

## Chapitre 3

# Pondération des termes

1

### 1. Introduction

- ✓ La pondération des termes est une étape importante dans le processus d'indexation dans la RI
- ✓ Elle permet d'attribuer une valeur à chaque terme pour mesurer son importance dans le document
- ✓ La 1<sup>ère</sup> valeur qui peut être utilisée c'est la fréquence du terme dans le document
- ✓ Un terme est important dans le document s'il est fréquent dans le document et moins fréquent dans la collection
- ✓ Une autre mesure qui peut être ajoutée à la fréquence du terme pour bien mesurer son poids c'est sa fréquence dans la collection. Cette valeur s'appelle la valeur de discrimination

2

## 2. Pondération basée sur la valeur de discrimination

L'approche TF\*IDF est une approche très connue dans le domaine de la RI, qui se base sur la pondération par valeur de discrimination.

TF : signifie « term frequency » (Fréquence d'un terme dans le document)

IDF : signifie « inverse document frequency » (Fréquence inverse du document)

Donc, TF mesure la fréquence du terme dans le document et IDF mesure la valeur de discrimination du poids du terme.

3

## 3. Calcul de TF

Il existe plusieurs formules pour calculer TF, les plus connues sont :

$$tf_{ij} = \text{fréq}_{ij} \quad tf_{ij} = 1 + \log(\text{fréq}_{ij}) \quad tf_{ij} = \frac{\text{fréq}_{ij}}{\max(\text{fréq}_j)}$$

$$tf_{ij} = \frac{\text{freq}_{ij}}{(K + \text{freq}_{ij})}$$

K introduit pour tenir compte de la longueur des documents

$$tf_{ij} = \frac{\text{fréq.}}{(\text{fréq}_{ij.} + 0.5 + 1.5 * \frac{\text{longueur}_{doc}}{\text{longueur}_{moy}_{doc}})}$$

4

### 3. Calcul de l'IDF

Il y a deux formules pour calculer l'IDF d'un terme  $t_i$  :

$$\text{IDF}_i = \log\left(\frac{N}{n_i}\right) \qquad \text{IDF}_i = \log\left(\frac{N}{n_i} + 1\right)$$

Avec :

N : le nombre de documents de la collection,  
 $n_i$  le nombre de documents contenant le terme  $t_i$

#### Remarque :

La 1<sup>ère</sup> formule retourne zéro pour un terme qui existe dans tous les documents de la collection. Pour éviter ce problème nous utilisons la 2<sup>ème</sup> formule dans la suite de nos cours.

5

### 5. Calcul du poids d'un terme par TF\*IDF

Pour calculer le poids d'un terme par TF\*IDF, il suffit de choisir une formule de TF est al multiplier par une formule de l'IDF.

**Exercice :** Soit la collection suivante de 3 documents :

D1 : { langage de programmation python est très utilisé pour le traitement de texte }

D2 : { le langage JAVA est basé sur le langage C++ }

D3 : { un langage de programmation est un langage utilisé pour traduire un algorithme en un programme }

Motsvides: { de, est, très, pour, le, un, en }

Donner le fichier inverse de la collection avec la formule :

$$\text{poids}(t_i, d_j) = (\text{freq}(t_i, d_j) / \max(\text{freq}(d_j))) * \log(N/n_i + 1)$$

6

**Solution :**

<b>Terme \ doc</b>	<b>D1</b>	<b>D2</b>	<b>D3</b>
<b>langage</b>	0.301	0.301	0.301
<b>programmation</b>	0.397		0.198
<b>python</b>	0.602		
<b>utilisé</b>	0.397		0.198
<b>traitement</b>	0.602		
<b>texte</b>	0.602		
<b>java</b>		0.301	
<b>basé</b>		0.301	
<b>c++</b>		0.301	
<b>traduire</b>			0.301
<b>algorithme</b>			0.301
<b>programme</b>			0.301

7