

Esercizio 1 Relevance Feedback –

Considerate la query: “*cheap CDs cheap DVDs extremely cheap CDs*”

L'utente esamina i risultati e giudica rilevante

d1: *CDs cheap software cheap CDs*

e non rilevante

d2: *cheap thrills DVDs*

Assumendo di usare come pesi solo i valori di frequenza, e di assumere $\alpha = 1$, $\beta = 0.75$, $\gamma = 0.25$ espandete la query usando la formula di Rocchio

Esercizio 2 LSI

Considerate la matrice L (termini documenti)

0.0 2.0 1.0

0.0 3.0 0.0

2.0 1.0 0.0

E la sua decomposizione SVD

Matrice U

-0.545 -0.238 -0.804

-0.762 -0.259 0.593

-0.349 0.936 -0.040

Matrice V^T

-0.182 -0.973 -0.142

0.978 -0.165 -0.124

-0.097 0.162 -0.982

Matrice Sigma

3.830 0.000 0.000

0.000 1.914 0.000

0.000 0.000 0.819

Calcolate la matrice L' che si ottiene con una riduzione di rango 2 della matrice sigma, e confrontatela con la matrice L originaria.

Esercizio 3

q	1	1	1	1	0	0	0	0	0
D1	0	1	1	0	0	0	1	1	1
D2	1	1	0	0	0	0	0	1	1
D3	1	0	1	1	1	1	1	1	1
D4	1	0	1	0	0	0	1	0	0
D5	1	0	1	0	0	0	1	0	0
D6	1	0	0	1	0	0	0	1	0
D7	1	0	0	0	0	1	0	0	0
D8	0	0	0	1	0	1	0	1	0
D9	0	1	0	1	1	1	0	1	1

Ordinate i documenti rispetto alla query utilizzando il probabilistic retrieval con il Binary Independence Model, e metodo incrementale per il calcolo del Retrieval Status Value. Ponete $|V|=4$.

Riassunto delle formule:

$$RSD(D_j) = \sum_{x_{ji}=1, q_i=1} c_i$$

$$c_i = \log \frac{p_i(1-r_i)}{r_i(1-p_i)}$$

Usate le stime di p_i e r_i con “adjustement” per basse frequenze, semplificato come segue:

$$p_i = \frac{V_i + 0,5}{V + 1}, \quad r_i = \frac{n_i - V_i + 0,5}{N - V + 1}$$

Dove V_i è il numero dei documenti in V che contiene la parola i , n_i è il numero dei documenti in N che contiene la parola i , V è usato al posto di $|V|$ ed N è la dimensione della collezione completa (9 nel caso in esame).

Soluzione Esercizio 1

word	query q	d_1	d_2
CDs	2	2	0
cheap	3	2	1
DVDs	1	0	1
extremely	1	0	0
software	0	1	0
thrills	0	0	1

$$1.0 q + 0.75 d_1 - 0.25 d_2$$

si ottiene la query di Rocchio:

$$Q: (3.5 \ 4.25 \ 0.75 \ 1 \ 0.75 \ 0)$$

Soluzione <http://www.bluebit.gr/matrix-calculator/>

$$U_2 \Sigma_2$$

$$\begin{pmatrix} -2.087 & -0.456 \\ -2.918 & -0.496 \\ -1.337 & 1.792 \end{pmatrix}$$

$$\text{Matrice } L' \ U_2 \Sigma_2 V_2^T$$

$$\begin{pmatrix} -0.066 & 2.106 & 0.353 \\ 0.046 & 2.921 & 0.476 \\ 1.996 & 1.005 & -0.032 \end{pmatrix}$$

la nuova matrice delle co-occorrenze fra termini LL^T è

$$\begin{pmatrix} t1 & 4.564 & 6.317 & 1.973 \\ t2 & 6.317 & 8.761 & 3.012 \\ t3 & 1.973 & 3.012 & 4.995 \end{pmatrix}$$

che rispetto alla matrice delle co-occorrenze originaria

$$\begin{pmatrix} t1 & 5.000 & 6.000 & 0.000 \\ t2 & 6.000 & 9.000 & 0.000 \\ t3 & 2.000 & 3.000 & 4.000 \end{pmatrix}$$

assegna una probabilità di co-occorrenza anche a t1-t3 e t2-t3

Soluzione

Come spiegato a lezione (si vedano le formule), al primo passo si assume $p_i=0,5$, e, per $N \gg S$, $r_i \approx n_i/N$

Con queste semplificazioni, $c_i \approx \text{idf}(w_i)$

Dunque, al primo passo:

- 1) Stima approssimata delle r_i : $r_i = n_i/N$, $N=9$ (si noti che in realtà in questo esempio N è molto piccolo, la semplificazione non dovrebbe essere valida..)

Si ha: $r_1=6/9$ $r_2=3/9$ $r_3=3/9$ $r_4=4/9$ $r_5=2/9$ $r_6=4/9$ $r_7=4/9$ $r_8=6/9$ $r_9=4/9$

Al passo 1, $p_i=0,5$ per ogni i . e dunque:

$$\text{RTV}(1) = \log 9/3 + \log 9/3 = 2 \times 0,477 = 0,954$$

$$\text{RTV}(2) = \log 9/6 + \log 9/3 = 0,107 + 0,477 = 0,584$$

$$\text{RTV}(3) = \log 0,107 + 0,477 + \log(9/4) = 0,584 + 0,352 = 0,936$$

$$\text{RTV}(4) = 0,584$$

$$\text{RTV}(5) = 0,584$$

$$\text{RTV}(6) = 0,107 + 0,352 = 0,459$$

$$\text{RTV}(7) = 0,107$$

$$\text{RTV}(8) = 0,352$$

$$\text{RTV}(9) = 0,477 + 0,352 = 0,829$$

Poichè abbiamo stabilito $|V|=4$, vengono presentati: D1, D3, D9, D2

2) Seconda iterazione. Si ri-stimano le p_i e r_i a partire dai 4 documenti selezionati, usando per p_i e r_i le formule nel testo dell'esercizio.

Le nuove stime dei p_i al passo 2 sono:

$$p_1 = \frac{2 + \frac{1}{2}}{4 + 1} = 0,5 \quad p_2 = \frac{3 + 0,5}{5} = 0,7 \quad p_3 = p_1 \quad p_4 = p_1$$

e per r_i si ha:

$$r_1 = \frac{6 - 2 + \frac{1}{2}}{9 - 4 + 1} = \frac{4 + 0,5}{6} = 0,75, \quad r_2 = \frac{3 - 3 + 0,5}{6} = 0,0883, \quad r_3 = \frac{4 - 2 + 0,5}{6} = 0,416,$$

$$r_4 = \frac{4 - 2 + 0,5}{6} = 0,416$$

$$\text{Dunque avrò: } c_1 = \log \frac{p_1(1-r_1)}{r_1(1-p_1)} = \log \frac{0,5 \cdot 0,25}{0,75 \cdot 0,5} = -0,481, \quad c_2 = \log \frac{0,7 \cdot 0,91}{0,09 \cdot 0,3} = 1,37$$

$$c_3 = \log \frac{0,5 \cdot 0,58}{0,42 \cdot 0,5} = \log 1,38 = 0,14, \quad c_4 = c_3$$

Notate come la presenza del primo termine risulta penalizzante: infatti è poco discriminante (elevato IDF) e si distribuisce equamente nei documenti rilevanti e non rilevanti (nello step 1).

Al contrario, il termine più discriminante è c_2 , che ha un'alta probabilità di comparire fra i rilevanti (in 3 documenti di V su 4).

Ricalcoliamo dunque RTV per i vari documenti:

$$\text{RTV}(1)=c_2+c_3=1,37+0,14=1,51$$

$$\text{RTV}(2)=c_1+c_2=-0,481+1,37=0,889$$

$$\text{RTV}(3)=c_1+c_3+c_4=-0,481+0,14+0,14=0,201$$

...

$$\text{RTV}(8)=c_4=0,14$$

$$\text{RTV}(9)=c_2+c_4=1,37+0,14=1,77$$

Verrebbero presentati: D_9, D_1, D_2, D_3

Come si vede, i documenti restano gli stessi, ma con diverso ranking. Poichè i documenti sono gli stessi, ovviamente ci si arresta, la terza iterazione genererebbe gli stessi valori di RTV.