

LSI using SVD

# Task

- SVD
- Examples on SVD
- LSI using SVD
- Code of LSI

# Singular Value Decomposition

For an  $m \times n$  matrix  $\mathbf{A}$  of rank  $r$  there exists a factorization (Singular Value Decomposition = **SVD**) as follows:

$$A = U \Sigma V^T$$

$m \times m$     $m \times n$     $V$  is  $n \times n$

The columns of  $\mathbf{U}$  are orthogonal eigenvectors of  $\mathbf{A}\mathbf{A}^T$ .

The columns of  $\mathbf{V}$  are orthogonal eigenvectors of  $\mathbf{A}^T\mathbf{A}$ .

Eigenvalues  $\lambda_1 \dots \lambda_r$  of  $\mathbf{A}\mathbf{A}^T$  are the eigenvalues of  $\mathbf{A}^T\mathbf{A}$ .

$$\sigma_i = \sqrt{\lambda_i}$$

$$\Sigma = \text{diag}(\sigma_1 \dots \sigma_r)$$

Singular values.

# • Singular Value Decomposition

- Illustration of SVD dimensions and sparseness

The top diagram illustrates the SVD of a 5x3 matrix  $A$ . The matrix  $A$  is shown as a 5x3 grid of asterisks. It is equal to the product of three matrices:  $U$  (5x5),  $\Sigma$  (5x5), and  $V^T$  (5x3). The  $U$  matrix has a shaded 5x2 region in its last two columns. The  $\Sigma$  matrix has three non-zero singular values (represented by dots) on its diagonal, with a shaded 5x2 region below the first two rows. The  $V^T$  matrix has a shaded 5x2 region in its last two columns.

The bottom diagram illustrates the SVD of a 5x5 matrix  $A$ . The matrix  $A$  is shown as a 5x5 grid of asterisks. It is equal to the product of three matrices:  $U$  (5x3),  $\Sigma$  (5x5), and  $V^T$  (5x5). The  $U$  matrix has a shaded 5x2 region in its last two columns. The  $\Sigma$  matrix has three non-zero singular values (represented by dots) on its diagonal, with a shaded 5x2 region to its right. The  $V^T$  matrix has a shaded 5x2 region in its last two columns.

# Steps to Compute SVD

1) Compute **U** matrix

1) Compute  $\mathbf{AA}^T$

2) Compute Eigen Values for  $\mathbf{AA}^T$  and arrange in decreasing order.

3) For each Eigen value compute Eigen vector and normalize it.

4) Place each Eigen vector in **U** in order of decreasing Eigen values.

2) Compute **V** matrix

$$\mathbf{V} = \mathbf{A}^T \mathbf{U} \mathbf{S}^{-1}$$

# SVD example

$$\text{Let } A = \begin{bmatrix} 1 & -1 \\ 0 & 1 \\ 1 & 0 \end{bmatrix}$$

Thus  $m=3$ ,  $n=2$ . Its SVD is

$$\begin{bmatrix} 0 & 2/\sqrt{6} & 1/\sqrt{3} \\ 1/\sqrt{2} & -1/\sqrt{6} & 1/\sqrt{3} \\ 1/\sqrt{2} & 1/\sqrt{6} & -1/\sqrt{3} \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & \sqrt{3} \\ 0 & 0 \end{bmatrix} \begin{bmatrix} 1/\sqrt{2} & 1/\sqrt{2} \\ 1/\sqrt{2} & -1/\sqrt{2} \end{bmatrix}$$

Typically, the singular values arranged in decreasing order.

# Low-rank Approximation

- SVD can be used to compute optimal **low-rank approximations**.
- Approximation problem: Find  $\mathbf{A}_k$  of rank  $k$  such that

$$\mathbf{A}_k = \min_{X: \text{rank}(X)=k} \|\mathbf{A} - X\|_F \longleftarrow \text{Frobenius norm}$$

$$\|\mathbf{A}\|_F \equiv \sqrt{\sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2}.$$

$\mathbf{A}_k$  and  $X$  are both  $m \times n$  matrices.

Typically, want  $k \ll r$ .

# LSI using SVD

- Consider 3 documents and Query:
- d1: Shipment of gold damaged in a fire.
- d2: Delivery of silver arrived in a silver truck.
- d3: Shipment of gold arrived in a truck.
- q: gold silver truck



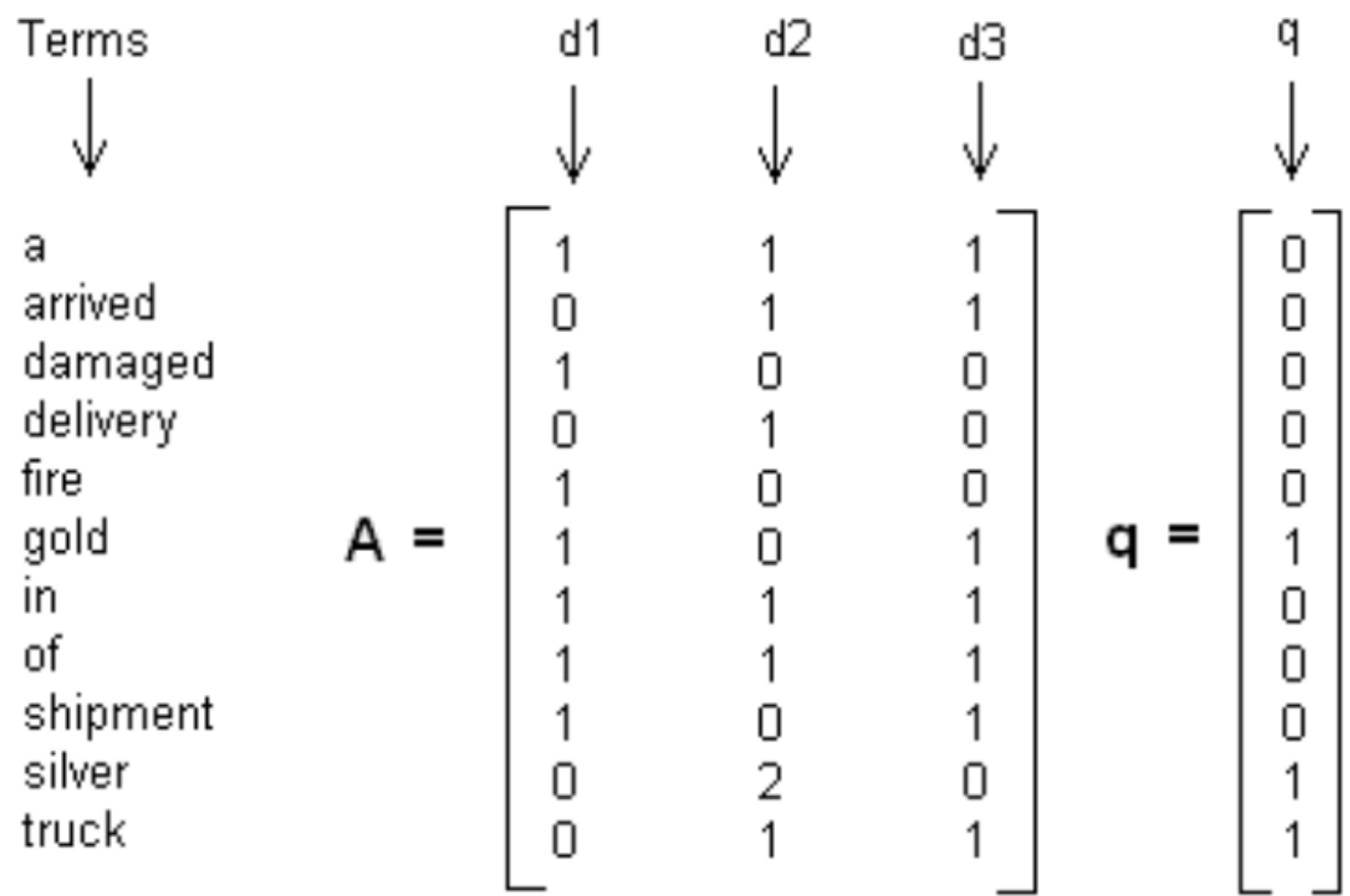


Figure 2. Term-document matrix and query matrix example.

$$\mathbf{U} = \begin{bmatrix}
 -0.4201 & 0.0748 & -0.0460 \\
 -0.2995 & -0.2001 & 0.4078 \\
 -0.1206 & 0.2749 & -0.4530 \\
 -0.1576 & -0.3046 & -0.2006 \\
 -0.1206 & 0.2749 & -0.4538 \\
 -0.2626 & 0.3794 & 0.1547 \\
 -0.4201 & 0.0748 & -0.0460 \\
 -0.4201 & 0.0748 & -0.0460 \\
 -0.2626 & 0.3794 & 0.1547 \\
 -0.3151 & -0.6093 & -0.4013 \\
 -0.2995 & -0.2001 & 0.4078
 \end{bmatrix}
 \quad
 \mathbf{S} = \begin{bmatrix}
 4.0989 & 0.0000 & 0.0000 \\
 0.0000 & 2.3616 & 0.0000 \\
 0.0000 & 0.0000 & 1.2737
 \end{bmatrix}$$
  

$$\mathbf{V} = \begin{bmatrix}
 -0.4945 & 0.6492 & -0.5780 \\
 -0.6458 & -0.7194 & -0.2556 \\
 -0.5817 & 0.2469 & 0.7750
 \end{bmatrix}
 \quad
 \mathbf{V}^T = \begin{bmatrix}
 -0.4945 & 0.6458 & -0.5817 \\
 0.6492 & -0.7194 & 0.2469 \\
 -0.5780 & -0.2556 & 0.7750
 \end{bmatrix}$$

**Figure 3. SVD results from the Bluebit Matrix Calculator.**

$$\begin{aligned}
 \mathbf{U} \approx \mathbf{U}_k &= \begin{bmatrix} -0.4201 & 0.0740 \\ -0.2995 & -0.2001 \\ -0.1206 & 0.2749 \\ -0.1576 & -0.3046 \\ -0.1206 & 0.2749 \\ -0.2626 & 0.3794 \\ -0.4201 & 0.0748 \\ -0.4201 & 0.0748 \\ -0.2626 & 0.3794 \\ -0.3151 & -0.6093 \\ -0.2995 & -0.2001 \end{bmatrix} & \mathbf{S} \approx \mathbf{S}_k &= \begin{bmatrix} 4.0989 & 0.0000 \\ 0.0000 & 2.3616 \end{bmatrix} & k = 2 \\
 \\
 \mathbf{V} \approx \mathbf{V}_k &= \begin{bmatrix} -0.4945 & 0.6492 \\ -0.6458 & -0.7194 \\ -0.5817 & 0.2469 \end{bmatrix} & \mathbf{V}^T \approx \mathbf{V}_k^T &= \begin{bmatrix} -0.4945 & -0.6458 & -0.5817 \\ 0.6492 & -0.7194 & 0.2469 \end{bmatrix}
 \end{aligned}$$

Figure 4. A Rank 2 Approximation.

# Compute Vector for Document and Query in new Vector Space

- As we know  $\mathbf{V} = \mathbf{A}^T \mathbf{U} \mathbf{S}^{-1}$
- Similarly we can compute vector for document and query as

$$\mathbf{d} = \mathbf{d}^T \mathbf{U}_k \mathbf{S}_k^{-1} \quad \mathbf{q} = \mathbf{q}^T \mathbf{U}_k \mathbf{S}_k^{-1}$$

- Compute similarity as

$$\mathbf{sim}(\mathbf{q}, \mathbf{d}) = \mathbf{sim}(\mathbf{q}^T \mathbf{U}_k \mathbf{S}_k^{-1}, \mathbf{d}^T \mathbf{U}_k \mathbf{S}_k^{-1})$$

$$\begin{aligned}
 \mathbf{q} &= \mathbf{q}^T \mathbf{U}_k \mathbf{S}_k^{-1} \\
 \mathbf{q} &= \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 1 & 1 \end{bmatrix} \begin{bmatrix} -0.4201 & 0.0748 \\ -0.2995 & -0.2001 \\ -0.1206 & 0.2749 \\ -0.1576 & -0.3046 \\ -0.1206 & 0.2749 \\ -0.2626 & 0.3794 \\ -0.4201 & 0.0748 \\ -0.4201 & 0.0748 \\ -0.2626 & 0.3794 \\ -0.3151 & -0.6093 \\ -0.2995 & -0.2001 \end{bmatrix} \begin{bmatrix} 1 & \\ 4.0989 & 0.0000 \\ & 1 \\ 0.0000 & 2.3616 \end{bmatrix} \\
 \mathbf{q} &= \begin{bmatrix} -0.2140 & -0.1821 \end{bmatrix}
 \end{aligned}$$

$k = 2$

Figure 5. Computing the query vector

Code LSI