# Latent Semantic Analysis

Boston University CS 506 - Lance Galletti

doc-to-term similarity **X** term-to-concept similarity **=** doc-to-concept similarity

doc-to-term similarity      X      term-to-concept similarity      =      doc-to-concept similarity

# Latent Semantic Analysis

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

- The presence of each word (0 / 1)

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

doc-to-term similarity $\quad$ X $\quad$ term-to-concept similarity $\quad$ = $\quad$ doc-to-concept similarity
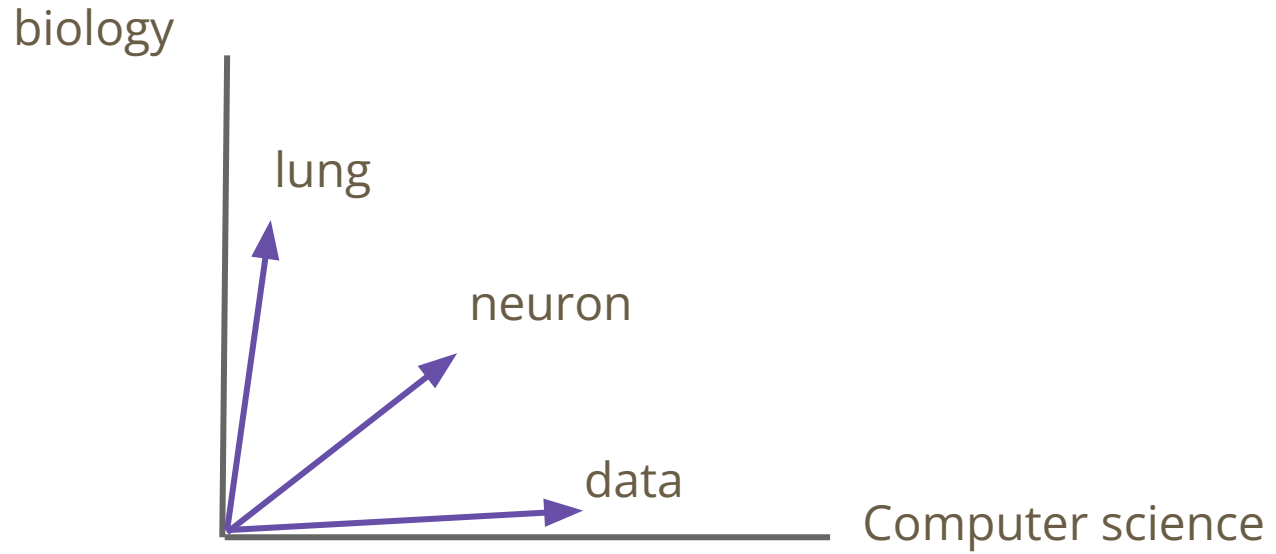
doc-to-term similarity    X    term-to-concept similarity    =    doc-to-concept similarity

# In theory

biology

lung

neuron

data

Computer science

# In practice

# In practice



Words with similar semantic meanings should be close

# In actual practice



king

queen

car  bottle

?  ?

Words with similar semantic meanings should be close

Lots of ways to generate embeddings. SVD is one of them

term-to-concept similarity

| 1 | 1 | 1 | 0 | 0 |

X

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

= 1.74

embedding

doc-to-concept similarity
/ CS feature

# Latent Semantic Analysis

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

- The presence of each word (0 / 1)
- 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 |

X

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

=     3.48

doc-to-concept similarity

# 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 |

# Latent Semantic Analysis

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

# Latent Semantic Analysis

| | | | | |
|---|---|---|---|---|
| 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 |

← Doc to term similarity

# Latent Semantic Analysis

CS concept        MD concept

| | |
|---|---|
| 0.18 | 0 |
| 0.36 | 0 |
| 0.18 | 0 |
| 0.90 | 0 |
| 0 | 0.53 |
| 0 | 0.80 |
| 0 | 0.27 |

**X**

| | |
|---|---|
| 9.64 | 0 |
| 0 | 5.29 |

**X**

| | | | | |
|---|---|---|---|---|
| 0.58 | 0.58 | 0.58 | 0 | 0 |
| 0 | 0 | 0 | 0.71 | 0.71 |

# Latent Semantic Analysis

CS concept  MD concept

doc-to-concept similarity

| | |
|---|---|
| 0.18 | 0 |
| 0.36 | 0 |
| 0.18 | 0 |
| 0.90 | 0 |
| 0 | 0.53 |
| 0 | 0.80 |
| 0 | 0.27 |

**X**

| | |
|---|---|
| 9.64 | 0 |
| 0 | 5.29 |

**X**

| | | | | |
|---|---|---|---|---|
| 0.58 | 0.58 | 0.58 | 0 | 0 |
| 0 | 0 | 0 | 0.71 | 0.71 |

# Latent Semantic Analysis

doc-to-concept
similarity matrix

| | |
|------|------|
| 0.18 | 0 |
| 0.36 | 0 |
| 0.18 | 0 |
| 0.90 | 0 |
| 0 | 0.53 |
| 0 | 0.80 |
| 0 | 0.27 |

**X**

| | |
|------|------|
| 9.64 | 0 |
| 0 | 5.29 |

**X**

| | | | | |
|------|------|------|------|------|
| 0.58 | 0.58 | 0.58 | 0 | 0 |
| 0 | 0 | 0 | 0.71 | 0.71 |

# Latent Semantic Analysis

doc-to-concept
similarity matrix

| | |
|------|------|
| 0.18 | 0 |
| 0.36 | 0 |
| 0.18 | 0 |
| 0.90 | 0 |
| 0 | 0.53 |
| 0 | 0.80 |
| 0 | 0.27 |

**X**

"strength" of the CS concept

| | |
|------|------|
| 9.64 | 0 |
| 0 | 5.29 |

**X**

| | | | | |
|------|------|------|------|------|
| 0.58 | 0.58 | 0.58 | 0 | 0 |
| 0 | 0 | 0 | 0.71 | 0.71 |

# Latent Semantic Analysis

doc-to-concept
similarity matrix

| | |
|------|------|
| 0.18 | 0 |
| 0.36 | 0 |
| 0.18 | 0 |
| 0.90 | 0 |
| 0 | 0.53 |
| 0 | 0.80 |
| 0 | 0.27 |

**X**

"strength" of the
each concept

| | |
|------|------|
| 9.64 | 0 |
| 0 | 5.29 |

**X**

| | | | | |
|------|------|------|------|------|
| 0.58 | 0.58 | 0.58 | 0 | 0 |
| 0 | 0 | 0 | 0.71 | 0.71 |

# Latent Semantic Analysis

doc-to-concept
similarity matrix

| | |
|------|------|
| 0.18 | 0 |
| 0.36 | 0 |
| 0.18 | 0 |
| 0.90 | 0 |
| 0 | 0.53 |
| 0 | 0.80 |
| 0 | 0.27 |

**X**

"strength" of the
each concept

| | |
|------|------|
| 9.64 | 0 |
| 0 | 5.29 |

**X**

term-to-concept similarity

| | | | | |
|------|------|------|------|------|
| 0.58 | 0.58 | 0.58 | 0 | 0 |
| 0 | 0 | 0 | 0.71 | 0.71 |

# Latent Semantic Analysis

doc-to-concept
similarity matrix

| | |
|------|------|
| 0.18 | 0 |
| 0.36 | 0 |
| 0.18 | 0 |
| 0.90 | 0 |
| 0 | 0.53 |
| 0 | 0.80 |
| 0 | 0.27 |

**X**

"strength" of the
each concept

| | |
|------|------|
| 9.64 | 0 |
| 0 | 5.29 |

**X**

term-to-concept similarity
matrix

| | | | | |
|------|------|------|------|------|
| 0.58 | 0.58 | 0.58 | 0 | 0 |
| 0 | 0 | 0 | 0.71 | 0.71 |

# Latent Semantic Analysis

We can better represent each document by:

- Frequency of the word ($n_i$ / $\Sigma n_i$ )
- TfiDf

$$tf \cdot idf$$

Term frequency in the document

$$\log \left( \frac{\text{number of documents}}{\substack{\text{number of documents} \\ \text{that contain the term}}} \right)$$