# Seminar 4

**Definition 1 (Term weight)** *Weight of a term t in a document d is counted as*

$$w_{t,d} = \begin{cases} 1 + \log\left(tf_{t,d}\right) & \text{if } n > 0 \\ 0 & \text{otherwise} \end{cases}$$

*where $tf_{t,d}$ is the number of terms t in a document d.*

**Definition 2 (Inverse document frequency)** *Inverse document frequency of a term t is defined as*

$$idf_t = \log\left(\frac{N}{df_t}\right)$$

*where N is the number of all documents and $df_t$ (document frequency) is the number of documents that contain t.*

**Definition 3 (tf-idf weighting scheme)** *In the tf-idf weighting scheme, a term t in a document d has weight*

$$tf\text{-}idf_{t,d} = tf_{t,d} \cdot idf_t$$

**Definition 4 (Cosine (Euclidean) normalization)** *A vector v is cosine-normalized by*

$$v_j = \frac{v_j}{||v||} = \frac{v_j}{\sqrt{\sum_{k=1}^{|v|} v_k{}^2}}$$

*where $v_j$ be the number on the j-th position in v.*

## Exercise 1

Consider the frequency table of the words of three documents $doc_1$, $doc_2$, $doc_3$ below. Calculate the *tf-idf* weight of the terms *car, auto, insurance, best* for each document. *idf* values of terms are in the table.

|           | $doc_1$ | $doc_2$ | $doc_3$ | $idf$ |
|-----------|---------|---------|---------|-------|
| car       | 27      | 4       | 24      | 1.65  |
| auto      | 3       | 33      | 0       | 2.08  |
| insurance | 0       | 33      | 29      | 1.62  |
| best      | 14      | 0       | 17      | 1.5   |

Table 1: Exercise.

---

After counting *tf-idf* weights by Definition 3 individually for each term we get the following table

|  | tf-idf | | |
|---|---|---|---|
|  | $doc_1$ | $doc_2$ | $doc_3$ |
| car | 44.55 | 6.6 | 39.6 |
| auto | 6.24 | 68.64 | 0 |
| insurance | 0 | 53.46 | 46.98 |
| best | 21 | 0 | 25.5 |

<div align="center">Table 2: Solution.</div>

## Exercise 2

Count document representations as normalized Euclidean weight vectors for each document from the previous exercise. Each vector has four components, one for each term.

Normalized Euclidean weight vectors are counted by Definition 4. Denominators $m_{doc_n}$ for the individual documents are

$$m_{doc_1} = \sqrt{44.55^2 + 6.24^2 + 21^2} = 49.6451$$

$$m_{doc_2} = \sqrt{6.6^2 + 68.64^2 + 53.46^2} = 87.2524$$

$$m_{doc_3} = \sqrt{39.6^2 + 46.98^2 + 25.5^2} = 66.5247$$

and the document representations are

$$d_1 = \left( \frac{44.55}{49.6451}; \frac{6.24}{49.6451}; \frac{0}{49.6451}; \frac{21}{49.6451} \right) = (0.8974; 0.1257; 0; 0.423)$$

$$d_2 = \left( \frac{6.6}{87.2524}; \frac{68.64}{87.2524}; \frac{53.46}{87.2524}; \frac{0}{87.2524} \right) = (0.0756; 0.7876; 0.6127; 0)$$

$$d_3 = \left( \frac{39.6}{66.5247}; \frac{0}{66.5247}; \frac{46.98}{66.5247}; \frac{25.5}{66.5247} \right) = (0.5953; 0; 0.7062; 0.3833)$$

## Exercise 3

Based on the weights from the last exercise, compute the relevance scores of the three documents for the query *car insurance*. Use each of the two weighting schemes:

a) Term weight is 1 if the query contains the word and 0 otherwise.

b) Euclidean normalized *tf-idf*.

Please note that a document and a representation of this document are different things. Document is always fixed but the representations may vary under different settings and conditions. In this exercise we fix document representations from the last exercises and will count relevance scores for query and documents under two different representations of the query. It might be helpful to view on a query as on another document, as it is a sequence of words.

We count the relevance scores for **a)** as the scalar products of the representation of the query $q = (1, 0, 1, 0)$ with representations of the documents $d_n$ from the last exercise:

$$q \cdot d_1 = 1 \cdot 0.8974 + 0 \cdot 0.1257 + 1 \cdot 0 + 0 \cdot 0.423 = 0.8974$$

$$q \cdot d_2 = 1 \cdot 0.0756 + 0 \cdot 0.7876 + 1 \cdot 0.6127 + 0 \cdot 0 = 0.6883$$

$$q \cdot d_3 = 1 \cdot 0.5953 + 0 \cdot 0 + 1 \cdot 0.7062 + 0 \cdot 0.3833 = 1.3015$$

For **b)** we first need the normalized *tf-idf* vector $q$, which is obtained by dividing each component of the query by the length of *idf* vector $\sqrt{1.65^2 + 0^2 + 1.62^2 + 0^2} = 2.3123$

| | tf | idf | tf-idf | q |
|---:|---|---|---|---|
| car | 1 | 1.65 | 1.65 | 0.7136 |
| auto | 0 | 2.08 | 0 | 0 |
| insurance | 1 | 1.62 | 1.62 | 0.7006 |
| best | 0 | 1.5 | 0 | 0 |

Table 3: Process of finding the Euclidean normalized *tf-idf*.

Now we multiply $q$ with the document vectors and we obtain the relevance scores:

$$q \cdot d_1 = 0.7136 \cdot 0.8974 + 0 \cdot 0.1257 + 0.7006 \cdot 0 + 0 \cdot 0.423 = 0.6404$$

$$q \cdot d_2 = 0.7136 \cdot 0.0756 + 0 \cdot 0.7876 + 0.7006 \cdot 0.6127 + 0 \cdot 0 = 0.4832$$

$$q \cdot d_3 = 0.7136 \cdot 0.5953 + 0 \cdot 0 + 0.7006 \cdot 0.7062 + 0 \cdot 0.3833 = 0.9196$$

## Exercise 4

Calculate the vector-space similarity between the query *digital cameras* and a document containing *digital cameras and video cameras* by filling in the blank columns in the table below. Assume $N = 10000000$, logarithmic term weighting (columns $w$) for both query and documents, *idf* weighting only for the query and cosine normalization only for the document. *and* is a STOP word.

| | df | Query | | | | Document | | | relevance |
|---|---|---|---|---|---|---|---|---|---|
| | | tf | w | idf | q | tf | w | d | $q \cdot d$ |
| digital | 10 000 | | | | | | | | |
| video | 100 000 | | | | | | | | |
| cameras | 50 000 | | | | | | | | |

Table 4: Exercise.

---

The *tf* value is filled according to the occurrences of the terms in both query and document.

$$\begin{aligned} \text{tf}_q &= \textit{digital cameras} & &= (1, 0, 1) \\ \text{tf}_d &= \textit{digital cameras and video cameras} & &= (1, 1, 2) \end{aligned}$$

Logarithmic weighting uses the Definition 1. For the query the values are

$$
\begin{aligned}
w_{digital} &= 1 + \log(1) & = 1 + 0 & = 1 \\
w_{video} &= 0 \\
w_{cameras} &= 1 + \log(1) & = 1 + 0 & = 1
\end{aligned}
$$

and for the document

$$
\begin{aligned}
w_{digital} &= 1 + \log(1) & = 1 + 0 & = 1 \\
w_{video} &= 1 + \log(1) & = 1 + 0 & = 1 \\
w_{cameras} &= 1 + \log(2) & = 1 + 0.301 & = 1.301
\end{aligned}
$$

Now we need to count the *idf* weights for the query. These are counted by Definition 2.

$$
\begin{aligned}
idf_{digital} &= \log\left(\frac{10^7}{10^4}\right) & = \log(10^3) & = 3 \\
idf_{video} &= \log\left(\frac{10^7}{10^5}\right) & = \log(10^2) & = 2 \\
idf_{cameras} &= \log\left(\frac{10^7}{5 \times 10^4}\right) & = \log(200) & = 2.301
\end{aligned}
$$

and $q = w \cdot idf$. Cosine normalization for the document is counted similarly as in the last exercises by Definition 4 using $w$.

$$
d_{digital} = \frac{1}{\sqrt{1^2 + 1^2 + 1.301^2}} = 0.5204
$$

$$
d_{video} = \frac{1}{\sqrt{1^2 + 1^2 + 1.301^2}} = 0.5204
$$

$$
d_{cameras} = \frac{1.301}{\sqrt{1^2 + 1^2 + 1.301^2}} = 0.677
$$

The score is the scalar multiple of $q$ and $d$. The final table is

| | | Query | | | | Document | | | relevance |
|---|---|---|---|---|---|---|---|---|---|
| | $df$ | $tf$ | $w$ | $idf$ | $q$ | $tf$ | $w$ | $d$ | $q \cdot d$ |
| digital | 10 000 | 1 | 1 | 3 | 3 | 1 | 1 | 0.5204 | 1.5612 |
| video | 100 000 | 0 | 0 | 2 | 0 | 1 | 1 | 0.5204 | 0 |
| cameras | 50 000 | 1 | 1 | 2.301 | 2.301 | 2 | 1.301 | 0.677 | 1.5578 |

Table 5: Solution.

and the similarity score is

$$
score(d, q) = \sum_{i=1}^{3} (d_i \cdot q_i) = 3.119.
$$

## Exercise 5

Show that for the query $q_1 =$ *affection* the documents in the table below are sorted by relevance in the opposite order as for the query $q_2 =$ *jealous gossip*. Query is *tf* weight normalized.

|           | SaS   | PaP   | WH    |
|-----------|-------|-------|-------|
| affection | 0.996 | 0.993 | 0.847 |
| jealous   | 0.087 | 0.120 | 0.466 |
| gossip    | 0.017 | 0     | 0.254 |

Table 6: Exercise.

We add queries to the original table:

|           | SaS   | PaP   | WH    | $q_1$ | $q_2$ |
|-----------|-------|-------|-------|-------|-------|
| affection | 0.996 | 0.993 | 0.847 | 1     | 0     |
| jealous   | 0.087 | 0.120 | 0.466 | 0     | 1     |
| gossip    | 0.017 | 0     | 0.254 | 0     | 1     |

Table 7: Exercise with queries.

Now we normalize the vectors $q_i$ by Definition 4 and get

|           | SaS   | PaP   | WH    | $q_1$ | $q_2$ | $q_{1n}$ | $q_{2n}$ |
|-----------|-------|-------|-------|-------|-------|----------|----------|
| affection | 0.996 | 0.993 | 0.847 | 1     | 0     | 1        | 0        |
| jealous   | 0.087 | 0.120 | 0.466 | 0     | 1     | 0        | 0.7071   |
| gossip    | 0.017 | 0     | 0.254 | 0     | 1     | 0        | 0.7071   |

Table 8: Exercise with queries after normalization.

In the last step we count the similarity score between the queries and documents by $score(d,q) = \sum_{i=1}^{|d|}(d_i \cdot q_i)$

$$
\begin{aligned}
score(SaS, q_1) &= 0.9961 \cdot 1 + 0.087 \cdot 0 + 0.017 \cdot 0 & = 0.9961 \\
score(PaP, q_1) &= 0.993 \cdot 1 + 0.120 \cdot 0 + 0 \cdot 0 & = 0.993 \\
score(WH, q_1) &= 0.847 \cdot 1 + 0.466 \cdot 0 + 0.254 \cdot 0 & = 0.847 \\
\\
score(SaS, q_2) &= 0.9961 \cdot 0 + 0.087 \cdot 0.7071 + 0.017 \cdot 0.7071 & = 0.0735 \\
score(PaP, q_2) &= 0.993 \cdot 0 + 0.120 \cdot 0.7071 + 0 \cdot 0.7071 & = 0.0849 \\
score(WH, q_2) &= 0.847 \cdot 0 + 0.466 \cdot 0.7071 + 0.254 \cdot 0.7071 & = 0.5091
\end{aligned}
$$

The ordering for $q_1$ is SaS > PaP > WH and for $q_2$ is WH > PaP > SaS, and we see that they are opposite.