

# Foundations of Information Retrieval

## Exercise 3

Issue date: 01.11.2012  
 Exercise session on: 08.11.2012, 15.11.2012  
 Questions to: Elena Demidova: [demidova@l3s.de](mailto:demidova@l3s.de)

### TF/IDF weighting and Vector Space Model

Given is a document collection with two documents:

**Document 1:**

*A good cook could cook as much cookies as a good cook who could cook cookies.*

**Document 2:**

*Best chocolate chip cookies recipe.*

- a. Index the both documents using an inverted index.

Apply case folding, use stop words list: {a, as}.

*good*(DF=1)            *D1*(2)  
*cook* (DF=1)            *D1*(4)  
*could*(DF=1)           *D1*(2)  
*much*(DF=1)            *D1*(1)  
*cookies*(DF=2)        *D1*(2), *D2*(1)  
*who* (DF=1)            *D1*(1)  
*best*(DF=1)            *D2*(1)  
*chocolate*(DF=1)    *D2*(1)  
*chip* (DF=1)            *D2*(1)  
*recipe* (DF=1)        *D2*(1)

Compute the cosine similarity between these documents and the following queries.  
 Apply the **l<sub>tc</sub>** weighting.

Q<sub>1</sub>= *cookies*

Q<sub>2</sub>= *how to cook chocolate cake*

	DF	idf= log(N/DF)	D <sub>1</sub> : tf-raw	D <sub>1</sub> : tf-wght (see slide 21)	D <sub>1</sub> : tf-idf weight	D <sub>1</sub> : n <sup>2</sup> lized (lengthD <sub>1</sub> =0.84)
<i>good</i>	1	log(2/1)=0.3	2	1+log(2)=1.3	1.3*0.3=0.39	0.39/0.84= <b>0.46</b>
<i>cook</i>	1	0.3	4	1.6	0.48	<b>0.56</b>
<i>could</i>	1	0.3	2	1.3	0.39	<b>0.46</b>
<i>much</i>	1	0.3	1	1	0.3	<b>0.35</b>
<i>cookies</i>	2	0	2	1.3	0	<b>0</b>
<i>who</i>	1	0.3	1	1	0.3	<b>0.35</b>
<i>best</i>	1	0.3	0	0	0	<b>0</b>
<i>chocolate</i>	1	0.3	0	0	0	<b>0</b>
<i>chip</i>	1	0.3	0	0	0	<b>0</b>
<i>recipe</i>	1	0.3	0	0	0	<b>0</b>

$$\text{length } D_1 = \text{sqrt}(0.39^2+0.48^2+0.39^2+0.3^2+0.3^2) = 0.84$$

	DF	idf= log(N/DF)	D <sub>2</sub> : tf-raw	D <sub>2</sub> : tf-wght (see slide 21)	D <sub>2</sub> : tf-idf weight	D <sub>2</sub> : n'lized (lengthD <sub>2</sub> =0.6)
<i>good</i>	1	log(2/1) = 0.3	0	0	0	0
<i>cook</i>	1	0.3	0	0	0	0
<i>could</i>	1	0.3	0	0	0	0
<i>much</i>	1	0.3	0	0	0	0
<i>cookies</i>	2	log(2/2) = 0	1	1	0	0
<i>who</i>	1	0.3	0	0	0	0
<i>best</i>	1	0.3	1	1	0.3	0.3/0.6=0.5
<i>chocolate</i>	1	0.3	1	1	0.3	0.5
<i>chip</i>	1	0.3	1	1	0.3	0.5
<i>recipe</i>	1	0.3	1	1	0.3	0.5

$$\text{length } D_2 = \sqrt{4 * 0.3^2} = 0.6$$

	Q1: tf-raw, tf-wght	Q1: tf-idf, n.length	Q2: tf-raw, tf-wght	Q2: tf-idf weight	Q2: n'lized (length=0.42)	D1: n'lized (see above)	D2: n'lized (see above)
<i>good</i>	0	0	0	0	0	0.46	0
<i>cook</i>	0	0	1	0.3*1=0.3	0.3/0.42=0.71	0.56	0
<i>could</i>	0	0	0	0	0	0.46	0
<i>much</i>	0	0	0	0	0	0.35	0
<i>cookies</i>	1	0	0	0	0	0	0
<i>who</i>	0	0	0	0	0	0.35	0
<i>best</i>	0	0	0	0	0	0	0.5
<i>chocolate</i>	0	0	1	0.3	0.71	0	0.5
<i>chip</i>	0	0	0	0	0	0	0.5
<i>recipe</i>	0	0	0	0	0	0	0.5

$$\text{length } Q_2 = \sqrt{2 * 0.3^2} = 0.42$$

$$\cos(d1, q1) = 0$$

$$\cos(d2, q1) = 0$$

→ The both documents are not relevant for Q1 under ltc weighting (as IDF of Q1 is 0).

$$\cos(d1, q2) = 0.71 * 0.56 = 0.3976$$

$$\cos(d2, q2) = 0.71 * 0.5 = 0.355$$

→ D1 is more relevant for Q2 under ltc weighting.

Explain the similarity scores! How would the result change with the lnc weighting?

As Q<sub>1</sub> contains only one term whose IDF is 0, there will be no relevant results under the ltc weighting. Lnc does not take document frequency into account, such that the both documents would have non-zero scores for Q<sub>1</sub>. D<sub>1</sub> would be ranked ahead of D<sub>2</sub>, as it has a higher TF. No change for Q<sub>2</sub>.