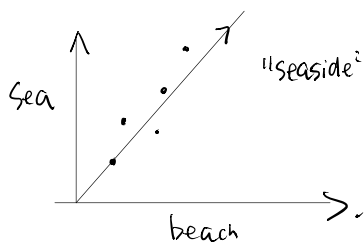


Latent Semantic Analysis (LSA)

- Discover latent topics
- Dimension Reduction
- Represent words in terms of topics

Latent Semantic Analysis



beach is often with sea

• word document matrix $A = \begin{bmatrix} w_{t1}d_1 & \dots & w_{t1}d_N \\ \vdots & & \vdots \\ v_{t1}d_1 & \dots & v_{t1}d_N \end{bmatrix}$ ↑ documents
↓

• Singular Value Decomposition (SVD) to A
 "similar to eigenvector decomposition"

$$A = USV^T$$

strength of most significant correlation

$$A = \begin{bmatrix} u_{11} & u_{12} & \dots & u_{1N} \\ \vdots & \vdots & & \vdots \\ u_{M1} & u_{M2} & \dots & u_{MN} \end{bmatrix} \begin{bmatrix} s_1 & 0 & \dots & 0 \\ 0 & s_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & \dots & 0 & s_N \end{bmatrix} \begin{bmatrix} v_{11} & \dots & v_{1N} \\ \vdots & & \vdots \\ v_{M1} & \dots & v_{MN} \end{bmatrix}$$

this is the TRANSPOSE

U, V are orthogonal matrices

$$U \cdot U^T = I$$

$$V \cdot V^T = I$$

s_1, \dots, s_N : singular values

$$s_1 \geq s_2 \geq \dots \geq s_N$$

Consider V

Columns: V dimensional unit vectors.

orthogonal to each other.

$$V = \begin{bmatrix} v_{11} & \dots & v_{1N} \\ \vdots & & \vdots \\ v_{M1} & \dots & v_{MN} \end{bmatrix} \Rightarrow \text{orthonormal basis.}$$

v_i : a document vector. \Rightarrow semantic class.

[topic]

Importance of the class: indicated by s_n .

$$v_i = \begin{bmatrix} v_{1i} \\ \vdots \\ v_{ji} \\ \vdots \\ v_{Mi} \end{bmatrix}$$

TF-IDF weight of the j^{th} term.

if v_{ji} is large

\Rightarrow this term in vocab. is significant

Consider. U

• $A v_n = U S V^T \cdot \underline{v_n} = \sum_n U_n$ Combination of documents
the n th topic

slides 15 20

topic-based representation.

$\text{vec}(d) \cdot v_n$: magnitude of the component.
in direction of v_n .

$$\text{top}(d) = \begin{bmatrix} \text{vec}(d) \cdot v_1 \\ \vdots \\ \text{vec}(d) \cdot v_n \\ \hline \text{vec}(d) \cdot v_r \end{bmatrix} \quad \left. \begin{array}{l} \text{truncate.} \\ \approx \text{top}(d) \end{array} \right\} = \vec{v(n)} \cdot \text{vec}(d)$$

↓
Reduced Dimension