

INFORMATION RETRIEVAL – SHORT EXERCISES II – VECTOR SPACE MODEL AND LATENT SEMANTIC INDEXING

I. Consider a set of terms $T = \{t_1, t_2, t_3, t_4\}$ and the following collection of two documents: $D1 = \{t_1, t_2, t_1, t_2, t_3\}$ and $D2 = \{t_4, t_2, t_2, t_3\}$. Consider query $Q = \{t_1, t_4\}$. Represent $D1$, $D2$, and Q using TF (normalized Bag-Of-Words).

TF	t_1	t_2	t_3	t_4
D1	2/2	2/2	1/2	0
D2	0	2/2	1/2	1/2
Q	1/1	0	0	1/1

$max = 2$

Compute IDFs for all four terms (note that only $D1$ and $D2$ are included in the collection).

	t_1	t_2	t_3	t_4
IDF	$\log 2$	$\log 1 = 0$	$\log 1 = 0$	$\log 2$

$$idf_i = \log\left(\frac{N}{df_i}\right)$$

II. Consider the below term-document matrix C for the bag-of-words representation of five documents $D1-D5$ in the space of six terms t_1-t_6 . Using the SVD factorization method, matrix C has been decomposed into matrices K , S , and D^T given below. The rank of C is 4 ($4 \leq \min\{6,5\}$), so 4 concepts (semantic dimensions) were discovered.

						concept				document					
							c_1	c_2	c_3	c_4	D_1	D_2	D_3	D_4	D_5
$C =$	t_1	5	5	0	0	1	-0.27	0.55	-0.78	0	13.74	0	0	0	0
	t_2	4	5	1	1	0	-0.29	0.47	0.44	-0.71	0	10.88	0	0	0
	t_3	5	4	1	1	0	-0.29	0.47	0.44	0.71	0	0	1.36	0	0
	t_4	0	0	4	4	4	-0.45	-0.29	-0.01	0	0	0	0	1	0
	t_5	0	0	5	5	5	-0.56	-0.36	-0.02	0	0	0	0	0	1
	t_6	1	1	4	4	4	-0.50	-0.18	-0.05	0	0	0	0	0	0