

Latent Semantic Analysis

1. Latent Semantic Analysis

a. Inputs are documents. Each word is a feature. We can represent each document by

i. The presence of each word (0 / 1)

| | data | information | retrieval | brain | lung |
|-------------------|------|-------------|-----------|-------|------|
| CS-paper-1 | 1 | 1 | 1 | 0 | 0 |

term-to-concept similarity

| | | | | |
|---|---|---|---|---|
| 1 | 1 | 1 | 0 | 0 |
|---|---|---|---|---|

 \times

| |
|-----|
| .58 |
| .58 |
| .58 |
| 0 |
| 0 |

 $= 1.74$

doc-to-concept similarity / CS feature

ii. Count of the word (0, 1, ...)

| | data | information | retrieval | brain | lung |
|-------------------|------|-------------|-----------|-------|------|
| CS-paper-1 | 2 | 2 | 2 | 0 | 0 |

term-to-concept similarity

| | | | | |
|---|---|---|---|---|
| 2 | 2 | 2 | 0 | 0 |
|---|---|---|---|---|

 \times

| |
|-----|
| .58 |
| .58 |
| .58 |
| 0 |
| 0 |

 $= 3.48$

doc-to-concept similarity

2. Latent Semantic Analysis

| | data | information | retrieval | brain | lung |
|-------------|------|-------------|-----------|-------|------|
| CS-paper-1 | 1 | 1 | 1 | 0 | 0 |
| CS-paper-2 | 2 | 2 | 2 | 0 | 0 |
| CS-paper-3 | 1 | 1 | 1 | 0 | 0 |
| CS-paper-4 | 5 | 5 | 5 | 0 | 0 |
| Med-paper-1 | 0 | 0 | 0 | 2 | 2 |
| Med-paper-2 | 0 | 0 | 0 | 3 | 3 |
| Med-paper-3 | 0 | 0 | 0 | 1 | 1 |

a.

$$\begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 2 & 2 & 2 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 \\ 5 & 5 & 5 & 0 & 0 \\ 0 & 0 & 0 & 2 & 2 \\ 0 & 0 & 0 & 3 & 3 \\ 0 & 0 & 0 & 1 & 1 \end{bmatrix} = \begin{bmatrix} 0.18 & 0 \\ 0.36 & 0 \\ 0.18 & 0 \\ 0.90 & 0 \\ 0 & 0.53 \\ 0 & 0.80 \\ 0 & 0.27 \end{bmatrix} \times \begin{bmatrix} 9.64 & 0 \\ 0 & 5.29 \end{bmatrix} \times \begin{bmatrix} 0.58 & 0.58 & 0.58 & 0 & 0 \\ 0 & 0 & 0 & 0.71 & 0.71 \end{bmatrix}$$

b.

doc-to-concept
similarity matrix

$$\begin{bmatrix} 0.18 & 0 \\ 0.36 & 0 \\ 0.18 & 0 \\ 0.90 & 0 \\ 0 & 0.53 \\ 0 & 0.80 \\ 0 & 0.27 \end{bmatrix} \times \begin{bmatrix} 9.64 & 0 \\ 0 & 5.29 \end{bmatrix} \times \begin{bmatrix} 0.58 & 0.58 & 0.58 & 0 & 0 \\ 0 & 0 & 0 & 0.71 & 0.71 \end{bmatrix}$$

c.

d. Each document can be better represented by

- i. Frequency of the word ($n_i / \sum n_i$)
- ii. TfIdf
 1. tf: term frequency in the document
 2. idf: $\log(\text{number of documents} / \text{number of documents that contain the term})$