

TD 2

Exo 4:

- Soit la collection (C) suivante :

D_1 = "été soleil brille température augmente"

D_2 = "hiver neige tombe température diminue"

D_3 = "soleil fait fondre neige"

D_4 = "montagne neige tombe été"

- Considérant le modèle booléen, répondre aux requêtes suivantes :

R_{t_i} = "ensbl de documents qui contient le terme t_i "

$$R_{été} = \{D_1, D_4\}$$

$$R_{neige} = \{D_2, D_3, D_4\}$$

$$R_{soleil} = \{D_1, D_3\}$$

$$R_{montagne} = \{D_4\}$$

$$R_{température} = \{D_1, D_2\}$$

$$q_1 = \text{soleil} \wedge \text{neige}$$

$$R_{q_1} = R_{soleil} \cap R_{neige} = \{D_1, D_3\} \cap \{D_2, D_3, D_4\} = \{D_3\}$$

$$q_2 = \text{température} \vee \text{été}$$

$$R_{q_2} = R_{température} \cup R_{été} = \{D_1, D_2\} \cup \{D_1, D_4\} = \{D_1, D_2, D_4\}$$

$$q_3 = \text{neige} \wedge \neg (\text{montagne} \wedge \text{été})$$

$$R_{q_3} = R_{neige} \cap \neg (R_{montagne} \cap R_{été}) = \{D_2, D_3, D_4\} \cap \neg (\{D_4\} \cap \{D_1, D_4\})$$

$$= \{D_2, D_3, D_4\} \cap (\{D_4\})^c = \{D_2, D_3, D_4\} \cap \{D_1, D_2, D_3\}$$

$$= \{D_2, D_3\}$$

Exo 2:

Avec la même collection. Bon s'agissant le modèle victorien

1) Construire la matrice des poids doc-termes. $N=4$

	été	soleil	brûle	temp	augmente	hiver	neige	tombe	clin	fait	fond	lente
ni	2	2	1	2	1	1	3	2	1	1	1	1
id_1	1	1	2	1	2	2	0.44	1	2	2	2	2
tg_1	1	1	1	1	1	0	0	0	0	0	0	0
tg_2	0	0	0	1	0	1	1	1	1	0	0	0
tg_3	0	1	0	0	0	0	1	0	0	1	1	0
tg_4	0	0	0	0	0	0	1	1	0	0	0	1
W_{i1}	1	1	2	1	2	0	0	0	0	0	0	0
W_{i2}	0	0	0	1	0	2	0.44	1	2	0	0	0
W_{i3}	0	1	0	0	0	0	0.44	0	0	2	2	0
W_{i4}	1	0	0	0	0	0	0.44	1	0	0	0	2
w_q	1	1	1	1	1	1	1	1	1	1	1	1

2) Répondre à la requête en utilisant les mesures de similarité:

- le produit scalaire: $Score_{ps}(d_j, q) = \sum_{i=1}^n W_{ij} \times q_i$

$$Score_{ps}(d_1, q) = (1 \times 1) + (1 \times 0) + (1 \times 0) + (1 \times 0) = 1$$

$$Score_{ps}(d_2, q) = (1 \times 0) + (1 \times 0.44) + (1 \times 0) + (1 \times 0) = 0.44$$

$$Score_{ps}(d_3, q) = (1 \times 1) + (1 \times 0.44) + (1 \times 2) + (1 \times 0) = 3.44$$

$$Score_{ps}(d_4, q) = (1 \times 0) + (1 \times 0.44) + (1 \times 0) + (1 \times 2) = 2.44$$

* Classement:

$$D_3 > D_4 > D_1 > D_2$$

- la mesure de cosinus :

Score_{cos} (d_i, q)

$$\frac{\|d_i\| \times \|q\|}{\|d_i\| \times \|q\|}$$

$$\text{Score}_{\cos}(d_1, q) = \frac{1}{2 \times \sqrt{16}} \approx 0,125$$

$$\text{Score}_{\cos}(d_2, q) = \frac{0,44}{2 \times \sqrt{10,46}} \approx 0,064$$

$$\text{Score}_{\cos}(d_3, q) = \frac{3,44}{2 \times \sqrt{9,46}} \approx 0,56$$

$$\text{Score}_{\cos}(d_4, q) = \frac{2,44}{2 \times \sqrt{6,16}} \approx 0,48$$

* le classement : $D_3 > D_4 > D_1 > D_2$

Exemple de cours :

- Soit la requête suivante : $q = \text{"Système recherche information"}$

- Soit les documents : $D_1 = \text{"la conception de système d'information"}$
 $D_2 = \text{"recherche dans la recherche d'information"}$
 $D_3 = \text{"la connaissance sur système de gestion de BDD"}$

- Modèle de pondération locale utilisée : Modèle simple.

Déterminer le classement des documents par rapport à la requête donnée en utilisant un modèle vectoriel avec les mesures de similarité suivantes :

- produit scalaire
- cosinus.

$$N=3$$

- classement g-v: $D_3 > D_1 > D_5 > D_4 > D_2$

* Dirichlet: avec $M = 40$ $Score_{Dir}(D_i, \varphi) = \prod [\lambda_{Dir} P(q_i | D_i) + (1 - \lambda_{Dir}) P(q_i | \varphi_0)]$

$$\lambda_{Dir} = \frac{|D_i|}{M + |D_i|}$$

$$1 - \lambda_{Dir} = \frac{M}{|D_i| + M}$$

$$Score_{Dir}(D_1, \varphi) = \left[0.144 \times \frac{3}{7} + 0.59 \times \frac{8}{33} \right] \times \left[0.144 \times \frac{2}{7} + 0.59 \times \frac{4}{33} \right] \times \left[0.144 \times 0 + 0.59 \times \frac{3}{33} \right]$$

$$= 0.0032$$

$$Score_{Dir}(D_2, \varphi) = \left[0.33 \times 0 + 0.66 \times \frac{8}{33} \right] \times \left[0.33 \times 0 + 0.66 \times \frac{4}{33} \right] \times \left[0.33 \times 0 + 0.66 \times \frac{3}{33} \right]$$

$$= 0.00077$$

$$Score_{Dir}(D_3, \varphi) = 0.0038$$

$$Score_{Dir}(D_5, \varphi) = 0.0048$$

$$Score_{Dir}(D_4, \varphi) = 0.0024$$

- classement Dir: $D_3 > D_1 > D_5 > D_4 > D_2$

$\varphi_0 = \{t_1, t_2\}$ $\xrightarrow[\text{recherche}]{\text{Dir}}$ $D_3 > D_1 > D_5 > D_4 > D_2 \rightarrow$ Modèle de langue de R

$$P(t_i | R) = \sum_{i=1}^{|R|} P(t_i | D_i) \times \frac{Score(D_i)}{\sum_{D \in R} Score(D_i)}$$

$$Score(D_3) = 0.0038$$

$$Score(D_1) = 0.0032$$

$t_i = Recherche$

$$P(Recherche | R) = \frac{P(Recherche | D_3) \times Score(D_3)}{Score(D_3) + Score(D_1)} + \frac{P(Recherche | D_1) \times Score(D_1)}{Score(D_3) + Score(D_1)}$$

$$\frac{Score(D_1)}{Score(D_3) + Score(D_1)} = \frac{1}{7} \times \frac{38}{38+32} + \frac{2}{7} \times \frac{32}{38+32} = 0.23$$

TD 3

Exo 1:

on a: $C = \text{doc doc}$

$R = \{ \textcircled{D_2}, D_{15}, D_{16}, D_{18}, \textcircled{D_{29}}, D_{42}, D_{58}, \textcircled{D_{77}}, D_{83}, D_{85}, D_{94}, D_{99} \}$ list de doc
retourés pour la
requête q .

$P = \{ \textcircled{D_2}, D_{15}, \textcircled{D_{29}}, D_{47}, D_{67}, \textcircled{D_{77}}, D_{93}, D_{97} \}$ doc pertinents pour la q
dans la collection.

4) calculer: $RP = \{ D_2, D_{29}, D_{77} \} = 3$; $|R| = 42$; $|P| = 8$

rappel: $\frac{|RP|}{|P|} = \frac{3}{8} = \underline{0,375}$

précision = $\frac{|RP|}{|R|} = \frac{3}{42} = \underline{0,25}$

2) Calculer la F -mesure

$\alpha = 0,5$ $F = \frac{2PR}{P+R} = \frac{2(0,25)(0,375)}{0,25+0,375} = \underline{0,3}$

$\alpha = 0,8$ $F = \frac{4}{0,8 \cdot \frac{1}{0,25} + (1-0,8) \cdot \frac{1}{0,375}} = \underline{0,267}$

Exo 2:

1) calculons le rappel et la précision pour chaque requête.

$|P_{q_1}| = 8$ $|R_{q_1}| = 10$ $|RP_{q_1}| = 4$

précision(q_1) = $\frac{4}{10} = \underline{0,4}$; rappel(q_1) = $\frac{4}{8} = \underline{0,5}$

$|P_{q_2}| = 8$ $|R_{q_2}| = 10$ $|RP_{q_2}| = 5$

précision(q_2) = $\frac{5}{10} = \underline{0,5}$; rappel(q_2) = $\frac{5}{8} = \underline{0,625}$

2) Calculons:

$P@k = \frac{\text{nbr de doc pertinents au rang } k}{k}$

1x1

$Q_1 \quad P@5 = \frac{2}{5}$

$P@10 = \frac{4}{10}$

$Q_2 \quad P@5 = \frac{3}{5}$

$P@10 = \frac{5}{10}$

3) Calculer (MAP):

$AP_1 = \frac{1}{8} \left[1 + \frac{2}{5} + \frac{3}{6} + \frac{4}{10} \right] = 0,2875$

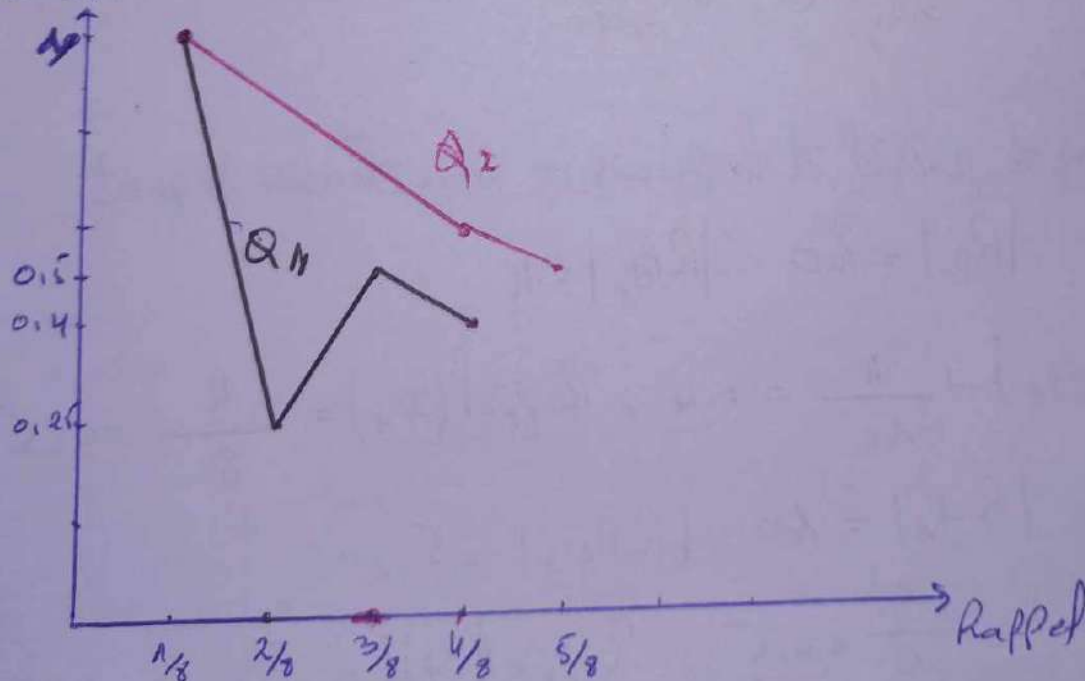
$AP_2 = \frac{1}{8} \left[3 + \frac{4}{7} + \frac{5}{10} \right] = 0,5$

$MAP = \frac{1}{2} [0,2875 + 0,5] = 0,39$

$R_{1j} = 1/8, 1, 1, 2/8, 3/8, 1, 1, 1, 4/8$

$R_{2j} = 1/8, 1/8, 1/8, 1, 1, 1, 4/8, 1, 1, 5/8$

Precision



Ex 3: (Modèle probabiliste).

	recherche	base	donnée	coordonnée	logique	math	persistance	stat
D ₁	1	0	0	1	0	0	1	0
D ₂	0	1	1	0	0	1	0	0
D ₃	1	0	0	1	1	1	1	1
D ₄	0	0	0	0	1	1	0	1
D ₅	1	1	1	1	0	0	1	0
n _i	3	2	2	3	2	3	3	2
N	5	5	5	5	5	5	5	5
S	3	3	3	3	3	3	3	3
s _i	2	1	1	2	1	2	2	1
P _i	2/3	1/3	1/3	2/3	1/3	2/3	2/3	1/3
U _i	1/2	1/2	1/2	1/2	1/2	1/2	1/2	1/2

$$P_i = \frac{S_i}{S}$$

$$U_i = \frac{n_i - s_i}{N - S}$$

$$\text{Score}(D, q) = \sum_{\substack{i \in q \\ \text{tied}}} \log_2 \left[\frac{P_i}{U_i} \times \frac{1 - U_i}{1 - P_i} \right]$$

$$\begin{aligned} \text{Score}(D_1, q) &= \log_2 \left[\frac{P_{\text{rech}}}{U_{\text{rech}}} \times \frac{1 - U_{\text{rech}}}{1 - P_{\text{rech}}} \right] + \log_2 \left[\frac{P_{\text{info}}}{U_{\text{info}}} \times \frac{1 - U_{\text{info}}}{1 - P_{\text{info}}} \right] \\ &+ \log_2 \left[\frac{P_{\text{pert}}}{U_{\text{pert}}} \times \frac{1 - U_{\text{pert}}}{1 - P_{\text{pert}}} \right] = \log_2 \left[\frac{2/3}{1/2} \times \frac{1 - 1/2}{1 - 2/3} \right] + \log_2 \left[\frac{2/3}{1/2} \times \frac{1 - 1/2}{1 - 2/3} \right] + \\ &\log_2 \left[\frac{2/3}{1/2} \times \frac{1 - 1/2}{1 - 2/3} \right] = 3 \log_2 \left(\frac{2/3}{1/2} \times \frac{1 - 1/2}{1 - 2/3} \right) = 3 \log_2 \left(\frac{4}{3} \times \frac{3}{2} \right) \\ &= 3 \log_2 \left(\frac{4}{2} \right) = \boxed{3} \end{aligned}$$

+ log₂

EX04: (Modèle de langue)

1) construction de Modèle de langue des (5) documents.

	recherche	info	pertinence	base	donnée	math	logique	flux	
MD ₁	2/7	3/7	2/7	0	0	0	0	0	7
MD ₂	0	0	0	2/5	2/5	4/5	0	0	5
MD ₃	1/7	2/7	4/7	0	0	1/7	1/7	1/7	7
MD ₄	0	0	0	0	0	2/6	2/6	2/6	6
MD ₅	2/8	3/8	1/8	1/8	1/8	0	0	0	8
MC	5/33	8/33	4/33	3/33	3/33	4/33	3/33	3/33	33

2) Calculons le score des doc pour la requête: "info - pertinence - flux" en utilisant les méthodes deissage suivantes:

* Jelinek Mercer: Pour $\lambda = 0.8$
$$\text{score}(D, q) = \prod_{q_i \in q} [\lambda \cdot P(q_i | D) + (1-\lambda) P(q_i | MC)]$$

$$\text{score}(D_1, q) = [0.8 \times \frac{3}{7} + 0.2 \times \frac{8}{33}] \times [0.8 \times \frac{2}{7} + 0.2 \times \frac{4}{33}] \times [0.8 \times 0 + 0.2 \times \frac{3}{33}]$$

$$= 0.0047$$

$$\text{score}(D_2, q) = [0.8 \times 0 + 0.2 \times \frac{8}{33}] \times [0.8 \times 0 + 0.2 \times \frac{4}{33}] \times [0.8 \times 0 + 0.2 \times \frac{3}{33}]$$

$$= 0.00024$$

$$\text{score}(D_3, q) = [0.8 \times \frac{2}{7} + 0.2 \times \frac{8}{33}] \times [0.8 \times \frac{4}{7} + 0.2 \times \frac{4}{33}] \times [0.8 \times \frac{1}{7} + 0.2 \times \frac{3}{33}]$$

$$= 0.0054$$

$$\text{score}(D_4, q) = [0.8 \times 0 + 0.2 \times \frac{8}{33}] \times [0.8 \times 0 + 0.2 \times \frac{4}{33}] \times [0.8 \times \frac{2}{6} + 0.2 \times \frac{3}{33}]$$

$$= 0.00033$$

$$\text{score}(D_5, q) = [0.8 \times \frac{3}{8} + 0.2 \times \frac{8}{33}] \times [0.8 \times \frac{1}{8} + 0.2 \times \frac{4}{33}] \times [0.8 \times 0 + 0.2 \times \frac{3}{33}]$$

$$= 0.00078$$

$$\text{score}(D_2, q) = \log_2 \left(\frac{1/3}{4/2} \times \frac{1 - 1/2}{1 - 1/3} \right) + \log_2 \left(\frac{1/3}{1/2} \times \frac{1 - 1/2}{1 - 1/3} \right) + \log_2 \left(\frac{2/3}{1/2} \times \frac{1 - 1/2}{1 - 1/3} \right) = 2 \log_2 \left(\frac{2}{3} \times \frac{3}{4} \right) + \log_2 \left(\frac{4}{3} \times \frac{3}{2} \right) = 2 \log_2 \left(\frac{2}{4} \right) + \log_2 \left(\frac{4}{2} \right) = -1$$

$$\text{score}(D_3, q) = 4 \log_2 \left(\frac{4}{2} \right) + 2 \log_2 \left(\frac{2}{4} \right) = 2$$

$$\text{score}(D_4, q) = \log_2 \left(\frac{4}{2} \right) + 2 \log_2 \left(\frac{2}{4} \right) = -1$$

$$\text{score}(D_5, q) = 3 \log_2 \left(\frac{4}{2} \right) + 2 \log_2 \left(\frac{2}{4} \right) = 1$$

suite exo 4 :

$$P(\text{info} | R) = 0.33 \checkmark$$

$$P(\text{pert} | R) = 0.23 \checkmark$$

$$P(\text{Basse} | R) = 0.$$

$$P(\text{donné} | R) = 0$$

$$P(\text{long} | R) = 0.07$$

$$P(\text{matti} | R) = 0.07$$

$$P(\text{flou} | R) = 0.07 \checkmark$$

EX04

4) cons

rec

MD₁

MD₂

MD₃

MD₄

MD₅

MC

e/catent
en uti

*Teb

score(D

score(D

score(D

score(D

score

	Conception	de	système	d	info	recherche	obus	la	Connaiss	sur	gestion	BDD
n_i	1	2	2	2	2	1	1	3	1	1	1	1
$colg_i$	1.58	0	0.58	0.58	0.58	0.58	1.58	0	1.58	1.58	1.58	1.58
tf_{i1}	1	1	1	1	1	0	0	1	0	0	0	0
tf_{i2}	0	0	0	1	1	2	1	1	0	0	0	0
tf_{i3}	0	2	1	0	0	0	0	1	0	0	0	0
W_{i1}	1.58	0	0.58	0.58	0.58	0	0	1	1	1	1	1
W_{i2}	0	0	0	0.58	0.58	0	0	0	0	0	0	0
W_{i3}	0	0	0.58	0	0	3.146	0.58	0	0	0	0	0
W_{iq}	0	0	<u>1</u>	0	<u>1</u>	<u>1</u>	0	0	1.58	1.58	1.58	1.58
							0	0	0	0	0	0

$$\text{Score}_{P-S}(d_1, q) = (1 \times 0.58) + (1 \times 0.58) + (1 \times 0) = \boxed{1.16}$$

$$\text{Score}_{P-S}(d_2, q) = (1 \times 0) + (1 \times 0.58) + (1 \times 1.16) = \boxed{1.74}$$

$$\text{Score}_{P-S}(d_3, q) = (1 \times 0.58) + (1 \times 0) + (1 \times 0) = \boxed{0.58}$$

$$D_2 > D_1 > D_3$$

$$\text{Score}_{\cos}(d_1, q) = \frac{1.16}{\sqrt{3} \times \sqrt{3.5}}$$