

Un texte , une modélisation

Laurent Dumas

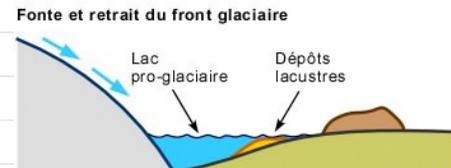
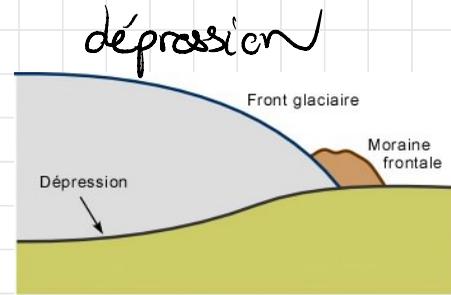
Texte 4: une requête bibliographique

* Objectif: classer la pertinence de documents par rapport à une requête par mots clés.

* Exemple étudié: classer la pertinence de 12 documents par rapport à une requête polysémique.

* Outils mathématiques:

Algèbre linéaire (décomposition en valeurs singulières).



Etape 1 : modélisation d'une requête bibliographique

* On dispose de :

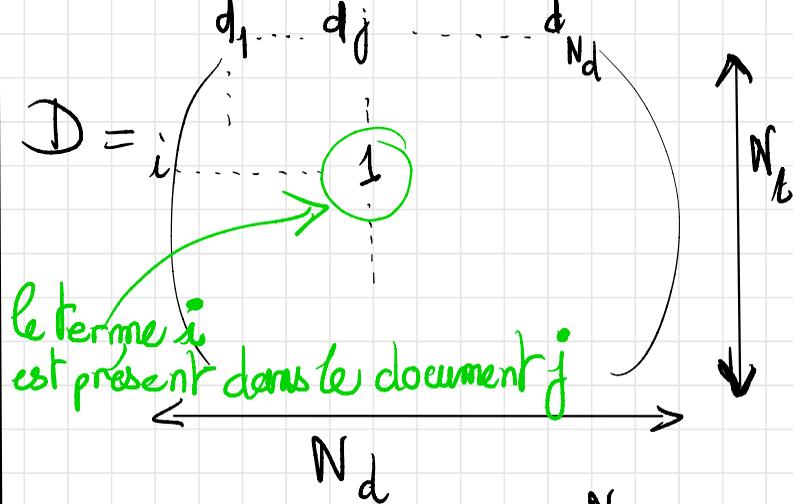
→ N_d documents
contenant :

→ N_t termes

et on cherche la pertinence d'une
requête choisie parmi ces termes.

* On construit :

→ $D \in \mathbb{R}^{N_t \times N_d}$ (IR) la matrice
termes - documents :



→ la requête $q \in \mathbb{R}^{N_t}$:

$q =$

le terme i est
présent dans la
requête q

→ la réponse sous la forme d'un
vecteur de score : $s \in \mathbb{R}^{N_d}$:

$$s_j = \frac{\langle q, d_j \rangle}{\|q\| \cdot \|d_j\|}$$

Etape 3 : réduction de rang et calcul d'un score approché

* Un aspect important de la décomposition SVD est le suivant :

Théorème 2 : la meilleure approximation de A parmi les matrices de rang k est donnée par :

$$A_k = U \begin{pmatrix} \sigma_1 & & & & \\ & \ddots & & & \\ & & \sigma_k & & \\ & & & 0 & \\ & & & & \ddots \\ & & & & & 0 \end{pmatrix}^k V$$

(au sens de la norme $\| \cdot \|_2$)

* A partir de la décomposition SVD de la matrice $D \in \mathbb{R}^{N_t \times N_d}$

on calcule une matrice approchée :

$$D_k = \underbrace{U_k}_{N_t \times k} \underbrace{\begin{pmatrix} \sigma_1 & & & 0 \\ & \ddots & & \\ & & \sigma_k & \\ & & & 0 \end{pmatrix}}_{\Sigma_k} \underbrace{V_k}_{k \times N_d}$$

puis un score approché \tilde{s} :

$$\tilde{s}_j = \frac{\langle q, D_k e_j \rangle}{\|q\| \cdot \|D_k e_j\|}$$
$$= \frac{\langle {}^t U_k q, \Sigma_k {}^t V_k e_j \rangle}{\|{}^t U_k q\|_2 \|\Sigma_k {}^t V_k e_j\|_2}$$

Etape 4 : simulation avec le logiciel Scilab

DOCUMENTS ÉCONOMIQUES

- Doc. 1 : chômage, impôts, dépression, commerce
- Doc. 2 : économie, marché, commerce
- Doc. 3 : commerce, marché, production
- Doc. 4 : dépression, liquidation, emploi, redressement
- Doc. 5 : indemnisation, chômage, liquidation
- Doc. 6 : commerce, marché, prix, emploi
- Doc. 7 : bénéfices, indemnisation, chômage

DOCUMENTS GÉOLOGIQUES

- Doc. 8 : bassin, faille, dérive
- Doc. 9 : dépression, faille
- Doc. 10 : bassin, drainage, vallée
- Doc. 11 : dépression, drainage, érosion
- Doc. 12 : bassin, drainage, volcan

Gn prend $N_d = 12$, $N_r = 20$

* Gn teste les requêtes :

→ requête 1 : [dépression, commerce]

↘ requête 2 : [emploi]

Pour aller plus loin :

→ aspects algorithmiques de la SVD

→ lien avec la recherche de valeurs propres.

Référence :

www.agreg.org

(texte de modélisation)