

# Matrix Decompositions and Latent Semantic Indexing (LSI)

October 24, 2024

The chapter you're referring to appears to be about **Matrix Decompositions and Latent Semantic Indexing (LSI)**, as described in the context of information retrieval. Let's break it down into three parts, as you requested. The first part will cover the introduction, linear algebra review, and matrix decompositions.

## Part 1: Introduction to Matrix Decompositions and LSI

### 1. Term-Document Matrix (TDM):

- A **term-document matrix (TDM)** is a large matrix representing documents and their associated terms, where each row corresponds to a term and each column corresponds to a document.
- In information retrieval, even modest-sized collections produce matrices with thousands of rows and columns, making them quite large.
- This chapter introduces matrix decomposition techniques to manage these large matrices efficiently by reducing their dimensions.

### 2. Linear Algebra Review:

- **Basic Concepts:** The chapter starts with a refresher on linear algebra concepts, such as **matrix rank**, **eigenvalues**, and **eigenvectors**.
  - The **rank** of a matrix is the number of linearly independent rows or columns in the matrix.
  - An **eigenvalue** is a scalar that reflects how a matrix scales a particular vector (its eigenvector).
  - The **principal eigenvector** is the eigenvector corresponding to the eigenvalue with the largest magnitude.

- **Example 18.1** demonstrates how multiplying a vector by a matrix is primarily influenced by its largest eigenvalues, which is a key idea that leads to dimensionality reduction later in the chapter.

### 3. Matrix Decompositions:

- Matrix decompositions allow us to break down complex matrices into simpler components for better processing. The **matrix diagonalization theorem** is introduced as a key tool:
  - If a matrix has  $M$  linearly independent eigenvectors, it can be factorized into three matrices:  $S = U\Lambda U^{-1}$ , where  $U$  contains the eigenvectors of  $S$ , and  $\Lambda$  is a diagonal matrix with the eigenvalues of  $S$ .
- **Symmetric Diagonalization:** The theorem also covers symmetric matrices, where the decomposition becomes  $S = Q\Lambda Q^T$ , with **orthogonal** eigenvectors.

### 4. Relevance to Information Retrieval:

- These decompositions are the foundation of more advanced matrix techniques that help reduce the size of the term-document matrix without losing significant information.
- Matrix decompositions, such as **singular value decomposition (SVD)**, are applied to text analysis and clustering tasks.

In summary, **Part 1** of the chapter sets up the foundational mathematics required for understanding **matrix decompositions** and how they are applied to the **term-document matrix** in **information retrieval systems**. This groundwork leads into **singular value decomposition (SVD)** and its application to information retrieval through **latent semantic indexing**, which will be discussed in **Part 2**.

## Part 2: Term-Document Matrices, Singular Value Decompositions (SVD), and Low-Rank Approximations

This section dives into how matrix decompositions, particularly **singular value decomposition (SVD)**, are applied to **term-document matrices** to reduce their complexity, making information retrieval more efficient.

## 1. Term-Document Matrices and Singular Value Decompositions (SVD)

- In information retrieval, the **term-document matrix** ( $C$ ) is typically not square, meaning the number of terms (rows) does not equal the number of documents (columns). This makes simple matrix diagonalization (like those discussed in Part 1) inapplicable.
- **Singular Value Decomposition (SVD)** is introduced as an extension of matrix decomposition for non-square matrices. The idea behind SVD is that any matrix can be factored into three matrices:

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

Where:

- $U$  is an  $M \times M$  matrix containing the left singular vectors (orthogonal eigenvectors) of  $CC^T$  (corresponding to the terms).
- $V$  is an  $N \times N$  matrix containing the right singular vectors (orthogonal eigenvectors) of  $C^TC$  (corresponding to the documents).
- $\Sigma$  is an  $M \times N$  diagonal matrix that contains the **singular values** of  $C$ , which are square roots of the eigenvalues.

## 2. Low-Rank Approximations

- Once we have the singular value decomposition of the term-document matrix, we can approximate it by focusing only on the largest singular values and ignoring the smaller ones. This results in a **low-rank approximation**:

$$C_k = U_k \Sigma_k V_k^T$$

Where:

- $C_k$  is the approximation of the original matrix  $C$ , retaining only the **top  $k$  singular values**.
- This approximation compresses the information, keeping the most significant features while discarding noise or less important data.
- **Frobenius Norm:** To measure how good the approximation is, we use the **Frobenius norm**. The goal is to minimize the difference between the original matrix  $C$  and the approximated matrix  $C_k$ , as measured by the Frobenius norm:

$$\|C - C_k\|_F = \sqrt{\sum (C_{ij} - C_{k,ij})^2}$$

- The theorem by Eckart and Young states that  $C_k$  is the best possible rank- $k$  approximation to  $C$  that minimizes the Frobenius norm error.

### 3. Intuition Behind Low-Rank Approximations

- The main idea behind low-rank approximation is that the **largest singular values** capture the most essential structure of the term-document relationships, while the smaller singular values capture noise or less important information.
- By keeping only the top  $k$  singular values, we create a matrix that retains the important co-occurrence patterns between terms and documents, while the smaller details (captured by smaller singular values) are ignored.

In summary, **Part 2** shows how **SVD** is applied to the **term-document matrix**, allowing efficient storage, retrieval, and processing of large text collections. This prepares for **latent semantic indexing (LSI)**, discussed next.

## Part 3: Latent Semantic Indexing (LSI) and Applications

This section explores how the concepts of **SVD** and **low-rank approximations** improve retrieval systems using **LSI**.

### 1. Latent Semantic Indexing (LSI)

- **Latent Semantic Indexing (LSI)** uses **SVD** to reduce the dimensionality of the term-document matrix, capturing its underlying structure.
- The low-rank approximation  $C_k$  is used to represent both documents and queries in a reduced semantic space, enabling:
  - **Latent associations** between terms and documents.
  - **Noise reduction** due to natural language variability.
- **Synonymy and Polysemy**: LSI helps with **synonymy** (e.g., "car" and "automobile") by aligning similar terms in the reduced space, and **polysemy** (e.g., "bank") by using co-occurrence patterns to infer context.

### 2. How LSI Improves Retrieval Performance

- **Vector Space Model**: Both documents and queries are represented as vectors, and their similarity is calculated using **cosine similarity**:

$$\text{Similarity}(\mathbf{q}, \mathbf{d}) = \frac{\mathbf{q} \cdot \mathbf{d}}{\|\mathbf{q}\| \|\mathbf{d}\|}$$

- **Improvement**: LSI captures underlying semantic relationships, improving retrieval performance by better matching conceptually related documents and queries.

### 3. Practical Applications and Challenges

- **Performance:** LSI has been shown to improve precision and recall in tasks like TREC.
- **Challenges:** The computational cost of **SVD** limits the scalability of LSI, especially for large collections. One solution is random sampling.

### 4. LSI Beyond Text Retrieval

- LSI has applications in **soft clustering**, **cross-language information retrieval**, and other areas like **memory modeling** and **computer vision**.

In conclusion, **LSI** improves information retrieval by capturing latent relationships between terms, addressing challenges like synonymy and polysemy, and extending beyond text retrieval tasks.