

Partie 4

Recherche positionnelle et approchée

Vincent Labatut

Laboratoire Informatique d'Avignon – LIA EA 4128

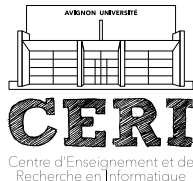
vincent.labatut@univ-avignon.fr

2019/20

M2 ILSEN

UE Ingénierie du document et de l'information

UCE3 Indexation & Recherche d'information



Plan de la séance

- 1 Traitement de groupes de mots
 - Groupes de mots
 - Index positionnel
- 2 Recherche approchée par jokers
 - Recherche approchée
 - Index Permuterm
 - Index de k -grammes
- 3 Recherche approchée par correction
 - Comparaison de chaînes de caractères
 - Correction phonétique & Soundex

Section 1

Traitement de groupes de mots

Traitement de groupes de mots

Problématique

- **Limitation :**
 - De nombreux concepts sont décrits par des groupes de mots
 - Ex. : Stanford University
 - Traitement sous forme de conjonction : sensible à The inventor **Stanford** Ovshinsky never went to **university**
 - Problème connexe : proximité des tokens
- **Solution :**
 - Requête d'**expression exacte** (phrase query)
→ index actuel insuffisant
- Deux **modifications** proposées :
 - Manipuler des **paires de mots**
 - Utiliser un **index positionnel**

Traitement de groupes de mots

Notion de k -gramme

k -gramme

Sous-séquence de k éléments dans une séquence plus longue.

En NLP, il s'agit :

- Soit de k -grammes **de mots**, par ex. :
 - Séquence : il pleuvait hier il fera beau demain
 - **Bigrammes** : il pleuvait ; pleuvait hier ; hier il ; il fera ; fera beau...
 - **Trigrammes** : il pleuvait hier ; pleuvait hier il ; hier il fera...
 - **Tetra**grammes : il pleuvait hier il ; pleuvait hier il fera...
- Soit de k -grammes **de caractères**, par ex. :
 - Séquence : sapristi
 - **Bigrammes** : sa ; ap ; pr ; ri ; st ; ti
 - **Trigrammes** : sap ; apr ; pri ; ris ; ist ; sti
 - **Tetra**grammes : sapr ; apri ; pris ; rist ; isti
 - **Penta**grammes : sapri ; apris ; prist ; risti

Traitement de groupes de mots

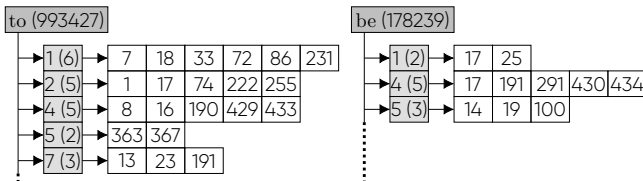
Paires de mots

- Principe :
 - Tokénisation par **bigrammes** de mots
 - Pour corpus **et** requêtes
 - Ex. : Friends, Romans, Countrymen
 - friends romans; romans countrymen
 - **Entrées** de l'index = **couples** de mots
- Limitations :
 - **Faux positifs**
 - Ex. : Stanford University Palo Alto
 - stanford university; university palo; palo alto
 - Contenir les 3 couples \neq contenir la séquence recherchée
 - Index beaucoup plus **grand**

Traitement de groupes de mots

Notion d'index positionnel

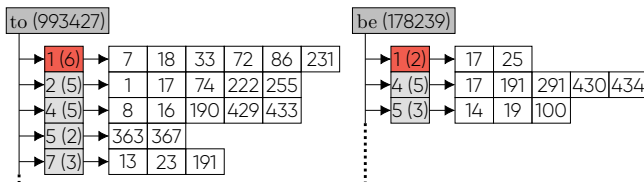
- Non-positionnel (jusqu'à présent) : postings = docID
- **Positionnel** : docID + liste de positions
 - Position exprimée en termes de numéro de token dans le texte
- Traitement des requêtes :
 - Similaire à l'approche classique...
 - ...mais en considérant en plus les positions **relatives** des mots
- Exemple : **to**₁ **be**₂ or₃ not₄ **to**₅ **be**₆



Traitement de groupes de mots

Notion d'index positionnel

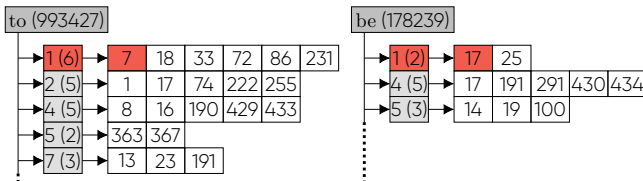
- Non-positionnel (jusqu'à présent) : postings = docID
- **Positionnel** : docID + liste de positions
 - Position exprimée en termes de numéro de token dans le texte
- Traitement des requêtes :
 - Similaire à l'approche classique...
 - ...mais en considérant en plus les positions **relatives** des mots
- Exemple : **to**₁ **be**₂ or₃ not₄ **to**₅ **be**₆



Traitement de groupes de mots

Notion d'index positionnel

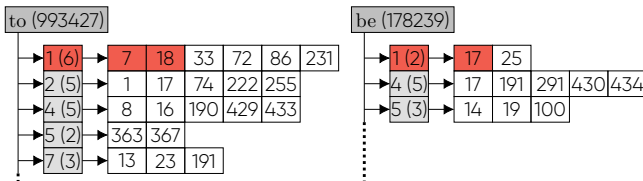
- Non-positionnel (jusqu'à présent) : postings = docID
- **Positionnel** : docID + liste de positions
 - Position exprimée en termes de numéro de token dans le texte
- Traitement des requêtes :
 - Similaire à l'approche classique...
 - ...mais en considérant en plus les positions **relatives** des mots
- Exemple : **to₁ be₂ or₃ not₄ to₅ be₆**



Traitement de groupes de mots

Notion d'index positionnel

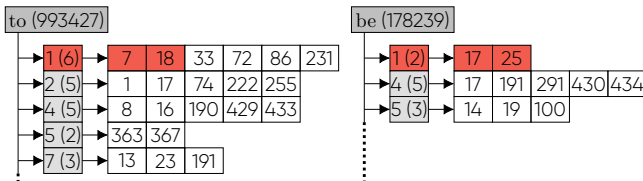
- Non-positionnel (jusqu'à présent) : postings = docID
- **Positionnel** : docID + liste de positions
 - Position exprimée en termes de numéro de token dans le texte
- Traitement des requêtes :
 - Similaire à l'approche classique...
 - ...mais en considérant en plus les positions **relatives** des mots
- Exemple : **to**₁ **be**₂ or₃ not₄ **to**₅ **be**₆



Traitement de groupes de mots

Notion d'index positionnel

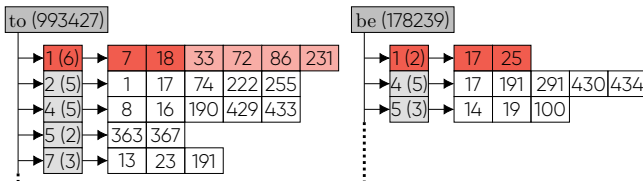
- Non-positionnel (jusqu'à présent) : postings = docID
- **Positionnel** : docID + liste de positions
 - Position exprimée en termes de numéro de token dans le texte
- Traitement des requêtes :
 - Similaire à l'approche classique...
 - ...mais en considérant en plus les positions **relatives** des mots
- Exemple : **to**₁ **be**₂ or₃ not₄ **to**₅ **be**₆



Traitement de groupes de mots

Notion d'index positionnel

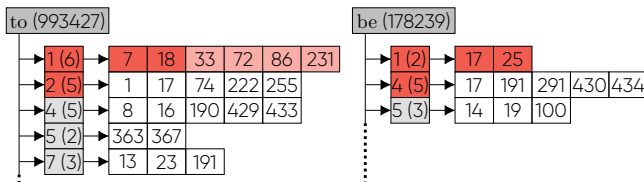
- Non-positionnel (jusqu'à présent) : postings = docID
- **Positionnel** : docID + liste de positions
 - Position exprimée en termes de numéro de token dans le texte
- Traitement des requêtes :
 - Similaire à l'approche classique...
 - ...mais en considérant en plus les positions **relatives** des mots
- Exemple : **to**₁ **be**₂ or₃ not₄ **to**₅ **be**₆



Traitement de groupes de mots

Notion d'index positionnel

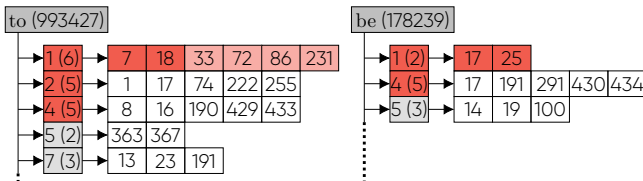
- Non-positionnel (jusqu'à présent) : postings = docID
- **Positionnel** : docID + liste de positions
 - Position exprimée en termes de numéro de token dans le texte
- Traitement des requêtes :
 - Similaire à l'approche classique...
 - ...mais en considérant en plus les positions **relatives** des mots
- Exemple : **to₁ be₂ or₃ not₄ to₅ be₆**



Traitement de groupes de mots

Notion d'index positionnel

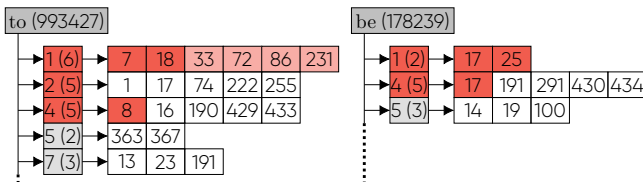
- Non-positionnel (jusqu'à présent) : postings = docID
- **Positionnel** : docID + liste de positions
 - Position exprimée en termes de numéro de token dans le texte
- Traitement des requêtes :
 - Similaire à l'approche classique...
 - ...mais en considérant en plus les positions **relatives** des mots
- Exemple : **to₁ be₂ or₃ not₄ to₅ be₆**



Traitement de groupes de mots

Notion d'index positionnel

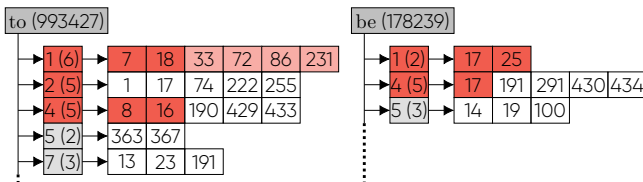
- Non-positionnel (jusqu'à présent) : postings = docID
- **Positionnel** : docID + liste de positions
 - Position exprimée en termes de numéro de token dans le texte
- Traitement des requêtes :
 - Similaire à l'approche classique...
 - ...mais en considérant en plus les positions **relatives** des mots
- Exemple : **to₁ be₂ or₃ not₄ to₅ be₆**



Traitement de groupes de mots

Notion d'index positionnel

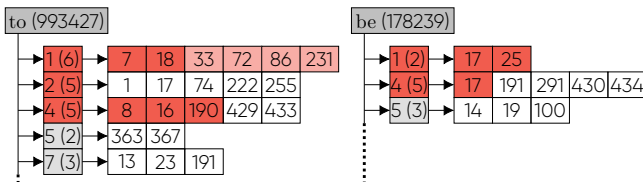
- Non-positionnel (jusqu'à présent) : postings = docID
- **Positionnel** : docID + liste de positions
 - Position exprimée en termes de numéro de token dans le texte
- Traitement des requêtes :
 - Similaire à l'approche classique...
 - ...mais en considérant en plus les positions **relatives** des mots
- Exemple : **to**₁ **be**₂ or₃ not₄ **to**₅ **be**₆



Traitement de groupes de mots

Notion d'index positionnel

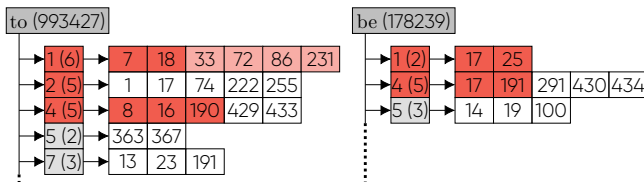
- Non-positionnel (jusqu'à présent) : postings = docID
- **Positionnel** : docID + liste de positions
 - Position exprimée en termes de numéro de token dans le texte
- Traitement des requêtes :
 - Similaire à l'approche classique...
 - ...mais en considérant en plus les positions **relatives** des mots
- Exemple : **to**₁ **be**₂ or₃ not₄ **to**₅ **be**₆



Traitement de groupes de mots

Notion d'index positionnel

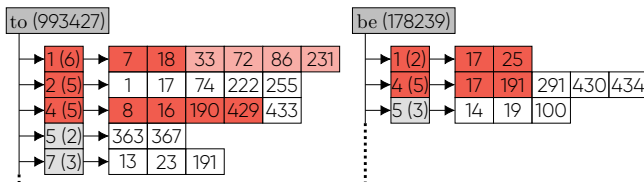
- Non-positionnel (jusqu'à présent) : postings = docID
- **Positionnel** : docID + liste de positions
 - Position exprimée en termes de numéro de token dans le texte
- Traitement des requêtes :
 - Similaire à l'approche classique...
 - ...mais en considérant en plus les positions **relatives** des mots
- Exemple : **to₁ be₂ or₃ not₄ to₅ be₆**



Traitement de groupes de mots

Notion d'index positionnel

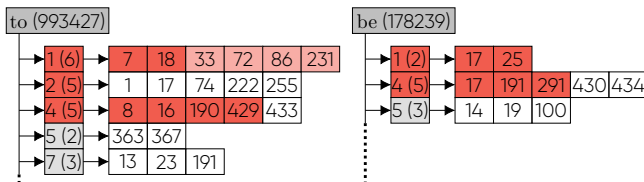
- Non-positionnel (jusqu'à présent) : postings = docID
- **Positionnel** : docID + liste de positions
 - Position exprimée en termes de numéro de token dans le texte
- Traitement des requêtes :
 - Similaire à l'approche classique...
 - ...mais en considérant en plus les positions **relatives** des mots
- Exemple : **to**₁ **be**₂ or₃ not₄ **to**₅ **be**₆



Traitement de groupes de mots

Notion d'index positionnel

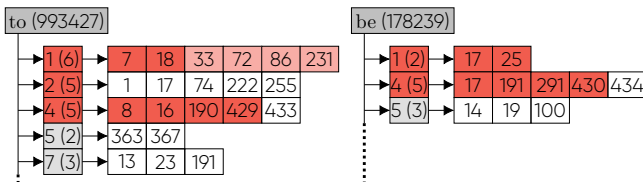
- Non-positionnel (jusqu'à présent) : postings = docID
- **Positionnel** : docID + liste de positions
 - Position exprimée en termes de numéro de token dans le texte
- Traitement des requêtes :
 - Similaire à l'approche classique...
 - ...mais en considérant en plus les positions **relatives** des mots
- Exemple : **to**₁ **be**₂ or₃ not₄ **to**₅ **be**₆



Traitement de groupes de mots

Notion d'index positionnel

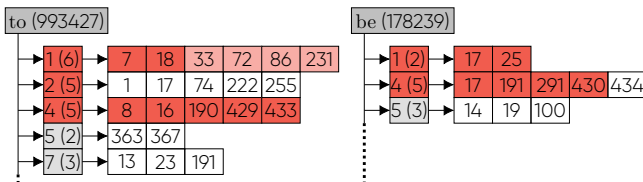
- Non-positionnel (jusqu'à présent) : postings = docID
- **Positionnel** : docID + liste de positions
 - Position exprimée en termes de numéro de token dans le texte
- Traitement des requêtes :
 - Similaire à l'approche classique...
 - ...mais en considérant en plus les positions **relatives** des mots
- Exemple : **to**₁ **be**₂ or₃ not₄ **to**₅ **be**₆



Traitement de groupes de mots

Notion d'index positionnel

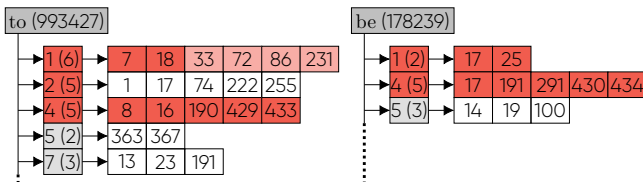
- Non-positionnel (jusqu'à présent) : postings = docID
- **Positionnel** : docID + liste de positions
 - Position exprimée en termes de numéro de token dans le texte
- Traitement des requêtes :
 - Similaire à l'approche classique...
 - ...mais en considérant en plus les positions **relatives** des mots
- Exemple : **to**₁ **be**₂ or₃ not₄ **to**₅ **be**₆



Traitement de groupes de mots

Notion d'index positionnel

- Non-positionnel (jusqu'à présent) : postings = docID
- **Positionnel** : docID + liste de positions
 - Position exprimée en termes de numéro de token dans le texte
- Traitement des requêtes :
 - Similaire à l'approche classique...
 - ...mais en considérant en plus les positions **relatives** des mots
- Exemple : **to**₁ **be**₂ or₃ not₄ **to**₅ **be**₆



Traitement de groupes de mots

Propriétés des index positionnels

- Plus **précis** :
 - Renvoie une liste de **positions** dans le texte
 - (par opposition à une liste de **documents**)
- **Opérateur de proximité** : même approche
 - Ex. : employment /3 place
 - I.e. : place situé 3 mots (ou moins) après employment
- **Propriétés** de l'index positionnel
 - Occupe plus d'**espace** mémoire
 - 2 ordres de grandeurs
 - Accès plus **lent** :
 - Linéaire relativement au nombre de tokens
 - (et non plus au nombre de documents)

Section 2

Recherche approchée par jokers

Recherche approchée par jokers

Limitations de la recherche exacte

- **Limitation** : rigidité orthographique
 - Requête **exacte** : pas de faute d'orthographe
 - Requête **complète** : mots décrits en entier
 - Correspondance **exacte** : même orthographe dans l'index
- **Solutions** :
 - **Corriger** les erreurs dans les requêtes
 - Utiliser des **jokers** (*, ?, etc.)
 - Comparer **phonétiquement**
 - nécessité d'une structure de données **appropriée** pour l'index

Recherche approchée par jokers

Structure du lexique

- Lexique : 2 alternatives
 - Table de hachage
 - Fonction de hachage : termes \rightarrow entiers
 - Recherche en temps **constant**
 - Impossible de trouver facilement des variantes lexicales du terme
 - Pas de recherche par préfixe
 - Fonction peut être insuffisante si le lexique grandit
 - Arbres de recherche
 - Label = préfixe (mot pour les feuilles)
 - Recherche en temps **logarithmique**
 - Possible de trouver des variantes
 - Possible de chercher par préfixe
 - Mais insertion d'élément coûteuse (\rightarrow *B-trees*)
 - Occupation mémoire supérieure à hachage

Recherche approchée par jokers

Notion de joker

Métacaractère, ou joker (eng : wild card)

Caractère spécial se substituant à une **sous-séquence** de caractères dans une **requête**.

La **syntaxe** exacte dépend du moteur de recherche, mais le caractère * est souvent utilisé pour représenter n'importe quelle sous-séquence (y compris \emptyset)

Les **jokers** sont utiles dans les cas suivants :

- Utilisateur **incertain** de l'orthographe
 - Ex. : silhouette vs. sihlouette
- Utilisateur **sait** qu'il existe plusieurs orthographes
 - Ex. : clé vs. clef
- Utilisateur **visé** plusieurs mots proches
 - Ex. : algorithmique et algorithme

Recherche approchée par jokers

Exploitation de l'arbre de recherche

- Recherche de **préfixe** sur arbre : part^*
 - Naviguer jusqu'au nœud correspondant au préfixe
 - Parcourir en partant de ce nœud, jusqu'à rencontrer un nœud ne correspondant plus au préfixe
 - Ex. : pour mon^* : aller à mon- et parcourir jusqu'à moo-
- Recherche de **suffixe** : $^*\text{lier}$
 - Même chose sur un arbre construit à l'envers (ordre alphabétique ascendant)
- Recherche d'**infixe** : part^*lier
 - Calculer l'ensemble des termes pour part^*
 - Calculer l'ensemble des termes pour $^*\text{lier}$
 - Prendre leur intersection

→ Approche coûteuse

Recherche approchée par jokers

Notion de Permuterm

- Alternative : index **Permuterm**
 - **Principe** : transformer la requête de manière à placer le * à la fin
 - Caractère \$ marque la **fin** du terme
 - Lexique contient toutes les **rotations** d'un terme
 - Ex. : racine \rightarrow racine\$, \$racine, e\$racin, ne\$raci, ine\$rac, cine\$ra, acine\$r
- **Transformation** de la requête :
 - Rotation amenant l'éventuel * à la fin
 - racine \rightarrow racine\$
 - *ine \rightarrow ine\$*; rac* \rightarrow \$rac*
 - r*e \rightarrow e\$r*; *ci* \rightarrow ci*
- Bilan :
 - **Plus rapide** qu'approche précédente
 - Mais lexique **beaucoup plus grand**
 - Ex. : $\times 10$ pour un corpus en anglais

Recherche approchée par jokers

Index de k -grammes de caractères

- **Alternative** : index de k -gramme de caractères
 - On utilise \$ pour noter le début ou la fin d'un terme.
 - Ex. pour $k = 3$: chat \rightarrow \$ch, cha, hat, at\$
 - Index supplémentaire :
 - Lexique = tous les k -grammes des termes originaux
 - Chacun pointe vers la liste des termes originaux le contenant
 - Ex. : ntr pointe vers **antre**, **entre**, **cintre**...
- **Utilisation** :
 - ra^* \rightarrow on cherche le trigramme \$ra
 - $*ne$ \rightarrow on cherche le trigramme ne\$
 - ra^*ne \rightarrow intersection des résultats de \$ra et ne\$
- **Complication** pour red^* :
 - \$re et red vont produire (entre autres) **retired**
 - \rightarrow besoin d'un filtrage supplémentaire (simple comparaison)

Section 3

Recherche approchée par correction

Recherche approchée par correction

Distance de Levenshtein

Distance édit

Permet de comparer deux **séquences** en déterminant le nombre **minimal** d'**opérations** nécessaires pour passer de l'une à l'autre.

Distance de Levenshtein

La distance de **Levenshtein** [Lev66] est une distance **édit** autorisant les opérations **insérer**, **supprimer** et **remplacer**.

Utilisée en NLP pour comparer les chaînes de caractères :

- dog vs. do $\rightarrow 1$
- cat vs. cart $\rightarrow 1$
- cat vs. act $\rightarrow 2$
- niche vs. chiens $\rightarrow 5$

Recherche approchée par correction

Notations

- Calcul de la distance de Levenshtein :
 - Méthode récursive naïve [Lev66]
 - Algorithme de Wagner–Fischer [WF74]
- Notations :
 - $|s|$: longueur de la chaîne s , ex. $|\text{"truc"}| = 4$
 - $\text{init}(s)$: chaîne privée de son dernier caractère
 - Ex. $\text{init}(\text{"truc"}) = \text{"tru"}$
 - $\text{last}(s)$: dernier caractère
 - Ex. $\text{last}(\text{"truc"}) = \text{"c"}$
 - $[P]$: crochets d'Iverson
 - Ex. $[1 > 4] = 0$, mais $[1 < 4] = 1$

Recherche approchée par correction

Méthode naïve pour calculer Levenshtein

```
1 Function lev( $s_1, s_2$ ) :  
    Input :  $s_1, s_2$  : String  
    Output :  $d$  : Integer  
  
2    if  $s_1 = \emptyset$  then  
3         $d \leftarrow |s_2|$   
4    else if  $s_2 = \emptyset$  then  
5         $d \leftarrow |s_1|$   
6    else  
7         $d \leftarrow \min(\text{lev}(\text{init}(s_1), s_2) + 1 ;$   
8             $\text{lev}(s_1, \text{init}(s_2)) + 1 ;$   
9             $\text{lev}(\text{init}(s_1), \text{init}(s_2)) + [\text{last}(s_1) \neq \text{last}(s_2)])$   
10    end if  
11    return  $d$   
12 return
```

Algorithme 1 – Méthode naïve pour calculer Levenshtein [Lev66].

Recherche approchée par correction

Limites de la méthode naïve

- Complexité **temporelle** : $O(3^{\min(|s_1|; |s_2|)})$
 - Cause de la complexité élevée : opérations redondantes
 - (de nombreux calculs sont effectués plusieurs fois)
- **Amélioration** : Algorithme de Wagner-Fischer
 - Méthode à base de **programmation dynamique**
- **Principe** :
 - Utiliser une matrice entière **D** de taille $|s_1| \times |s_2|$ pour stocker les résultats des calculs intermédiaires
 - Initialisée avec des valeurs négatives
 - Chaque élément **D**_{ij} contient la distance entre :
 - La sous-chaîne contenant les *i* premiers caractères de s_1
 - Et celle contenant les *j* premiers caractères de s_2
 - Si on a besoin d'une distance déjà calculée, on la récupère dans la matrice au lieu de la recalculer

Recherche approchée par correction

Algorithme de Wagner-Fisher

```
1 Function lev( $s_1, s_2, \mathbf{D}$ ) :  
   Input :  $s_1, s_2$  : String,  $\mathbf{D}$  : Integer Matrix  
   Output :  $d$  : Integer  
2   if  $s_1 = \emptyset$  then  
3      $d \leftarrow |s_2|$   
4   else if  $s_2 = \emptyset$  then  
5      $d \leftarrow |s_1|$   
6   else if  $\mathbf{D}_{|s_1|, |s_2|} \geq 0$  then  
7      $d \leftarrow \mathbf{D}_{|s_1|, |s_2|}$   
8   else  
9      $d \leftarrow \min(\text{lev}(\text{init}(s_1), s_2) + 1 ;$   
10       $\text{lev}(s_1, \text{init}(s_2)) + 1 ;$   
11       $\text{lev}(\text{init}(s_1), \text{init}(s_2)) + [\text{last}(s_1) \neq \text{last}(s_2)])$   
12      $\mathbf{D}_{|s_1|, |s_2|} \leftarrow d$   
13   end if  
14   return  $d$   
15 return
```

Algorithme 2 – Algorithme de Wagner-Fischer [WF74].

Recherche approchée par correction

Utilisation de Levenshtein pour la recherche approchée

- Correction :
 - Remplacer par le terme le plus proche dans le lexique

- Recherche exhaustive = trop coûteuse

→ Approche **heuristique**

- Restriction aux mots de même initiale
- Restriction aux termes commençant par une rotation du mot
- Restriction aux termes ayant un certain nombre de k -grammes en commun
 - Déterminer tous les k -grammes du mot requête
 - Obtenir tous les termes ayant des k -grammes communs (via l'index k -gramme des jokers)
 - Filtrer ceux qui ont n k -grammes communs (par intersection des ensembles obtenus)
 - Autre possibilité : **coefficient de Jaccard** n/U

Recherche approchée par correction

Correction phonétique

- Correction **phonétique** :
 - Requête tapée "à l'oreille"
 - Ex. d'utilisation : recherche de noms propres
 - Principe :
 - Chaque mot est remplacé par un **code** représentant sa prononciation
 - Deux **homophones** (mots de même prononciation) ont le même code
 - On compare le code du mot requête à ceux des termes
- Exemple : code **Soundex** [Rep02]
 - Composé d'1 lettre et 3 chiffres : **A123**

Recherche approchée par correction

Algorithme du Soundex

- **Algorithme :**

- ① Garder la première lettre du mot
- ② Remplacer toutes les voyelles et muettes (H, W) par zéro
- ③ Remplacer toutes les consonnes par leur numéro de classe :
 - B, F, P, V = 1
 - C, G, J, K, Q, S, X, Z = 2
 - D, T = 3
 - L = 4
 - N, M = 5
 - R = 6
- ④ Remplacer les valeurs consécutives égales par une seule occurrence
- ⑤ Supprimer les zéros
- ⑥ Renvoyer les 4 premiers caractères
- ⑦ Compléter par des zéros, si nécessaire

- **Exemple :** Hermann → H655

Section 4

Conclusion

Concepts abordés dans cette partie

- Index positionnel
- Opérateur de proximité
- Distance édit
- Distance de Levenshtein
- Algorithme de Wagner-Fisher
- Joker, métacaractère
- Permuterm
- k -gramme de mots
- k -gramme de caractères
- Soundex

Lectures recommandées

- [MRS08] *Introduction to Information Retrieval*, chapitre 3.
- [BCC10] *Information Retrieval : Implementing and Evaluating Search Engines*, chapitre 2.
- [BR11] *Modern Information Retrieval : The Concepts and Technology behind Search*, chapitre 9.
- [AG13] *Recherche d'information - Applications, modèles et algorithmes*, chapitre 3.
- [CMS15] *Search Engines : Information Retrieval in Practice*, chapitre 5.

Références bibliographiques I

- [AG13] M.-R. Amini et É. Gaussier. *Recherche d'information – Applications, modèles et algorithmes*. Paris, FR : Eyrolles, 2013. url : <https://www.eyrolles.com/Informatique/Livre/recherche-d-information-9782212673760/>.
- [BR11] R. Baeza-Yates et B. Ribeiro-Neto. *Modern Information Retrieval : The Concepts and Technology behind Search*. 2nd Edition. Boston, USA : Addison Wesley Longman, 2011. url : <http://people.ischool.berkeley.edu/~hearst/irbook/>.
- [BCC10] S. Büttcher, C. L. A. Clarke et G. V. Cormack. *Information Retrieval : Implementing and Evaluating Search Engines*. Cambridge, USA : MIT Press, 2010. url : <http://www.ir.uwaterloo.ca/book/>.
- [CMS15] W. B. Croft, D. Metzler et T. Strohman. *Search Engines : Information Retrieval in Practice*. Pearson, 2015. url : <http://www.search-engines-book.com/>.

Références bibliographiques II

- [Lev66] V. I. Levenshtein. « Binary codes capable of correcting deletions, insertions, and reversals ». In : *Soviet Physics Doklady* 10.8 (1966), p. 707–710. url : <https://nymity.ch/sybilhunting/pdf/Levenshtein1966a.pdf>.
- [MRS08] C. D. Manning, P. Raghavan et H. Schütze. *Introduction to Information Retrieval*. New York, USA : Cambridge University Press, 2008. url : <http://www-nlp.stanford.edu/IR-book/>.
- [Rep02] D. J. Repici. *Understanding Classic SoundEx Algorithms*. Creativyst. 2002. url : <http://www.creativyst.com/Doc/Articles/SoundEx1/SoundEx1.htm>.
- [WF74] R. A. Wagner et M. J. Fischer. « The String-to-String Correction Problem ». In : *Journal of the Association for Computing Machinery* 21.1 (1974), p. 168–173. doi : [10.1145/321796.321811](https://doi.org/10.1145/321796.321811).