 'économie'],
[4.58257569459584, 'économie'],
[4.242640687119285, 'économie'],
[4.242640687119285, 'sport'],
[4.242640687119285, 'sport'],
[4.242640687119285, 'économie'],
[4.123105625617661, 'économie'],
[3.872983346207417, 'politique'],
[3.872983346207417, 'sport'],
[3.7416573867739413, 'politique'],
[3.605551275463989, 'politique'],
[3.605551275463989, 'sport'],
[3.4641016151377544, 'sport'],
[3.3166247903354, 'politique']]
[2, 0, 2]

Parmi les 15 textes de la base, la DIST minimale (en utilisant T2V + MV + ELIM) corres-
```

10 - Quel modèle d'apprentissage artificiel avez-vous utilisé dans les questions 7 à 9 ? Donnez 1 avantage et 1 inconvénient de ce modèle

Le modèle d'apprentissage artificiel utilisé dans les questions 7 à 9 est : **KNN**

Un avantage :

- Il s'agit d'un modèle simple et facile à appliquer dans la pratique..

Un inconvénient :

- Le modèle est plus lent lorsque le nombre de ses échantillons d'apprentissage est plus élevé..