

✧

Vector Notation

$$\text{vec}(d) = \begin{bmatrix} w_{1,d} & \dots & w_{m,d} \end{bmatrix}$$

$$w_{i(d),d} = \text{freq}(i, d) \times \text{IDF}(i(m))$$

$$\leftarrow V \cdot \text{(term)} \rightarrow$$

• word document matrix A

$$= \begin{bmatrix} w_{1,d_1} & \dots & w_{1,d_2} \\ \vdots & & \vdots \\ w_{m,d_1} & & w_{m,d_2} \end{bmatrix}$$

↑ documents
↓

• $W = USV^T$

PCA and LSA

• Suppose $\text{vec}(D)$ zero mean

• Subtract mean from each doc. vec. $\text{vec}(d_n)$ is zero

→ So that the mean

• Cov. matrix Σ of $\text{vec}(D)$

$$\begin{aligned} \Sigma &= \frac{1}{N-1} W^T W = \frac{1}{N-1} (USV^T)^T (USV^T) \\ &= \frac{1}{N-1} V S^T U^T U S V^T \\ &= \frac{1}{N-1} V S^2 V^T \end{aligned}$$

$$\Sigma = V \frac{S^2}{N-1} V^T = V \bar{S} V^T$$

in the form of eigenvector decomposition
(the decomp. is unique)

→ but don't get u if using PCA

description of docs in terms of topics

Summary:

Document Vectors $\text{vec}(D)$ with 0 mean vector

LSA and PCA give same orthonormal basis (V)

⇒ LSA, PCA : Same set of topics

D denotes eigenvalue matrix

$$d_{ij} = \bar{S}_{ij} = \frac{S_{ij}}{N-1}$$

Effect of subtracting Mean document vector from $\text{vec}(D)$

Normalize the vector
then get the PCA = LSA