

UNIVERSITÉ DES SCIENCES ET DE LA  
TECHNOLOGIE HOUARI BOUMÉDIÈNE

MASTER SII

---

## Projet de TP en Recherche d'Information

---

AMMAR KHODJA Hichem  
BOUDJENIBA Oussama

6 mars 2021





# Table des matières

<b>I</b>	<b>Introduction</b>	<b>2</b>
<b>II</b>	<b>Lecture et prétraitement de la collection CACM</b>	<b>2</b>
II.1	Lecture et extraction . . . . .	2
II.2	Tokenisation . . . . .	2
II.3	Résultat . . . . .	2
<b>III</b>	<b>Modèle booléen</b>	<b>2</b>
III.1	Structure de données . . . . .	2
III.2	Évaluation des requêtes . . . . .	3
III.3	Temps d'exécution . . . . .	4
III.4	Exemple d'une requête . . . . .	4
<b>IV</b>	<b>Modèle vectoriel</b>	<b>5</b>
IV.1	Structure de données . . . . .	5
IV.2	Les fonctions d'accès . . . . .	6
IV.3	Pondération TF-IDF . . . . .	7
IV.4	Calcul des normes . . . . .	7
IV.5	Évaluation des requêtes . . . . .	8
IV.5.1	Calcul de la similarité . . . . .	8
IV.6	Temps d'exécution . . . . .	10
IV.7	Exemple d'une requête . . . . .	10
IV.8	Optimisation du modèle (Bonus/non demandé) . . . . .	11
<b>V</b>	<b>Phase d'indexation en chiffres</b>	<b>12</b>
<b>VI</b>	<b>Évaluation du modèle vectoriel</b>	<b>13</b>
VI.1	Modèle vectoriel : Juger la pertinence d'un item . . . . .	13
VI.1.1	Méthode des $k$ -meilleurs documents . . . . .	13
VI.1.2	Méthode de seuillage . . . . .	14
VI.2	Combinaison des méthodes . . . . .	14
VI.3	Moyenne des courbes précision-rappel interpolées . . . . .	14
VI.4	Interface graphique (Bonus/Non demandé) . . . . .	15

## Partie I : Introduction

Dans ce document, nous allons expliquer la procédure suivie pour construire un modèle de recherche d'information et l'appliquer sur la collection CACM. Ce document est divisé en plusieurs chapitres :

- Lecture et prétraitement de la collection CACM
- Conception du modèle Booléen
- Conception du modèle vectoriel
- Évaluation du modèle graphique
- Interface graphique pour tester les performances du modèle vectoriel.

## Partie II : Lecture et prétraitement de la collection CACM

### II.1 Lecture et extraction

On commence par lire le fichier *cacm.all*. Ensuite, Chaque document dans CACM commence par **.I** ce qui nous permet de détecter le début et fin de chaque document. Après, on extrait de chaque document son numéro d'identification (qui vient après le **.I**), le titre (**.T**), le nom des auteurs (**.A**) et éventuellement le résumé (**.W**).

### II.2 Tokenisation

Maintenant, on construit un dictionnaire qui a comme clé l'ID du document et la valeur contient le reste des informations (**.T**, **.A** et **.W**) regroupées en une seule chaîne de caractères qu'on met **en minuscule**. On procède à présent à la phase de tokenisation, on définit un token comme étant une suite de caractères alphanumériques ce qui est équivalent à l'expression régulière suivante "**\w+**". On utilise la **stopword liste de NLTK** pour éliminer les mots non-significatifs.

### II.3 Résultat

Le résultat est un dictionnaire contenant comme clé l'identificateur d'un document et la valeur correspondante est une liste des tokens contenus de ce document.

## Partie III : Modèle booléen

### III.1 Structure de données

Une structure de données plus appropriée est requise pour optimiser le temps d'exécution. On obtient cette structure en modifiant le dictionnaire obtenu dans la

section II.3 suivant l'algorithme :

---

**Algorithme 1** : Construction de la structure de données pour le modèle booléen

---

**Entrées** :  $D$  : Dictionnaire obtenu dans la section II.3  
**Sorties** :  $D'$  : Dictionnaire prêt pour le modèle booléen

```
1  $D' \leftarrow \text{Dictionnaire}()$ ;  
2 pour  $id, L_{tokens}$  dans  $D$  faire  
3    $S \leftarrow \text{ensemble}(L_{tokens})$ ;  
4   pour  $token$  dans  $S$  faire  
5     si  $token$  n'est pas dans  $D'$  alors  
6        $D'.\text{ajouter}(token, \{\})$ ;  
7     fin  
8      $D'[token].\text{ajouter}(id)$ ;  
9   fin  
10 fin  
11 retourner  $D'$ ;
```

---

On sauvegarde aussi la liste de tous les documents dans un ensemble, ce qui sera utile par la suite pour l'opérateur **not** (Voir ci-dessous).

Le résultat est un dictionnaire contenant comme clé un **token** et comme valeur l'ensemble des documents contenant ce token.

## III.2 Évaluation des requêtes

Une requête est une expression booléenne, pouvant contenir des parenthèses imbriquées et ayant comme opérateurs : **and**, **or** et **not**.

Un exemple d'une requête du modèle booléen :

('Mathematics' or 'Science') and not 'Algebra'

Comme on peut l'observer, les tokens sont entourés d'apostrophe et qui sont reliés entre-eux par des opérateurs booléens (sans apostrophes). Les trois opérateurs booléens acceptés sont : **and**, **or** et **not**. Les expressions booléennes peuvent être imbriquées.

Pour évaluer une requête, on commence par vérifier que la syntaxe est correcte. Ensuite, on remplace respectivement les opérateurs **and**, **or** et **not** par &, | et -. Si on suit l'exemple précédant, ça nous donne :

('Mathematics' | 'Science') & -'Algebra'

Maintenant, on entoure chaque token de la manière suivante : **'token'**  $\rightarrow$  **Token('token')**. Ce qui nous donne :

( Token('Mathematics') | Token('Science') ) & -Token('Algebra')

**Token** est une classe qu'on a codé dont le comportement ressemble à celui de la classe **set** dans Python. A l'initialisation d'une instance de cette classe de la manière suivante : **Token(T)** où **T** est une chaîne de caractères, on récupère l'ensemble des documents contenant le token **T**.

Petite remarque : Cette manoeuvre justifie le choix de la structure de données présentée dans la section III.1 car, elle nous garantit l'accès en  $O(1)$ .

Maintenant, il n'y a plus qu'à définir le comportement des opérateurs booléens :

- **La conjonction de deux Token (T1 and T2)** : retourne un nouveau **Token** contenant l'intersection de l'ensemble de documents de **T1** et l'ensemble de documents de **T2**.
- **La disjonction de deux Token (T1 or T2)** : retourne un nouveau **Token** contenant l'union de l'ensemble de documents de **T1** et l'ensemble de documents de **T2**.
- **La négation d'un Token T** : retourne un nouveau **Token** contenant le complément de l'ensemble des documents de **T** par rapport à l'ensemble de tous les documents.

Enfin, il suffit d'exécuter la commande **eval(R)** de Python, **R** étant la requête transformée, ce qui nous retourne le résultat de la requête par le modèle booléen.

### III.3 Temps d'exécution

Pour évaluer les performances du modèle booléen par rapport au temps d'exécution de l'évaluation d'une requête, on a fait l'expérience suivante :

- On génère 100000 requête aléatoires : 1000 requête contenant 1 token, 1000 requête contenant 2 tokens, ..., 1000 requête contenant 100 tokens. On relie les tokens aléatoirement par les opérateurs booléens (**and**, **or**, **not**).
- On évalue ces requêtes par le modèle booléen et on enregistre le temps d'exécution.
- On regroupe les résultats par nombre de tokens et on calcule le temps d'exécution moyen.

Maintenant, on affiche le graphe de la fonction **F(Nombre de tokens) = Temps d'exécution** dans la Figure 1. On observe que le temps d'exécution dépend linéairement du nombre de tokens dans la requête booléenne.

### III.4 Exemple d'une requête

On teste notre modèle sur la collection CACM en lui demandant d'évaluer la requête suivante :

**('science' or 'compiler') and not 'algebra' and 'code'**

Le modèle nous retourne une série de documents (leurs identifiants) : **123, 1223, 1234, 1542, 1551, 1613, 1807, 2064, 2423, 2433, 2897, 2968, 3080** comme on peut le voir dans la Figure 2.

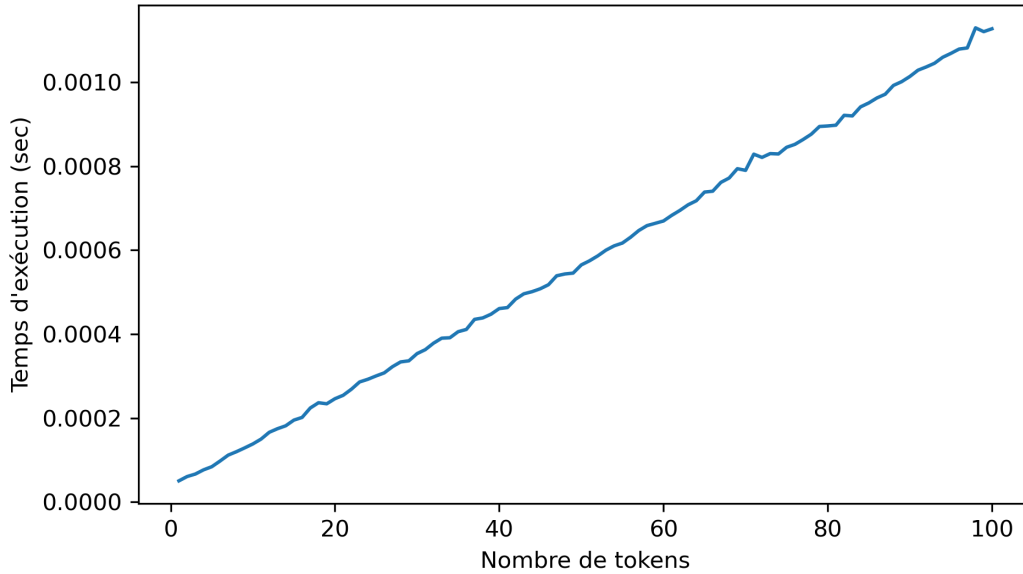


FIGURE 1 – Performances du modèle booléen en termes de temps d'exécution

## Tester une requête (Modèle booléen) ¶

```
: query = " ('science' or 'compiler') and not 'algebra' and 'code'"
bm.eval(query)

{123, 1223, 1234, 1542, 1551, 1613, 1807, 2064, 2423, 2433, 2897, 2968, 3080}
```

FIGURE 2 – Exemple d'une requête pour le modèle booléen

## Partie IV : Modèle vectoriel

### IV.1 Structure de données

Comme pour le modèle booléen, une structure de données plus appropriée est requise pour optimiser le temps d'exécution, à savoir, le **fichier inverse**. On obtient cette structure en modifiant le dictionnaire obtenu dans la section II.3 suivant l'algorithme :

**La fonction  $\text{Count}(\mathbf{L})$**  : Prend en entrée une liste  $\mathbf{L}$  de tokens et retourne un dictionnaire qui a comme clé un token et comme valeur le nombre d'occurrences de ce token dans la liste  $\mathbf{L}$ .

Exemple :  $\mathbf{L} = ['A', 'B', 'C', 'B', 'C', 'D'] \longrightarrow \text{Count}(\mathbf{L}) = \{'A' : 1, 'B' : 2, 'C' : 2, 'D' : 1\}$

On garde aussi la fréquence maximale de chaque document dans un dictionnaire qu'on nomme  $\text{max\_freq}$ . Ce dictionnaire nous sera utile par la suite pour la pondération du fichier inverse, notamment, pour l'optimisation du temps d'exécution. On

sauvegarde aussi le nombre total de documents  $N$  pour les mêmes raisons.

---

**Algorithme 2** : Construction de la structure de données pour le modèle vectoriel

---

```

1 Fonction Count(Liste  $L$ )
2 Debut
3    $D \leftarrow \text{Dictionnaire}()$ ;
4   pour token dans  $L$  faire
5     si token dans  $D.\text{keys}()$  alors
6        $D[\text{token}] \leftarrow 0$ ;
7     fin
8      $D[\text{token}] \leftarrow D[\text{token}] + 1$ ;
9   fin
10  retourner  $D$ ;
11 Fin
    /* On passe maintenant à la construction du fichier inverse */
Entrées :  $D$  : Dictionnaire obtenu dans la section II.3
Sorties :  $D'$  : Dictionnaire prêt pour le modèle vectoriel (Fichier inverse),
            $\text{max\_freq}$  : Dictionnaire contenant la fréquence de chaque
           document,  $N$  : le nombre total de documents
12  $D' \leftarrow \text{Dictionnaire}()$ ;
13  $\text{max\_freq} \leftarrow \text{Dictionnaire}()$ ;
14  $N \leftarrow \text{taille}(D)$ ;
15 pour  $\text{id}, L_{\text{tokens}}$  dans  $D$  faire
16    $S \leftarrow \text{Count}(L_{\text{tokens}})$ ;
17    $\text{max\_freq}[\text{id}] \leftarrow \text{max}(L_{\text{tokens}})$  pour token,  $n$  dans  $S$  faire
18     si token n'est pas dans  $D'$  alors
19        $D'.\text{ajouter}(\text{token}, \text{Dictionnaire}())$ ;
20     fin
21      $D'[\text{token}][\text{id}] \leftarrow n$ ;
22   fin
23 fin
24 retourner  $D', \text{max\_freq}, N$ ;

```

---

## IV.2 Les fonctions d'accès

- **Freq(token, id)** : Prend en entrée un token et l'identificateur d'un document (**id**) et retourne la fréquence du token dans ce document. Cette fonction est très simple : Si  $\mathbf{D}$  est le fichier inverse (dictionnaire de dictionnaires), cette fonction retourne  $\mathbf{D}[\text{token}][\text{id}]$ . Cette fonction est de complexité  $\mathbf{O}(1)$ . Cette fonction retourne 0 si la paire (**token, id**) n'existe pas dans le fichier inverse.
- **Documents(token)** : Prend en entrée un token et retourne les documents contenant ce document avec la fréquence du token dans ces documents. Si  $\mathbf{D}$  est le fichier inverse, cette fonction retourne  $\mathbf{D}[\text{token}]$ . Cette fonction est de



complexité  $O(1)$ .

- **Tokens(id)** : Prend en entrée un identificateur d'un document (**id**) et retourne les tokens qui sont présents dans ce document et leurs fréquences dans ce document. Cette fonction parcourt l'ensemble du fichier inverse et enregistre les tokens qui sont présents dans le document **id** avec leurs fréquences. Cette fonction est donc, de complexité  $O(n)$ , **n** étant le nombre de tokens dans le fichier inverse.

### IV.3 Pondération TF-IDF

Pour chaque paire (**token, document**) dans le fichier inverse, on calcule le poids de cette paire selon la formule TF-IDF suivante :

$$poids(t_i, d_j) = \frac{freq(t_i, d_j)}{max\_freq(d_j)} \log_{10} \left( \frac{N}{n_i} + 1 \right)$$

Tel que :

- $freq(t_i, d_j)$  est la fréquence du token  $t_i$  dans le document  $d_j$ .
- $max\_freq(d_j) = max(\{freq(t, d_j) : t \in d_j\})$  est la fréquence du token qui apparaît le plus dans le document  $d_j$ . Cet valeur a été calculée pour chaque document dans l'Algorithme 2 où les résultats sont stockées dans la variable  $max\_freq$ , ceci nous évite de recalculer à chaque itération.
- $N$  : le nombre total de documents qu'on a calculé dans l'Algorithme 2.
- $n_i$  : Le nombre de documents contenant le token  $t_i$ .

On déroule, donc, l'algorithme suivant pour pondérer le fichier inverse :

---

#### Algorithme 3 : Pondération du fichier inverse selon la formule TF-IDF

---

**Entrées** :  $D$  : Fichier inverse,  $max\_freq$  : la fréquence maximale de chaque document,  $N$  : le nombre total de documents

```

1 pour  $t$ ,  $D_{doc}$  dans  $D$  faire
2    $n \leftarrow taille(D_{doc})$ ;
3   pour  $d$  dans  $D_{doc}.keys()$  faire
4      $D_{doc}[d] \leftarrow D_{doc}[d] / max\_freq[d] * \log_{10}(N/n + 1)$ ;
5   fin
6 fin
```

---

### IV.4 Calcul des normes

Hors la fonction de similarité *produit interne*, les autres fonctions de similarité (Dice, Jaccard, Cosinus) requièrent le calcul de la norme euclidienne des documents. Pour optimiser la phase d'évaluation des requêtes, on calcule à l'avance cette norme pour tous les documents et on stock le résultat dans un dictionnaire (**document, norme**) qu'on appellera *doc\_norm*.

Soit  $d$  un document de la collection. On calcule  $doc\_norm[d]$  grâce à la formule suivante :

$$doc\_norm[d] = \sqrt{\sum_{i=1}^T d[t_i]^2}$$

Tel que :

- $T$  : est le nombre de tokens dans la collection
- $d[t_i]$  : est le poids de la paire  $(d, t_i)$  dans le fichier inverse pondéré. Ce poids vaut 0 si  $t_i \notin d$ .

On a besoin de prendre en considération que les tokens présents dans le document  $d$ . En optimisant la formule précédente, on obtient :

$$doc\_norm[d] = \sqrt{\sum_{t \in d} d[t]^2}$$

## IV.5 Évaluation des requêtes

D'abord, on commence par extraire les tokens de la requête en utilisant l'expression régulière " $\backslash \mathbf{w} +$ " qu'on stocke dans une liste. Ensuite, on mets en minuscule les tokens et on supprime de la liste les tokens qui sont présents dans la *stoplist* de NLTK. On peut légitimement supposer que les tokens qui apparaissent plusieurs fois dans la requête sont plus importants que ceux qui n'apparaissent qu'une seule fois. Donc, on applique la fonction **Count** décrite dans la Section IV.1 à cette liste, ce qui nous donne un dictionnaire (**token, fréquence**) qu'on nomme  $R$  et on divise chaque fréquence de  $R$  par la fréquence maximale de  $R$ . Le résultat est un dictionnaire (**token, valeur**) où **valeur** est un nombre entre 0 et 1. Cette valeur représente l'importance d'un token dans la requête. Si tous les tokens apparaissent le même nombre de fois dans la requête, on aura  $valeur = 1$  pour chaque token dans  $R$ . Le calcul de la norme euclidienne de la requête  $R$  se fait de la manière suivante :  $\|R\|_2 = \sqrt{\sum_{t \in R} R[t]^2}$

### IV.5.1 Calcul de la similarité

La similarité se calcule entre un document et une requête. Dans ce document, on va étudier quatre fonctions de similarité qui sont : *produit interne*, *Dice*, *cosinus* et *Jaccard*.

- **Produit interne** : Il se calcule avec la formule suivante :

$$PI(d, R) = \sum_{i=1}^T d[t_i] * R[t_i]$$

tel que :

- $d[t_i]$  : est le poids de la paire  $(t_i, d)$  dans le fichier inverse.
- $R[t_i]$  : est le poids du token  $t_i$  dans  $R$  (calculé plus haut). Ce poids vaut 0 si  $t_i \notin R$ .

—  $T$  : est le nombre total de tokens dans la collection.

On peut voir que si  $t_i \notin R$ , alors  $d[t_i] * R[t_i] = 0$ . Par conséquent, on peut simplifier et optimiser par la même occasion la formule précédente :

$$PI(d, R) = \sum_{t \in R} d[t] * R[t]$$

Voici l'algorithme pour calculer le produit interne de manière optimale pour tous les documents :

---

**Algorithme 4** : Calcul du produit interne entre la requête et tous les documents

---

**Entrées** :  $D$  : Fichier inverse pondéré,  $R$  : Dictionnaire représentant la requête,  $all\_documents$  : L'ensemble des identificateurs de tous les documents

**Sorties** :  $doc\_score$  : Dictionnaire contenant pour chaque document, son produit interne avec la requête  $R$

```

1  $doc\_score \leftarrow Dictionnaire()$ ;
2 pour  $d$  dans  $all\_documents$  faire
3    $doc\_score[d] \leftarrow 0$ ;
4   pour  $t, w$  dans  $R$  faire
5      $doc\_score[d] \leftarrow doc\_score[d] + D.poids(t, d) * w$ ;
6   fin
7 fin
8 retourner  $doc\_score$ ;

```

---

- **Dice** : Cette fonction de similarité se calcule avec la formule suivante :

$$Dice(d, R) = \frac{2 * PI(d, R)}{doc\_norm[d]^2 + \|R\|_2^2}$$

- **Cosinus** : Cette fonction de similarité se calcule avec la formule suivante :

$$Cosinus(d, R) = \frac{PI(d, R)}{doc\_norm[d] * \|R\|_2}$$

- **Jaccard** : Cette fonction de similarité se calcule avec la formule suivante :

$$Jaccard(d, R) = \frac{PI(d, R)}{doc\_norm[d]^2 + \|R\|_2^2 - PI(d, R)}$$

Après avoir calculé le score de chaque document avec la fonction de similarité choisie, on ordonne ces documents suivant leurs scores. On choisit les  $k$ -premiers documents et on les présente au client (choisir  $k$  selon l'application).

## IV.6 Temps d'exécution

Pour évaluer les performances du modèle vectoriel par rapport au temps d'exécution de l'évaluation d'une requête, on a fait l'expérience suivante :

- On génère 2000 requêtes aléatoires : 10 requêtes contenant 1 token, 10 requêtes contenant 2 tokens, ..., 10 requêtes contenant 200 tokens. Les requêtes sont composées de tokens tirés au hasard.
- On évalue ces requêtes par le modèle vectoriel et on enregistre le temps d'exécution.
- On regroupe les résultats par nombre de tokens et on calcule le temps d'exécution moyen.

Maintenant, on affiche le graphe de la fonction  $F(\text{Nombre de tokens}) = \text{Temps d'exécution}$  dans la Figure 3. On observe que le temps d'exécution dépend linéairement du nombre de tokens dans la requête.

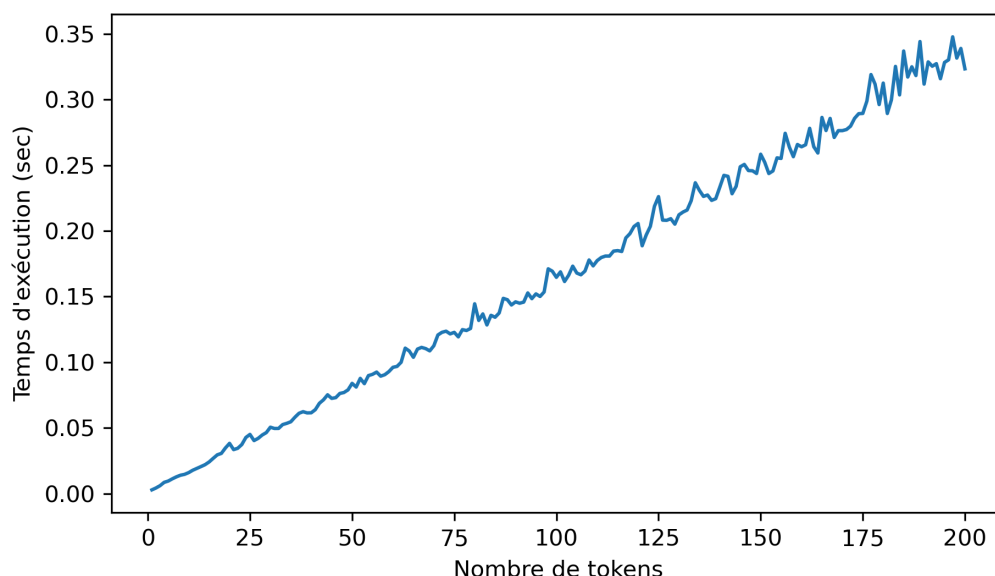


FIGURE 3 – Performances du modèle vectoriel en termes de temps d'exécution

## IV.7 Exemple d'une requête

On teste notre modèle sur la collection CACM en lui demandant d'évaluer la requête suivante :

**I want to consult a document about code optimization and compilers**

Le modèle nous retourne une série de documents et leurs scores respectifs suivant la fonction de similarité choisie (Figure 4).

## Tester une requête (Modèle vectoriel)

```
.2]: query = "I want to consult a document about code optimization and compilers"
      vm.eval(query, k=-1, sim=SimilarityFunctions.DOT)

: [(2233, 2.5436549019179804),
   (2575, 2.361458004084773),
   (2307, 2.361458004084773),
   (2423, 2.3409771342458017),
   (2897, 2.1261927617454233),
   (2658, 1.9663429036599014),
   (2551, 1.9663429036599014),
   (100, 1.9663429036599014),
   (1652, 1.882473941149483),
   (2835, 1.878013833046415),
   (2611, 1.878013833046415),
   (1231, 1.878013833046415),
   (1807, 1.7964931781383724),
```

FIGURE 4 – Exemple du résultat d'évaluation d'une requête utilisant le modèle vectoriel avec comme fonction de similarité le *produit interne*

### IV.8 Optimisation du modèle (Bonus/non demandé)

On a observé dans la section précédente que pour une requête de 100 tokens, l'évaluation de la requête prends 0.16s. Bien que ce temps paraît court, il commence à poser problème si, par exemple, notre modèle reçoit 1000 requêtes simultanément. Un utilisateur attendra en moyenne  $1000 \times 0.16 = 16s$  ce qui est inacceptable compte tenu des standards actuels. Heureusement, il y a un meilleur moyen pour calculer la similarité requête/document : les matrices creuses.

La quasi totalité du temps d'exécution s'écoule lors du calcul du produit interne entre la requête et les documents. On peut utiliser les matrices creuses pour accélérer cela. Donc, au lieu d'avoir un dictionnaire de dictionnaires, on aura une matrice creuse  $M$  de dimension  $T \times D$  ( $T$  : nombre de tokens,  $D$  : nombre de documents). La requête  $R$  est représentée par un vecteur ligne creux de dimension  $1 \times T$ . Ce vecteur contient 0 si le token n'appartient pas au vecteur sinon il contient le poids du token dans la requête.

Maintenant, pour calculer le produit interne entre  $R$  et tous les documents, il suffit de calculer  $R \times M$ , ce qui nous retourne un vecteur dont chaque composante est le produit interne entre un document et  $R$ . Ces matrices/vecteurs creux sont sous format **CSR (Compressed Sparse Row)** qui a été optimisé pour les opérations arithmétique (addition, soustraction, produit matriciel).

On refait la même expérience que dans la section IV.6 et cette fois-ci, les résultat sont nettement meilleur comme on peut le voir dans la Figure 5. On passe d'un temps d'exécution en moyenne de 0.16s pour une requête de 100 tokens à 0.001s, ce qui est **160 fois plus rapide**.

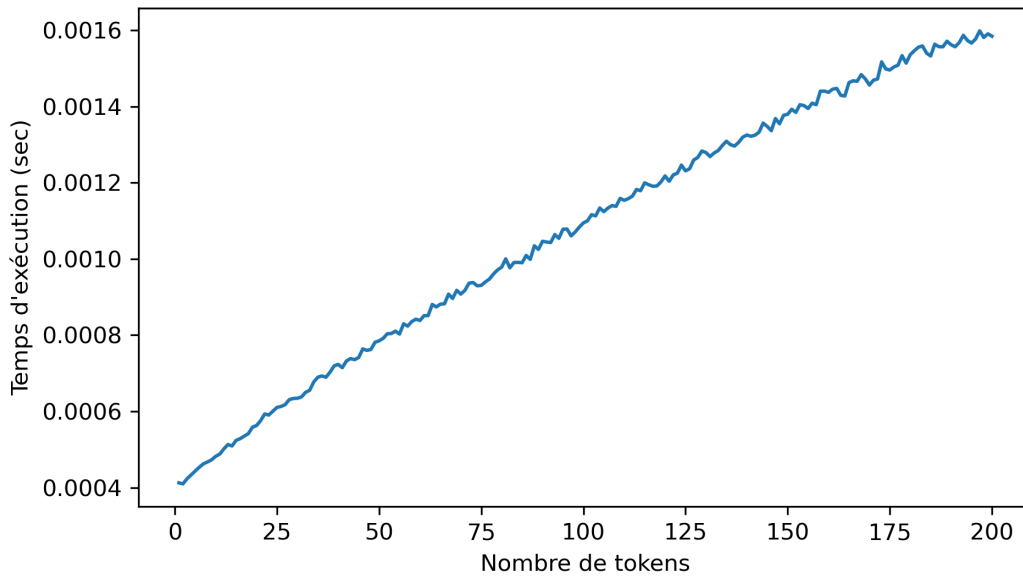


FIGURE 5 – Performances du modèle vectoriel avec matrices creuses en termes de temps d'exécution

## Partie V : Phase d'indexation en chiffres

- **Lecture et prétraitement de la collection CACM** : Cette phase dure **529 ms** (millisecondes). Le résultat est un dictionnaire (document, token) contenant 3204 clés (identificateurs des documents) et pesant **1 706 187 octets** si enregistré dans un fichier. Le fichier CACM de base (cacm.all) pesait lui **2 187 734 octets**.
- **Création du fichier inverse** : Cette phase dure **27 ms**. Le résultat est un dictionnaire (token, document) contenant 11406 clés (tokens) et pesant **760 984 octets** si enregistré dans un fichier.
- **Pondération du fichier inverse** : Cette phase dure **170 ms**. Le résultat est un dictionnaire (token, document) contenant 11406 clés (tokens) et pesant **1 397 739 octets** si enregistré dans un fichier. La taille a légèrement augmenté car la taille d'un entier est moins importante que celle d'un float après avoir été enregistré dans un fichier. Un test rapide permet de le vérifier.
- **Critiques de la phase d'indexation** : La structure de données utilisée pour le modèle vectoriel (dictionnaires de dictionnaires) n'est pas la meilleur en termes de performances pour effectuer rapidement des produits internes. Il existe une meilleur alternative comme décrit dans la Section IV.8

## Partie VI : Évaluation du modèle vectoriel

### VI.1 Modèle vectoriel : Juger la pertinence d'un item

Une fois une requête évaluée, le modèle vectoriel calcule le score de chaque document et les ordonne par rapport à ce score. Mais ensuite, quels sont les documents qu'on décide de garder (documents pertinents) et quels sont les documents qu'on ne doit pas garder (documents non-pertinents) ?

On propose deux méthodes pour résoudre ce problème. Premièrement, la méthode des  $k$ -meilleurs documents qui consiste à retourner seulement les  $k$  documents avec le plus grand score et deuxièmement, la méthode de seuillage qui consiste à retourner seulement les documents qui ont un score au dessus d'un certain seuil  $f$ . On cherche à trouver le meilleur compromis entre précision et rappel, ce qui est équivalent à maximiser le **F1-Score**.

#### VI.1.1 Méthode des $k$ -meilleurs documents

Cette méthode consiste à retourner seulement les  $k$  documents avec le plus grand score. Il suffit maintenant de choisir ce nombre  $k$ . Pour ce faire, on a calculé pour chaque  $k$  allant de 1 à 100 la moyenne des F1-Score. Cette procédure a été réalisée pour toutes les fonctions de similarité et obtient donc, la Figure 6. On trouve que le *produit interne* donne le meilleur résultat quand  $k = 10$  (27.61%), *Dice* et *Jaccard* quand  $k = 7$  (21.55%), *Cosinus* quand  $k = 9$  (24.42%).

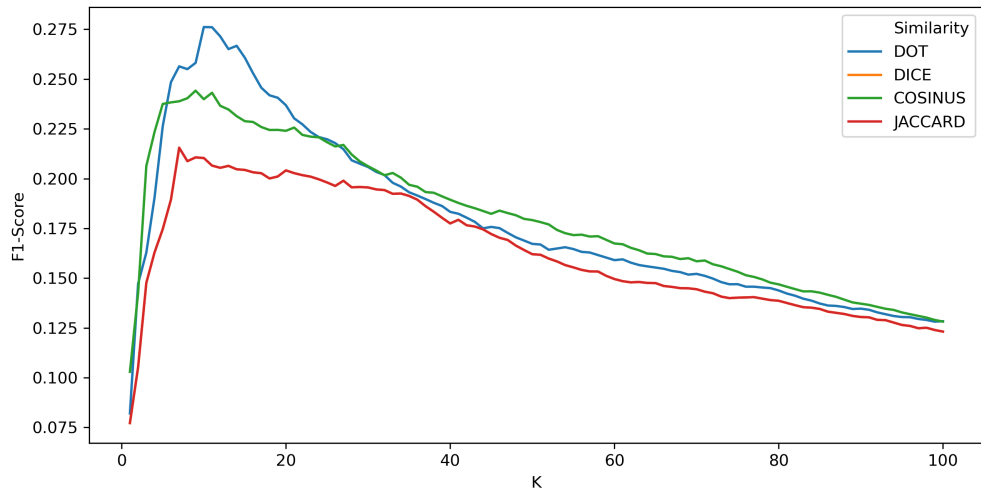


FIGURE 6 – Les performances du modèle vectoriel selon le choix du paramètre  $k$ . Le graphe de DICE est superposé sur le graphe de JACCARD car, ces fonctions de similarité sont équivalentes en termes d'ordre :  $DICE(R, D_1) > DICE(R, D_2)$  si et seulement si  $JACCARD(R, D_1) > JACCARD(R, D_2)$ .

### VI.1.2 Méthode de seuillage

Cette méthode consiste à retourner seulement les documents qui ont un score au dessus d'un certain seuil  $f$ . Il suffit maintenant de déterminer ce seuil  $f$ . On varie  $f$  de 0 à 3 par pas de 0.02 pour l'intervalle  $[0, 1]$  puis par pas de 0.05 pour l'intervalle  $[1, 3]$  et on enregistre à chaque fois le F1-Score (Figure 7). On trouve que le *produit interne* donne le meilleur résultat quand  $f = 1.55$  (26.65%), *Dice* quand  $f = 0.1$  (22.55%), *Jaccard* quand  $f = 0.04$  (20.56%) *Cosinus* quand  $f = 0.16$  (25.17%).

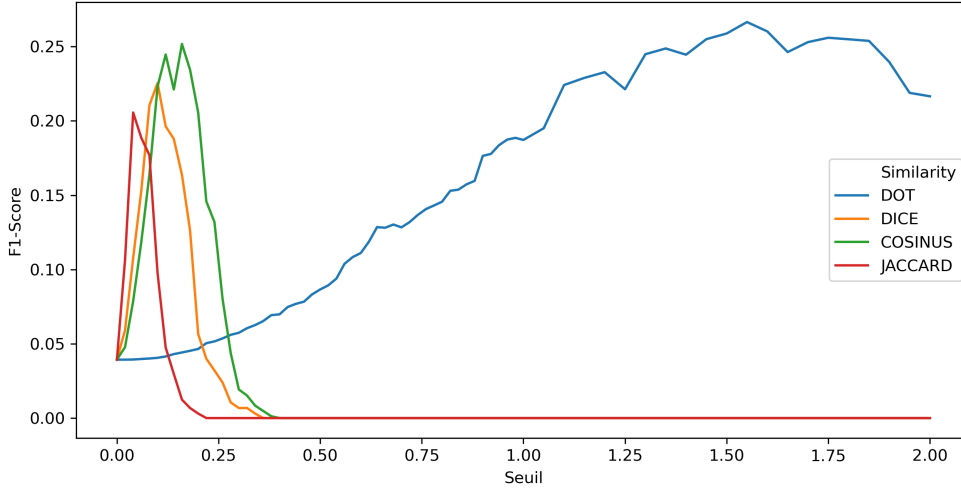


FIGURE 7 – F1-Score du modèle vectoriel selon le choix du paramètre  $f$  (Seuil). Le F1-Score du produit interne (DOT) continue de descendre après Seuil=2.

## VI.2 Combinaison des méthodes

Pour obtenir de meilleurs résultats, on peut combiner les deux méthodes précédentes. Le principe est simple : pour chaque fonction de similarité, on choisie la méthode qui a donnée les meilleurs résultats. La Table 1 montre les performances des méthodes et de leur combinaison. On observe que la similarité la plus performante est le *produit interne*, suivie du *Cosinus*, *Dice* puis *Jaccard*.

## VI.3 Moyenne des courbes précision-rappel interpolées

Une autre manière de comparer les performances des fonctions de similarité est d'utiliser la courbe précision-rappel interpolée (Figure 8). On calcule la courbe pour chaque requête de test, ensuite, on calcule la courbe moyenne. Une courbe au dessus d'une autre veut dire qu'elle est meilleure. Une fois encore, on voit que le produit interne donne de très bons résultats. Cependant, la fonction Cosinus reprend légèrement le dessus vers la deuxième moitié de la courbe.



	Méthode k-meilleurs documents	Méthode de seuillage	Meilleure méthode	Combinaison des méthodes
<b>Produit interne</b>	27.61%	26.65%	K-meilleur	27.61%
<b>Cosinus</b>	24.42%	25.17%	Seuillage	25.17%
<b>Dice</b>	21.55%	22.55%	Seuillage	22.55%
<b>Jaccard</b>	21.55%	20.56%	K-meilleur	21.55%

TABLE 1 – F1-Score des différentes fonction de similarité par rapport à la méthode choisie (Méthode des k-meilleurs documents, Méthode de seuillage et combinaison des deux).

## VI.4 Interface graphique (Bonus/Non demandé)

Pour l'évaluation, on a codé une interface graphique (Figure 9) pour afficher les différentes métriques de performances pour chaque requête et pour chaque fonction de la similarité (produit interne, cosinus, Dice et Jaccard). Ces métriques sont :

- Le rappel,
- La précision,
- Le F1-Score,
- La précision moyenne non-interpolée,
- La précision moyenne interpolée,
- La courbe précision-rappel non-interpolée,
- La courbe précision-rappel interpolée.

Pour calculer ces métriques pour une requête donnée, il suffit de cliquer sur la requête en question dans la table des requêtes.

Il y a trois modes de fonctionnement (Voir la Figure 9) :

- **Utiliser les meilleurs paramètres** : En sélectionnant cette option, l'application calcule les métriques en utilisant la combinaison des méthodes (Vu en Section VI.2).
- **Seuil** : Ici, on laisse le choix à l'utilisateur pour choisir son seuil afin de calculer les métriques. C'est la méthode de seuillage qui est utilisée.
- **K** : On laisse le choix à l'utilisateur pour choisir le  $k$  afin de calculer les métriques. C'est la méthode des  $k$ -meilleurs documents qui est utilisée.

Pour lancer l'application, il faut exécuter le script *app.py* situé dans le dossier *interface*. Il faut préalablement installer PyQt5 avec la commande suivante "pip install PyQt5" (dans la console Windows).

Vous pouvez tester les modèles avec le fichier *Présentation.ipynb* si vous avez Jupyter Notebook installé. Sinon, il y a le script *tester\_les\_modeles.py* dans lequel vous pouvez insérer les requêtes de votre choix pour le modèle booléen et vectoriel.

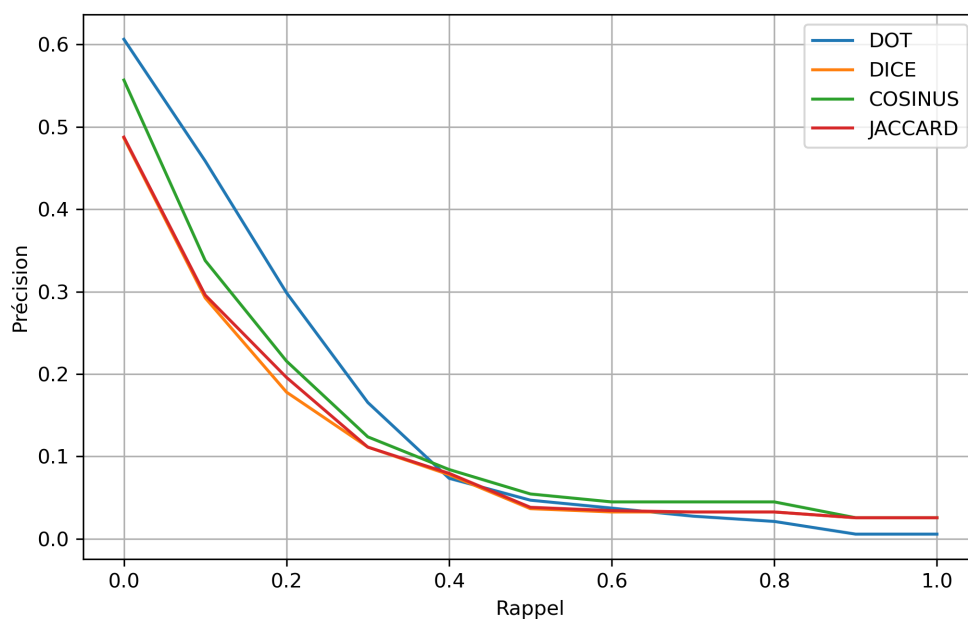


FIGURE 8 – Courbes précision-rappel interpolées des fonctions de similarité utilisant la combinaisons des méthodes vu en Section VI.2.

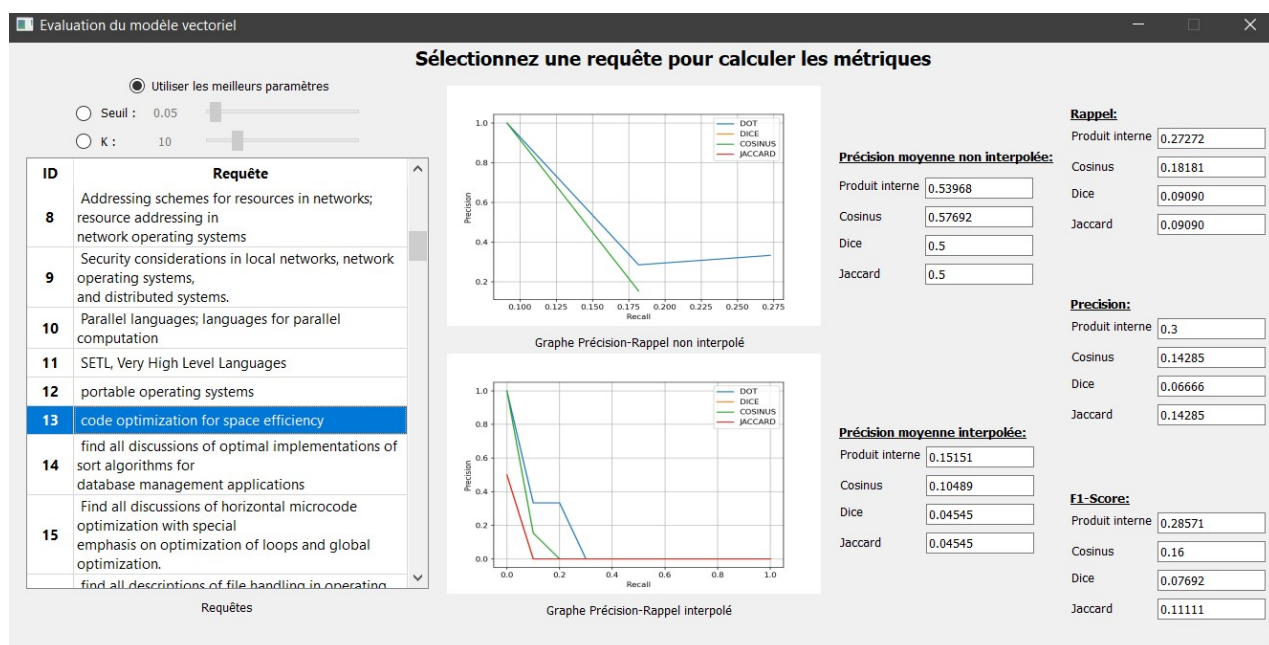


FIGURE 9 – Interface homme-machine