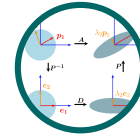


Matrix Decompositions



In Chapters 2 and 3, we studied ways to manipulate and measure vectors, projections of vectors and linear mappings. Mappings and transformations of vectors can be conveniently described as operations performed by matrices. Moreover, data is often represented in matrix form as well, e.g., where the rows of the matrix represent different people and the columns describe different features of the people, such as weight, height and socioeconomic status. In this chapter, we present three aspects of matrices: how to summarize matrices, how matrices can be decomposed, and how these decompositions can be used for matrix approximations.

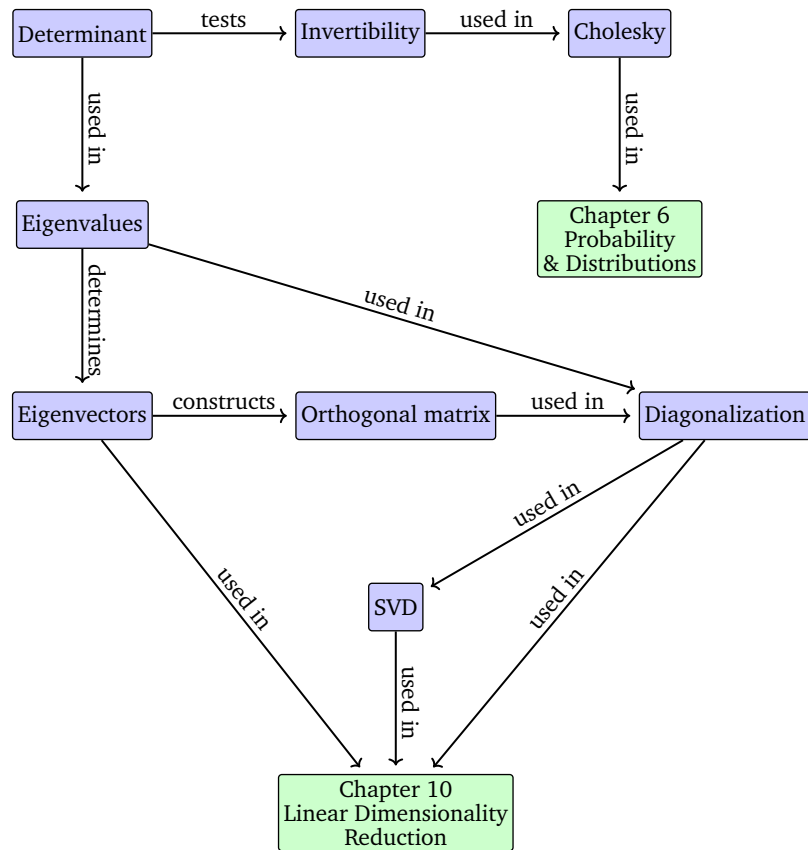
We first consider methods that allow us to describe matrices with just a few numbers that characterize the overall properties of matrices. We will do this in the sections on determinants (Section 4.1) and eigenvalues (Section 4.2) for the important special case of square matrices. These characteristic numbers have important mathematical consequences and allow us to quickly grasp what useful properties a matrix has. From here we will proceed to matrix decomposition methods: An analogy for matrix decomposition is the factoring of numbers, such as the factoring of 21 into prime numbers $7 \cdot 3$. For this reason matrix decomposition is also often referred to as *matrix factorization*. Matrix decompositions are used to describe a matrix by means of a different representation using factors of interpretable matrices.

matrix factorization

We will first cover a square-root-like operation for symmetric, positive definite matrices, the Cholesky decomposition (Section 4.3). From here we will look at two related methods for factorizing matrices into canonical forms. The first one is known as matrix diagonalization (Section 4.4), which allows us to represent the linear mapping using a diagonal transformation matrix if we choose an appropriate basis. The second method, singular value decomposition (Section 4.5), extends this factorization to non-square matrices, and it is considered one of the fundamental concepts in linear algebra. These decompositions are helpful as matrices representing numerical data are often very large and hard to analyze. We conclude the chapter with a systematic overview of the types of matrices and the characteristic properties that distinguish them in form of a matrix taxonomy (Section 4.7).

The methods that we cover in this chapter will become important in

Figure 4.1 A mind map of the concepts introduced in this chapter, along with where they are used in other parts of the book.



both subsequent mathematical chapters, such as Chapter 6 but also in applied chapters, such as dimensionality reduction in Chapters 10 or density estimation in Chapter 11. This chapter's overall structure is depicted in the mind map of Figure 4.1.

4.1 Determinant and Trace

The determinant notation $|\mathbf{A}|$ must not be confused with the absolute value.

Determinants are important concepts in linear algebra. A determinant is a mathematical object in the analysis and solution of systems of linear equations. Determinants are only defined for square matrices $\mathbf{A} \in \mathbb{R}^{n \times n}$, i.e., matrices with the same number of rows and columns. In this book, we write the determinant as $\det(\mathbf{A})$ or sometimes as $|\mathbf{A}|$ so that

$$\det(\mathbf{A}) = \begin{vmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \vdots & & \ddots & \vdots \\ a_{n1} & a_{n2} & \dots & a_{nn} \end{vmatrix}. \quad (4.1)$$

determinant

The *determinant* of a square matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ is a function that maps \mathbf{A}

onto a real number. Before providing a definition of the determinant for general $n \times n$ matrices let us have a look at some motivating examples, and define determinants for some special matrices.

Example 4.1 (Testing for Matrix Invertibility)

Let us begin with exploring if a square matrix \mathbf{A} is invertible (see Section 2.2.2). For the smallest cases, we already know when a matrix is invertible. If \mathbf{A} is a 1×1 matrix, i.e., it is a scalar number, then $\mathbf{A} = a \implies \mathbf{A}^{-1} = \frac{1}{a}$. Thus $a \frac{1}{a} = 1$ holds, if and only if $a \neq 0$.

For 2×2 matrices, by the definition of the inverse (Definition 2.3), we know that $\mathbf{A}\mathbf{A}^{-1} = \mathbf{I}$. Then, with (2.24), the inverse of \mathbf{A} is

$$\mathbf{A}^{-1} = \frac{1}{a_{11}a_{22} - a_{12}a_{21}} \begin{bmatrix} a_{22} & -a_{12} \\ -a_{21} & a_{11} \end{bmatrix}. \quad (4.2)$$

Hence, \mathbf{A} is invertible if and only if

$$a_{11}a_{22} - a_{12}a_{21} \neq 0. \quad (4.3)$$

This quantity is the determinant of $\mathbf{A} \in \mathbb{R}^{2 \times 2}$, i.e.,

$$\det(\mathbf{A}) = \begin{vmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{vmatrix} = a_{11}a_{22} - a_{12}a_{21}. \quad (4.4)$$

The example above points already at the relationship between determinants and the existence of inverse matrices. The next theorem states the same result for $n \times n$ matrices.

Theorem 4.1. *For any square matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ it holds that \mathbf{A} is invertible if and only if $\det(\mathbf{A}) \neq 0$.*

We have explicit (closed form) expressions for determinants of small matrices in terms of the elements of the matrix. For $n = 1$,

$$\det(\mathbf{A}) = \det(a_{11}) = a_{11}. \quad (4.5)$$

For $n = 2$,

$$\det(\mathbf{A}) = \begin{vmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{vmatrix} = a_{11}a_{22} - a_{12}a_{21}, \quad (4.6)$$

which we have observed in the example above. For $n = 3$ (known as Sarrus' rule),

$$\begin{vmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{vmatrix} = a_{11}a_{22}a_{33} + a_{21}a_{32}a_{13} + a_{31}a_{12}a_{23} \\ - a_{31}a_{22}a_{13} - a_{11}a_{32}a_{23} - a_{21}a_{12}a_{33}. \quad (4.7)$$

For a memory aid of the product terms in Sarrus' rule, try tracing the elements of the triple products in the matrix.

We call a square matrix T an *upper triangular matrix* if $T_{ij} = 0$ for $i > j$, that is the matrix is zero below its diagonal. Analogously, we define a *lower triangular matrix* as a matrix with zeros above its diagonal. For an upper/lower triangular matrix $T \in \mathbb{R}^{n \times n}$, the determinant is the product of the diagonal elements, i.e.,

$$\det(T) = \prod_{i=1}^n T_{ii}. \quad (4.8)$$

upper triangular
matrix
lower triangular
matrix

The determinant is the signed volume of the parallelepiped formed by the columns of the matrix.

Figure 4.2 The area of the parallelogram (shaded region) spanned by the vectors \mathbf{b} and \mathbf{g} is $|\det([\mathbf{b}, \mathbf{g}])|$.

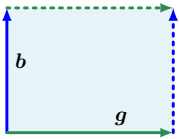
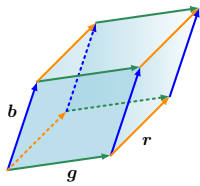


Figure 4.3 The volume of the parallelepiped (shaded volume) spanned by vectors $\mathbf{r}, \mathbf{b}, \mathbf{g}$ is $|\det([\mathbf{r}, \mathbf{b}, \mathbf{g}])|$.



The sign of the determinant indicates the orientation of the spanning vectors.

Example 4.2 (Determinants as Measures of Volume)

The notion of a determinant is natural when we consider it as a mapping from a set of n vectors spanning an object in \mathbb{R}^n . It turns out that the determinant $\det(\mathbf{A})$ is the signed volume of an n -dimensional parallelepiped formed by columns of the matrix \mathbf{A} .

For $n = 2$ the columns of the matrix form a parallelogram, see Figure 4.2. As the angle between vectors gets smaller the area of a parallelogram shrinks, too. Consider two vectors \mathbf{b}, \mathbf{g} that form the columns of a matrix $\mathbf{A} = [\mathbf{b}, \mathbf{g}]$. Then, the absolute value of the determinant of \mathbf{A} is the area of the parallelogram with vertices $0, \mathbf{b}, \mathbf{g}, \mathbf{b} + \mathbf{g}$. In particular, if \mathbf{b}, \mathbf{g} are linearly dependent so that $\mathbf{b} = \lambda \mathbf{g}$ for some $\lambda \in \mathbb{R}$ they no longer form a two-dimensional parallelogram. Therefore, the corresponding area is 0. On the contrary, if \mathbf{b}, \mathbf{g} are linearly independent and are multiples of the canonical basis vectors $\mathbf{e}_1, \mathbf{e}_2$ then they can be written as $\mathbf{b} = \begin{bmatrix} b \\ 0 \end{bmatrix}$ and

$\mathbf{g} = \begin{bmatrix} 0 \\ g \end{bmatrix}$, and the determinant

$$\begin{vmatrix} b & 0 \\ 0 & g \end{vmatrix} = bg - 0 = bg \quad (4.9)$$

The sign of the determinant indicates the orientation of the spanning vectors \mathbf{b}, \mathbf{g} with respect to the standard basis $(\mathbf{e}_1, \mathbf{e}_2)$. In our figure, flipping the order to \mathbf{g}, \mathbf{b} swaps the columns of \mathbf{A} and reverses the orientation of the shaded area. becomes the familiar formula: area = height \times length. This intuition extends to higher dimensions. In \mathbb{R}^3 , we consider three vectors $\mathbf{r}, \mathbf{b}, \mathbf{g} \in \mathbb{R}^3$ spanning the edges of a parallelepiped, i.e., a solid with faces that are parallel parallelograms (see Figure 4.3). The absolute value of the determinant of the 3×3 matrix $[\mathbf{r}, \mathbf{b}, \mathbf{g}]$ is the volume of the solid. Thus, the determinant acts as a function that measures the signed volume formed by column vectors composed in a matrix.

Consider the three linearly independent vectors $\mathbf{r}, \mathbf{g}, \mathbf{b} \in \mathbb{R}^3$ given as

$$\mathbf{r} = \begin{bmatrix} 2 \\ 0 \\ -8 \end{bmatrix}, \quad \mathbf{g} = \begin{bmatrix} 6 \\ 1 \\ 0 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} 1 \\ 4 \\ -1 \end{bmatrix}. \quad (4.10)$$

Writing these vectors as the columns of a matrix

$$\mathbf{A} = [\mathbf{r}, \mathbf{g}, \mathbf{b}] = \begin{bmatrix} 2 & 6 & 1 \\ 0 & 1 & 4 \\ -8 & 0 & -1 \end{bmatrix} \quad (4.11)$$

allows us to compute the desired volume as

$$V = |\det(\mathbf{A})| = 186. \quad (4.12)$$

Computing the determinant of an $n \times n$ matrix requires a general algorithm to solve the cases for $n > 3$, which we are going to explore in the following. The theorem below reduces the problem of computing the determinant of an $n \times n$ matrix to computing the determinant of $(n-1) \times (n-1)$ matrices. By recursively applying the following Laplace expansion we can therefore compute determinants of $n \times n$ matrices by ultimately computing determinants of 2×2 matrices.

Theorem 4.2 (Laplace Expansion). Consider a matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$. Then, for all $j = 1, \dots, n$:

1 Expansion along column j

$$\det(\mathbf{A}) = \sum_{k=1}^n (-1)^{k+j} a_{kj} \det(\mathbf{A}_{k,j}). \quad (4.13)$$

$\det(\mathbf{A}_{k,j})$ is called a *minor* and $(-1)^{k+j} \det(\mathbf{A}_{k,j})$ a *cofactor*.

2 Expansion along row j

$$\det(\mathbf{A}) = \sum_{k=1}^n (-1)^{k+j} a_{jk} \det(\mathbf{A}_{j,k}). \quad (4.14)$$

Here $\mathbf{A}_{k,j} \in \mathbb{R}^{(n-1) \times (n-1)}$ is the submatrix of \mathbf{A} that we obtain when deleting row k and column j .

Example 4.3 (Laplace Expansion)

Let us compute the determinant of

$$\mathbf{A} = \begin{bmatrix} 1 & 2 & 3 \\ 3 & 1 & 2 \\ 0 & 0 & 1 \end{bmatrix} \quad (4.15)$$

using the Laplace expansion along the first row. By applying (4.14) we get

$$\begin{vmatrix} 1 & 2 & 3 \\ 3 & 1 & 2 \\ 0 & 0 & 1 \end{vmatrix} = (-1)^{1+1} \cdot 1 \begin{vmatrix} 1 & 2 \\ 0 & 1 \end{vmatrix} \\ + (-1)^{1+2} \cdot 2 \begin{vmatrix} 3 & 2 \\ 0 & 1 \end{vmatrix} + (-1)^{1+3} \cdot 3 \begin{vmatrix} 3 & 1 \\ 0 & 0 \end{vmatrix}. \quad (4.16)$$

We use (4.6) to compute the determinants of all 2×2 matrices and obtain

$$\det(\mathbf{A}) = 1(1 - 0) - 2(3 - 0) + 3(0 - 0) = -5. \quad (4.17)$$

For completeness we can compare this result to computing the determinant using Sarrus' rule (4.7):

$$\det(\mathbf{A}) = 1 \cdot 1 \cdot 1 + 3 \cdot 0 \cdot 3 + 0 \cdot 2 \cdot 2 - 0 \cdot 1 \cdot 3 - 1 \cdot 0 \cdot 2 - 3 \cdot 2 \cdot 1 = 1 - 6 = -5. \quad (4.18)$$

For $\mathbf{A} \in \mathbb{R}^{n \times n}$ the determinant exhibits the following properties:

- The determinant of a matrix product is the product of the corresponding determinants, $\det(\mathbf{AB}) = \det(\mathbf{A})\det(\mathbf{B})$.
- Determinants are invariant to transposition, i.e., $\det(\mathbf{A}) = \det(\mathbf{A}^\top)$.
- If \mathbf{A} is regular (invertible) then $\det(\mathbf{A}^{-1}) = \frac{1}{\det(\mathbf{A})}$.
- Similar matrices (Definition 2.21) possess the same determinant. Therefore, for a linear mapping $\Phi : V \rightarrow V$ all transformation matrices \mathbf{A}_Φ of Φ have the same determinant. Thus, the determinant is invariant to the choice of basis of a linear mapping.
- Adding a multiple of a column/row to another one does not change $\det(\mathbf{A})$.
- Multiplication of a column/row with $\lambda \in \mathbb{R}$ scales $\det(\mathbf{A})$ by λ . In particular, $\det(\lambda \mathbf{A}) = \lambda^n \det(\mathbf{A})$.
- Swapping two rows/columns changes the sign of $\det(\mathbf{A})$.

Because of the last three properties, we can use Gaussian elimination (see Section 2.1) to compute $\det(\mathbf{A})$ by bringing \mathbf{A} into row-echelon form. We can stop Gaussian elimination when we have \mathbf{A} in a triangular form where the elements below the diagonal are all 0. Recall from (4.8) that the determinant of a triangular matrix is the product of the diagonal elements.

Theorem 4.3. *A square matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ has $\det(\mathbf{A}) \neq 0$ if and only if $\text{rk}(\mathbf{A}) = n$. In other words, \mathbf{A} is invertible if and only if it is full rank.*

When mathematics was mainly performed by hand, the determinant calculation was considered an essential way to analyze matrix invertibility. However, contemporary approaches in machine learning use direct numerical methods that superseded the explicit calculation of the determinant. For example, in Chapter 2, we learned that inverse matrices can

be computed by Gaussian elimination. Gaussian elimination can thus be used to compute the determinant of a matrix.

Determinants will play an important theoretical role for the following sections, especially when we learn about eigenvalues and eigenvectors (Section 4.2) through the characteristic polynomial.

Definition 4.4. The *trace* of a square matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ is defined as

trace

$$\text{tr}(\mathbf{A}) := \sum_{i=1}^n a_{ii}, \quad (4.19)$$

i.e., the trace is the sum of the diagonal elements of \mathbf{A} .

The trace satisfies the following properties:

- $\text{tr}(\mathbf{A} + \mathbf{B}) = \text{tr}(\mathbf{A}) + \text{tr}(\mathbf{B})$ for $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{n \times n}$
- $\text{tr}(\alpha \mathbf{A}) = \alpha \text{tr}(\mathbf{A})$, $\alpha \in \mathbb{R}$ for $\mathbf{A} \in \mathbb{R}^{n \times n}$
- $\text{tr}(\mathbf{I}_n) = n$
- $\text{tr}(\mathbf{AB}) = \text{tr}(\mathbf{BA})$ for $\mathbf{A} \in \mathbb{R}^{n \times k}$, $\mathbf{B} \in \mathbb{R}^{k \times n}$

It can be shown that only one function satisfies these four properties together – the trace (Gohberg et al., 2012).

The properties of the trace of matrix products are more general. Specifically, the trace is invariant under cyclic permutations, i.e.,

The trace is invariant under cyclic permutations.

$$\text{tr}(\mathbf{AKL}) = \text{tr}(\mathbf{KLA}) \quad (4.20)$$

for matrices $\mathbf{A} \in \mathbb{R}^{a \times k}$, $\mathbf{K} \in \mathbb{R}^{k \times l}$, $\mathbf{L} \in \mathbb{R}^{l \times a}$. This property generalizes to products of arbitrarily many matrices. As a special case of (4.20) it follows that for two vectors $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$

$$\text{tr}(\mathbf{xy}^\top) = \text{tr}(\mathbf{y}^\top \mathbf{x}) = \mathbf{y}^\top \mathbf{x} \in \mathbb{R}. \quad (4.21)$$

Given a linear mapping $\Phi : V \rightarrow V$, where V is a vector space, we define the trace of this map by using the trace of matrix representation of Φ . For a given basis of V we can describe Φ by means of the transformation matrix \mathbf{A} . Then, the trace of Φ is the trace of \mathbf{A} . For a different basis of V it holds that the corresponding transformation matrix \mathbf{B} of Φ can be obtained by a basis change of the form $\mathbf{S}^{-1} \mathbf{A} \mathbf{S}$ for suitable \mathbf{S} (see Section 2.7.2). For the corresponding trace of Φ this means

$$\text{tr}(\mathbf{B}) = \text{tr}(\mathbf{S}^{-1} \mathbf{A} \mathbf{S}) \stackrel{(4.20)}{=} \text{tr}(\mathbf{A} \mathbf{S} \mathbf{S}^{-1}) = \text{tr}(\mathbf{A}). \quad (4.22)$$

Hence, while matrix representations of linear mappings are basis dependent the trace of a linear mapping Φ is independent of the basis.

In this section, we covered determinants and traces as functions characterizing a square matrix. Taking together our understanding of determinants and traces we can now define an important equation describing a matrix \mathbf{A} in terms of a polynomial, which we will use extensively in the following sections.

Definition 4.5 (Characteristic Polynomial). For $\lambda \in \mathbb{R}$ and a square matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$

$$p_{\mathbf{A}}(\lambda) = \det(\mathbf{A} - \lambda \mathbf{I}) \quad (4.23a)$$

$$= c_0 + c_1 \lambda + c_2 \lambda^2 + \cdots + c_{n-1} \lambda^{n-1} + (-1)^n \lambda^n, \quad (4.23b)$$

characteristic
polynomial

$c_0, \dots, c_{n-1} \in \mathbb{R}$, is the *characteristic polynomial* of \mathbf{A} . In particular,

$$c_0 = \det(\mathbf{A}), \quad (4.24)$$

$$c_{n-1} = (-1)^{n-1} \text{tr}(\mathbf{A}). \quad (4.25)$$

The characteristic polynomial (4.23a) will allow us to compute eigenvalues and eigenvectors, covered in the next section.

4.2 Eigenvalues and Eigenvectors

We will now get to know a new way to characterize a matrix and its associated linear mapping. Recall from Section 2.7.1 that every linear mapping has a unique transformation matrix given an ordered basis. We can interpret linear mappings and their associated transformation matrices by performing an “eigen” analysis. As we will see, the eigenvalues of a linear mapping will tell us how a special set of vectors, the eigenvectors, are transformed by the linear mapping.

Eigen is a German word meaning “characteristic”, “self” or “own”.

Definition 4.6. Let $\mathbf{A} \in \mathbb{R}^{n \times n}$ be a square matrix. Then $\lambda \in \mathbb{R}$ is an *eigenvalue* of \mathbf{A} and $\mathbf{x} \in \mathbb{R}^n \setminus \{\mathbf{0}\}$ is the corresponding *eigenvector* of \mathbf{A} if

eigenvalue
eigenvector

$$\mathbf{A}\mathbf{x} = \lambda\mathbf{x}. \quad (4.26)$$

eigenvalue equation

We call (4.26) the *eigenvalue equation*.

Remark. In the linear algebra literature and software, it is often a convention that eigenvalues are sorted in descending order, so that the largest eigenvalue and associated eigenvector are called the first eigenvalue and its associated eigenvector, and the second largest called the second eigenvalue and its associated eigenvector, and so on. However, textbooks and publications may have different or no notion of orderings. We do not want to presume an ordering in this book if not stated explicitly. \diamond

The following statements are equivalent:

- λ is an eigenvalue of $\mathbf{A} \in \mathbb{R}^{n \times n}$
- There exists an $\mathbf{x} \in \mathbb{R}^n \setminus \{\mathbf{0}\}$ with $\mathbf{A}\mathbf{x} = \lambda\mathbf{x}$ or equivalently, $(\mathbf{A} - \lambda\mathbf{I}_n)\mathbf{x} = \mathbf{0}$ can be solved non-trivially, i.e., $\mathbf{x} \neq \mathbf{0}$.
- $\text{rk}(\mathbf{A} - \lambda\mathbf{I}_n) < n$
- $\det(\mathbf{A} - \lambda\mathbf{I}_n) = 0$

codirected
collinear

Definition 4.7 (Collinearity and Codirection). Two vectors that point in the same direction are called *codirected*. Two vectors are *collinear* if they point in the same or the opposite direction.

Remark (Non-uniqueness of eigenvectors). If \mathbf{x} is an eigenvector of \mathbf{A} associated with eigenvalue λ then for any $c \in \mathbb{R} \setminus \{0\}$ it holds that $c\mathbf{x}$ is an eigenvector of \mathbf{A} with the same eigenvalue since

$$\mathbf{A}(c\mathbf{x}) = c\mathbf{A}\mathbf{x} = c\lambda\mathbf{x} = \lambda(c\mathbf{x}). \quad (4.27)$$

Thus, all vectors that are collinear to \mathbf{x} are also eigenvectors of \mathbf{A} . \diamond

Theorem 4.8. $\lambda \in \mathbb{R}$ is eigenvalue of $\mathbf{A} \in \mathbb{R}^{n \times n}$ if and only if λ is a root of the characteristic polynomial $p_{\mathbf{A}}(\lambda)$ of \mathbf{A} .

Definition 4.9. Let a square matrix \mathbf{A} have an eigenvalue λ_i . The *algebraic multiplicity* of λ_i is the number of times the root appears in the characteristic polynomial. algebraic
multiplicity

Definition 4.10 (Eigenspace and Eigenspectrum). For $\mathbf{A} \in \mathbb{R}^{n \times n}$ the set of all eigenvectors of \mathbf{A} associated with an eigenvalue λ spans a subspace of \mathbb{R}^n , which is called the *eigenspace* of \mathbf{A} with respect to λ and is denoted by E_{λ} . The set of all eigenvalues of \mathbf{A} is called the *eigenspectrum*, or just *spectrum*, of \mathbf{A} . eigenspace
eigenspectrum
spectrum

If λ is an eigenvalue of $\mathbf{A} \in \mathbb{R}^{n \times n}$ then the corresponding eigenspace E_{λ} is the solution space of the homogeneous system of linear equations $(\mathbf{A} - \lambda\mathbf{I})\mathbf{x} = \mathbf{0}$. Geometrically, the eigenvector corresponding to a non-zero eigenvalue points in a direction that is stretched by the linear mapping. The eigenvalue is the factor by which it is stretched. If the eigenvalue is negative, the direction of the stretching is flipped.

Example 4.4 (The Case of the Identity Matrix)

The identity matrix $\mathbf{I} \in \mathbb{R}^{n \times n}$ has characteristic polynomial $p_{\mathbf{I}}(\lambda) = \det(\mathbf{I} - \lambda\mathbf{I}) = (1 - \lambda)^n = 0$, which has only one eigenvalue $\lambda = 1$ that occurs n times. Moreover, $\mathbf{I}\mathbf{x} = \lambda\mathbf{x} = 1\mathbf{x}$ holds for all vectors $\mathbf{x} \in \mathbb{R}^n \setminus \{\mathbf{0}\}$. Because of this, the sole eigenspace E_1 of the identity matrix spans n dimensions, and all n standard basis vectors of \mathbb{R}^n are eigenvectors of \mathbf{I} .

Useful properties regarding eigenvalues and eigenvectors include:

- A matrix \mathbf{A} and its transpose \mathbf{A}^{\top} possess the same eigenvalues, but not necessarily the same eigenvectors.
- The eigenspace E_{λ} is the null space of $\mathbf{A} - \lambda\mathbf{I}$ since

$$\mathbf{A}\mathbf{x} = \lambda\mathbf{x} \iff \mathbf{A}\mathbf{x} - \lambda\mathbf{x} = \mathbf{0} \quad (4.28a)$$

$$\iff (\mathbf{A} - \lambda\mathbf{I})\mathbf{x} = \mathbf{0} \iff \mathbf{x} \in \ker(\mathbf{A} - \lambda\mathbf{I}). \quad (4.28b)$$

- Similar matrices (see Definition 2.21) possess the same eigenvalues. Therefore, a linear mapping Φ has eigenvalues that are independent of the choice of basis of its transformation matrix. This makes eigenvalues,

together with the determinant and the trace, key characteristic parameters of a linear mapping as they are all invariant under basis change.

- Symmetric, positive definite matrices always have positive, real eigenvalues.

Example 4.5 (Computing Eigenvalues, Eigenvectors and Eigenspaces)

Let us find the eigenvalues and eigenvectors of the 2×2 matrix

$$\mathbf{A} = \begin{bmatrix} 4 & 2 \\ 1 & 3 \end{bmatrix}. \quad (4.29)$$

Step 1: Characteristic Polynomial. From our definition of the eigenvector $\mathbf{x} \neq \mathbf{0}$ and eigenvalue λ of \mathbf{A} there will be a vector such that $\mathbf{Ax} = \lambda\mathbf{x}$, i.e., $(\mathbf{A} - \lambda\mathbf{I})\mathbf{x} = \mathbf{0}$. Since $\mathbf{x} \neq \mathbf{0}$ this requires that the kernel (null space) of $\mathbf{A} - \lambda\mathbf{I}$ contains more elements than just $\mathbf{0}$. This means that $\mathbf{A} - \lambda\mathbf{I}$ is not invertible and therefore $\det(\mathbf{A} - \lambda\mathbf{I}) = 0$. Hence, we need to compute the roots of the characteristic polynomial (4.23a) to find the eigenvalues.

Step 2: Eigenvalues. The characteristic polynomial is

$$p_{\mathbf{A}}(\lambda) = \det(\mathbf{A} - \lambda\mathbf{I}) \quad (4.30a)$$

$$= \det\left(\begin{bmatrix} 4 & 2 \\ 1 & 3 \end{bmatrix} - \begin{bmatrix} \lambda & 0 \\ 0 & \lambda \end{bmatrix}\right) = \begin{vmatrix} 4 - \lambda & 2 \\ 1 & 3 - \lambda \end{vmatrix} \quad (4.30b)$$

$$= (4 - \lambda)(3 - \lambda) - 2 \cdot 1. \quad (4.30c)$$

We factorize the characteristic polynomial and obtain

$$p(\lambda) = (4 - \lambda)(3 - \lambda) - 2 \cdot 1 = 10 - 7\lambda + \lambda^2 = (2 - \lambda)(5 - \lambda) \quad (4.31)$$

giving the roots $\lambda_1 = 2$ and $\lambda_2 = 5$.

Step 3: Eigenvectors and Eigenspaces. We find the eigenvectors that correspond to these eigenvalues by looking at vectors \mathbf{x} such that

$$\begin{bmatrix} 4 - \lambda & 2 \\ 1 & 3 - \lambda \end{bmatrix} \mathbf{x} = \mathbf{0}. \quad (4.32)$$

For $\lambda = 5$ we obtain

$$\begin{bmatrix} 4 - 5 & 2 \\ 1 & 3 - 5 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} -1 & 2 \\ 1 & -2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \mathbf{0}. \quad (4.33)$$

We solve this homogeneous system and obtain a solution space

$$E_5 = \text{span}\left[\begin{bmatrix} 2 \\ 1 \end{bmatrix}\right]. \quad (4.34)$$

This eigenspace is one-dimensional as it possesses a single basis vector.

Analogously, we find the eigenvector for $\lambda = 2$ by solving the homogeneous system of equations

$$\begin{bmatrix} 4-2 & 2 \\ 1 & 3-2 \end{bmatrix} \mathbf{x} = \begin{bmatrix} 2 & 2 \\ 1 & 1 \end{bmatrix} \mathbf{x} = \mathbf{0}. \quad (4.35)$$

This means any vector $\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$ where $x_2 = -x_1$, such as $\begin{bmatrix} 1 \\ -1 \end{bmatrix}$ is an eigenvector with eigenvalue 2. The corresponding eigenspace is given as

$$E_2 = \text{span}\left[\begin{bmatrix} 1 \\ -1 \end{bmatrix}\right]. \quad (4.36)$$

The two eigenspaces E_5 and E_2 in Example 4.5 are one-dimensional as they are each spanned by a single vector. However, in other cases we may have multiple identical eigenvalues (see Definition 4.9) and the eigenspace may have more than one dimension.

Definition 4.11. Let λ_i be an eigenvalue of a square matrix A . Then the *geometric multiplicity* of λ_i is the number of linearly independent eigenvectors associated with λ_i . In other words, it is the dimensionality of the eigenspace spanned by the eigenvectors associated with λ_i .

geometric
multiplicity

Remark. A specific eigenvalue's geometric multiplicity must be at least one because every eigenvalue has at least one associated eigenvector. An eigenvalue's geometric multiplicity cannot exceed its algebraic multiplicity, but it may be lower. \diamond

Example 4.6

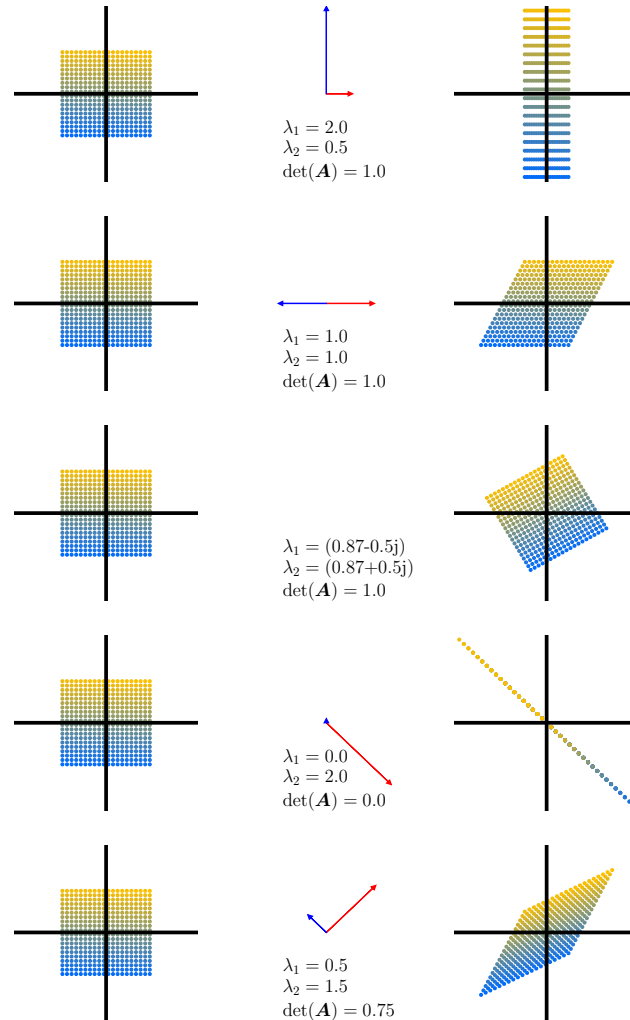
The matrix $A = \begin{bmatrix} 2 & 1 \\ 0 & 2 \end{bmatrix}$ has two repeated eigenvalues $\lambda_1 = \lambda_2 = 2$ and an algebraic multiplicity of 2. The eigenvalue has, however, only one distinct eigenvector $\mathbf{x}_1 = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$ and, thus, geometric multiplicity 1.

Graphical Intuition in Two Dimensions

Let us gain some intuition for determinants, eigenvectors, eigenvalues using different linear mappings. Figure 4.4 depicts five transformation matrices A_1, \dots, A_5 and their impact on a square grid of points, centered at the origin:

- $A_1 = \begin{bmatrix} \frac{1}{2} & 0 \\ 0 & 2 \end{bmatrix}$. The direction of the two eigenvectors correspond to the canonical basis vectors in \mathbb{R}^2 , i.e., to two cardinal axes. The vertical axis is extended by a factor of 2 (eigenvalue $\lambda_1 = 2$), and the horizontal axis

Figure 4.4
Determinants and eigenspaces. Overview of five linear mappings and their associated transformation matrices $\mathbf{A}_i \in \mathbb{R}^{2 \times 2}$ projecting 400 color-coded points $\mathbf{x} \in \mathbb{R}^2$ (left column) onto target points $\mathbf{A}_i \mathbf{x}$ (right column). The central column depicts the *first eigenvector*, stretched by its associated eigenvalue λ_1 , and the *second eigenvector* stretched by its eigenvalue λ_2 . Each row depicts the effect of one of five transformation matrices \mathbf{A}_i with respect to the standard basis.



In geometry, the area-preserving properties of this type of shearing parallel to an axis is also known as Cavalieri's principle of equal areas for parallelograms (Katz, 2004).

is compressed by factor $\frac{1}{2}$ (eigenvalue $\lambda_2 = \frac{1}{2}$). The mapping is area preserving ($\det(\mathbf{A}_1) = 1 = 2 \cdot \frac{1}{2}$).

- $\mathbf{A}_2 = \begin{bmatrix} 1 & \frac{1}{2} \\ 0 & 1 \end{bmatrix}$ corresponds to a shearing mapping, i.e., it shears the points along the horizontal axis to the right if they are on the positive half of the vertical axis, and to the left vice versa. This mapping is area preserving ($\det(\mathbf{A}_2) = 1$). The eigenvalue $\lambda_1 = 1 = \lambda_2$ is repeated and the eigenvectors are collinear (drawn here for emphasis in two opposite directions). This indicates that the mapping acts only along one direction (the horizontal axis).
- $\mathbf{A}_3 = \begin{bmatrix} \cos(\frac{\pi}{6}) & -\sin(\frac{\pi}{6}) \\ \sin(\frac{\pi}{6}) & \cos(\frac{\pi}{6}) \end{bmatrix} = \frac{1}{2} \begin{bmatrix} \sqrt{3} & -1 \\ 1 & \sqrt{3} \end{bmatrix}$ The matrix \mathbf{A}_3 rotates the points by $\frac{\pi}{6}$ rad = 30° anti-clockwise and has only complex eigenvalues,

reflecting that the mapping is a rotation (hence, no eigenvectors are drawn). A rotation has to be volume preserving, and so the determinant is 1. For more details on rotations we refer to Section 3.9.

- $A_4 = \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix}$ represents a mapping in the standard basis that collapses a two-dimensional domain onto one dimension. Therefore, the area of the image is 0. We can see this because one eigenvalue is 0, collapsing the space in direction of the (red) eigenvector corresponding to $\lambda_1 = 0$, while the orthogonal (blue) eigenvector stretches space by a factor $\lambda_2 = 2$.
- $A_5 = \begin{bmatrix} 1 & \frac{1}{2} \\ \frac{1}{2} & 1 \end{bmatrix}$ is a shear-and-stretch mapping that shrinks space space by 75% since $|\det(A_5)| = \frac{3}{4}$. It stretches space along the (blue) eigenvector of λ_2 a factor 1.5 and compresses it along the orthogonal (blue) eigenvector by a 0.5.

Example 4.7 (Eigspectrum of a Biological Neural Network)

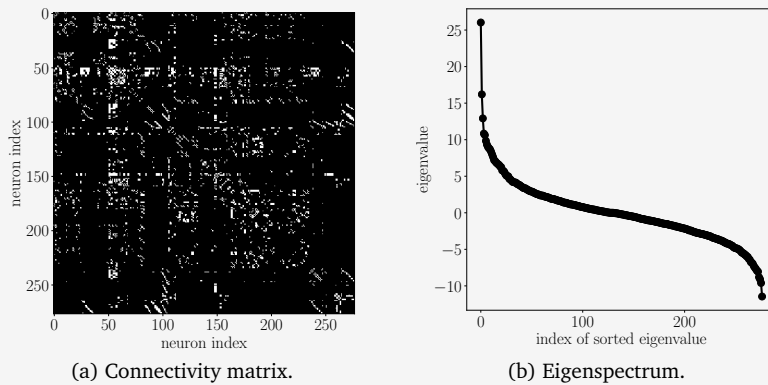


Figure 4.5
C.Elegans neural network (Kaiser and Hilgetag, 2006).
(a) Symmetrized connectivity matrix;
(b) Eigenspectrum.

Methods to analyze and learn from network data are an essential component of machine learning methods. The key to understanding networks is the connectivity between network nodes, especially if two nodes are connected to each other or not. In data science applications, it is often useful to study the matrix that captures this connectivity data.

We build a connectivity/adjacency matrix $A \in \mathbb{R}^{277 \times 277}$ of the complete neural network of the worm *C.Elegans*. Each row/column represents one of the 277 neurons of this worm's brain. The connectivity matrix A has a value of $a_{ij} = 1$ if neuron i talks to neuron j through a synapse, and $a_{ij} = 0$ otherwise. The connectivity matrix is not symmetric, which implies that eigenvalues may not be real valued. Therefore, we compute a symmetrized version of the connectivity matrix as $A_{sym} := A + A^T$. This new matrix A_{sym} is shown in Figure 4.6(a) and has a non-zero value a_{ij}

if and only if two neurons are connected (white pixels), irrespective of the direction of the connection. In Figure 4.6(b), we show the corresponding eigenspectrum of \mathbf{A}_{sym} . The horizontal axis shows the index of the eigenvalues, sorted in descending order. The vertical axis shows the corresponding eigenvalue. The S -like shape of this eigenspectrum is typical for many biological neural networks. The underlying mechanism responsible for this is an area of active neuroscience research.

Theorem 4.12. *The eigenvectors $\mathbf{x}_1, \dots, \mathbf{x}_n$ of a matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ with n distinct eigenvalues $\lambda_1, \dots, \lambda_n$ are linearly independent.*

This theorem states that eigenvectors of a matrix with n distinct eigenvalues form a basis of \mathbb{R}^n .

defective

Definition 4.13. A square matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ is *defective* if it possesses fewer than n linearly independent eigenvectors.

A non-defective matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ does not necessarily require n distinct eigenvalues, but it does require that the eigenvectors form a basis of \mathbb{R}^n . Looking at the eigenspaces of a defective matrix, it follows that the sum of the dimensions of the eigenspaces less than n . Specifically, a defective matrix has at least one eigenvalue λ_i with an algebraic multiplicity $m > 1$ and a geometric multiplicity of less than m .

Remark. A defective matrix cannot have n distinct eigenvalues as distinct eigenvalues have linearly independent eigenvectors (Theorem 4.12). \diamond

Theorem 4.14. *Given a matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ we can always obtain a symmetric, positive semi-definite matrix $\mathbf{S} \in \mathbb{R}^{n \times n}$ by defining*

$$\mathbf{S} := \mathbf{A}^\top \mathbf{A}. \quad (4.37)$$

Remark. If $\text{rk}(\mathbf{A}) = n$ then $\mathbf{S} := \mathbf{A}^\top \mathbf{A}$ is positive definite. \diamond

Understanding why Theorem 4.14 holds is insightful for how we can use symmetrized matrices: Symmetry requires $\mathbf{S} = \mathbf{S}^\top$ and by inserting (4.37) we obtain $\mathbf{S} = \mathbf{A}^\top \mathbf{A} = \mathbf{A}^\top (\mathbf{A}^\top)^\top = (\mathbf{A}^\top \mathbf{A})^\top = \mathbf{S}^\top$. Moreover, positive semi-definiteness (Section 3.2.3) requires that $\mathbf{x}^\top \mathbf{S} \mathbf{x} \geq 0$ and inserting (4.37) we obtain $\mathbf{x}^\top \mathbf{S} \mathbf{x} = \mathbf{x}^\top \mathbf{A}^\top \mathbf{A} \mathbf{x} = (\mathbf{x}^\top \mathbf{A}^\top)(\mathbf{A} \mathbf{x}) = (\mathbf{A} \mathbf{x})^\top (\mathbf{A} \mathbf{x}) \geq 0$, because the dot product computes a sum of squares (which are themselves non-negative).

Theorem 4.15. (Spectral Theorem) *If $\mathbf{A} \in \mathbb{R}^{n \times n}$ is symmetric, there exists an orthonormal basis of the corresponding vector space V consisting of eigenvectors of \mathbf{A} , and each eigenvalue is real.*

A direct implication of the spectral theorem is that the eigendecomposition of a symmetric matrix \mathbf{A} exists (with real eigenvalues), and that we can find an ONB of eigenvectors so that $\mathbf{A} = \mathbf{P} \mathbf{D} \mathbf{P}^\top$, where \mathbf{D} is diagonal and the columns of \mathbf{P} contain the eigenvectors.

Example 4.8

Consider the matrix

$$\mathbf{A} = \begin{bmatrix} 3 & 2 & 2 \\ 2 & 3 & 2 \\ 2 & 2 & 3 \end{bmatrix}. \quad (4.38)$$

The characteristic polynomial of \mathbf{A} is

$$p_{\mathbf{A}}(\lambda) = (\lambda - 1)^2(\lambda - 7), \quad (4.39)$$

so that we obtain the eigenvalues $\lambda_1 = 1$ and $\lambda_2 = 7$, where λ_1 is a repeated eigenvalue. Following our standard procedure for computing eigenvectors, we obtain the eigenspaces

$$E_1 = \text{span} \left[\underbrace{\begin{bmatrix} -1 \\ 1 \\ 0 \end{bmatrix}}_{=: \mathbf{x}_1}, \underbrace{\begin{bmatrix} -1 \\ 0 \\ -1 \end{bmatrix}}_{=: \mathbf{x}_2} \right], \quad E_7 = \text{span} \left[\underbrace{\begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}}_{=: \mathbf{x}_3} \right]. \quad (4.40)$$

We see that \mathbf{x}_3 is orthogonal to both \mathbf{x}_1 and \mathbf{x}_2 . However, since $\mathbf{x}_1^\top \mathbf{x}_2 = 1 \neq 0$ they are not orthogonal. The spectral theorem (Theorem 4.15) states that there exists an orthogonal basis, but the one we have is not orthogonal. However, we can construct one.

To construct such a basis, we exploit the fact that $\mathbf{x}_1, \mathbf{x}_2$ are eigenvectors associated with the same eigenvalue λ . Therefore, for any $\alpha, \beta \in \mathbb{R}$ it holds that

$$\mathbf{A}(\alpha \mathbf{x}_1 + \beta \mathbf{x}_2) = \mathbf{A} \mathbf{x}_1 \alpha + \mathbf{A} \mathbf{x}_2 \beta = \lambda(\alpha \mathbf{x}_1 + \beta \mathbf{x}_2), \quad (4.41)$$

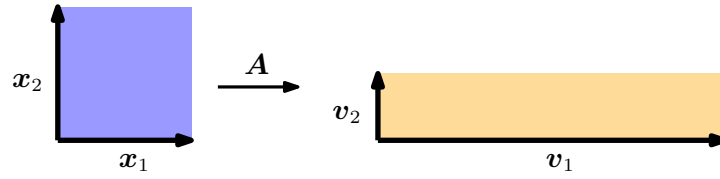
i.e., any linear combination of \mathbf{x}_1 and \mathbf{x}_2 is also an eigenvector of \mathbf{A} associated with λ . The Gram-Schmidt algorithm (Section 3.8.3) is a method for iteratively constructing an orthogonal/orthonormal basis from a set of basis vectors using such linear combinations. Therefore, even if \mathbf{x}_1 and \mathbf{x}_2 are not orthogonal, we can apply the Gram-Schmidt algorithm and find eigenvectors associated with $\lambda_1 = 1$ that are orthogonal to each other (and to \mathbf{x}_3). In our example, we will obtain

$$\mathbf{x}'_1 = \begin{bmatrix} -1 \\ 1 \\ 0 \end{bmatrix}, \quad \mathbf{x}'_2 = \frac{1}{2} \begin{bmatrix} -1 \\ -1 \\ 2 \end{bmatrix}, \quad (4.42)$$

which are orthogonal to each other, orthogonal to \mathbf{x}_3 and eigenvectors of \mathbf{A} associated with $\lambda_1 = 1$.

Before we conclude our considerations of eigenvalues and eigenvectors it is useful to tie these matrix characteristics together with the concepts of the determinant and the trace.

Figure 4.6
Geometric interpretation of eigenvalues. The eigenvectors of A get stretched by the corresponding eigenvalues. The area of the unit square changes by $|\lambda_1 \lambda_2|$, the circumference changes by a factor $(|\lambda_1| + |\lambda_2|)/2$.



Theorem 4.16. *The determinant of a matrix $A \in \mathbb{R}^{n \times n}$ is the product of its eigenvalues, i.e.,*

$$\det(A) = \prod_{i=1}^n \lambda_i, \quad (4.43)$$

where λ_i are (possibly repeated) eigenvalues of A .

Theorem 4.17. *The trace of a matrix $A \in \mathbb{R}^{n \times n}$ is the sum of its eigenvalues, i.e.,*

$$\text{tr}(A) = \sum_{i=1}^n \lambda_i, \quad (4.44)$$

where λ_i are (possibly repeated) eigenvalues of A .

Let us provide a geometric intuition of these two theorems. Consider a matrix $A \in \mathbb{R}^{2 \times 2}$ that possesses two linearly independent eigenvectors x_1, x_2 . For this example, we assume (x_1, x_2) are an ONB of \mathbb{R}^2 so that they are orthogonal and the area of the square they span is 1, see Figure 4.6. From Section 4.1 we know that the determinant computes the change of area of unit square under the transformation A . In this example, we can compute the change of area explicitly: Mapping the eigenvectors using A gives us vectors $v_1 = Ax_1 = \lambda_1 x_1$ and $v_2 = Ax_2 = \lambda_2 x_2$, i.e., the new vectors v_i are scaled versions of the eigenvectors x_i , and the scaling factors are the corresponding eigenvalues λ_i . v_1, v_2 are still orthogonal, and the area of the rectangle they span is $|\lambda_1 \lambda_2|$.

Given that x_1, x_2 (in our example) are orthonormal, we can directly compute the circumference of the unit square as $2(1 + 1)$. Mapping the eigenvectors using A creates a rectangle whose circumference is $2(|\lambda_1| + |\lambda_2|)$. Therefore, the sum of the absolute values of the eigenvalues tells us how the circumference of the unit square changes under the transformation matrix A .

Example 4.9 (Google's PageRank – Webpages as Eigenvectors)

Google uses the eigenvector corresponding to the maximal eigenvalue of a matrix A to determine the rank of a page for search. The idea for the PageRank algorithm, developed at Stanford University by Larry Page and Sergey Brin in 1996, was that the importance of any web page can be approximated by the importance of pages that link to it. For this, they write

down all websites as a huge directed graph that shows which page links to which. PageRank computes the weight (importance) $x_i \geq 0$ of a website a_i by counting the number of pages pointing to a_i . Moreover, PageRank takes into account the importance of the websites that link to a_i . The navigation behavior of a user is then modeled by a transition matrix \mathbf{A} of this graph that tells us with what (click) probability somebody will end up on a different website. The matrix \mathbf{A} has the property that for any initial rank/importance vector \mathbf{x} of a website the sequence $\mathbf{x}, \mathbf{Ax}, \mathbf{A}^2\mathbf{x}, \dots$ converges to a vector \mathbf{x}^* . This vector is called the *PageRank* and satisfies $\mathbf{Ax}^* = \mathbf{x}^*$, i.e., it is an eigenvector (with corresponding eigenvalue 1) of \mathbf{A} . After normalizing \mathbf{x}^* , such that $\|\mathbf{x}^*\| = 1$, we can interpret the entries as probabilities. More details and different perspectives on PageRank can be found in the original technical report (Page et al., 1999).

PageRank

4.3 Cholesky Decomposition

There are many ways to factorize special types of matrices that we encounter often in machine learning. In the positive real numbers, we have the square-root operation that gives us a decomposition of the number into identical components, e.g., $9 = 3 \cdot 3$. For matrices, we need to be careful that we compute a square-root like operation on positive quantities. For symmetric, positive definite matrices (see Section 3.2.3) we can choose from a number of square-root equivalent operations. The *Cholesky decomposition/Cholesky factorization* provides a square-root equivalent operation on symmetric, positive definite matrices that is useful in practice.

Cholesky
decomposition
Cholesky
factorization

Theorem 4.18. *Cholesky Decomposition: A symmetric, positive definite matrix \mathbf{A} can be factorized into a product $\mathbf{A} = \mathbf{LL}^\top$, where \mathbf{L} is a lower-triangular matrix with positive diagonal elements:*

$$\begin{bmatrix} a_{11} & \cdots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{n1} & \cdots & a_{nn} \end{bmatrix} = \begin{bmatrix} l_{11} & \cdots & 0 \\ \vdots & \ddots & \vdots \\ l_{n1} & \cdots & l_{nn} \end{bmatrix} \begin{bmatrix} l_{11} & \cdots & l_{n1} \\ \vdots & \ddots & \vdots \\ 0 & \cdots & l_{nn} \end{bmatrix}. \quad (4.45)$$

\mathbf{L} is called the *Cholesky factor* of \mathbf{A} , and \mathbf{L} is unique.

Cholesky factor

Example 4.10 (Cholesky Factorization)

Consider a symmetric, positive definite matrix $\mathbf{A} \in \mathbb{R}^{3 \times 3}$. We are interested in finding its Cholesky factorization $\mathbf{A} = \mathbf{LL}^\top$, i.e.,

$$\mathbf{A} = \begin{bmatrix} a_{11} & a_{21} & a_{31} \\ a_{21} & a_{22} & a_{32} \\ a_{31} & a_{32} & a_{33} \end{bmatrix} = \mathbf{LL}^\top = \begin{bmatrix} l_{11} & 0 & 0 \\ l_{21} & l_{22} & 0 \\ l_{31} & l_{32} & l_{33} \end{bmatrix} \begin{bmatrix} l_{11} & l_{21} & l_{31} \\ 0 & l_{22} & l_{32} \\ 0 & 0 & l_{33} \end{bmatrix}. \quad (4.46)$$

Multiplying out the right hand side yields

$$\mathbf{A} = \begin{bmatrix} l_{11}^2 & l_{21}l_{11} & l_{31}l_{11} \\ l_{21}l_{11} & l_{21}^2 + l_{22}^2 & l_{31}l_{21} + l_{32}l_{22} \\ l_{31}l_{11} & l_{31}l_{21} + l_{32}l_{22} & l_{31}^2 + l_{32}^2 + l_{33}^2 \end{bmatrix}. \quad (4.47)$$

Comparing the left hand side of (4.46) and the right hand side of (4.47) shows that there is a simple pattern in the diagonal elements l_{ii} :

$$l_{11} = \sqrt{a_{11}}, \quad l_{22} = \sqrt{a_{22} - l_{21}^2}, \quad l_{33} = \sqrt{a_{33} - (l_{31}^2 + l_{32}^2)}. \quad (4.48)$$

Similarly for the elements below the diagonal (l_{ij} , where $i > j$) there is also a repeating pattern:

$$l_{21} = \frac{1}{l_{11}}a_{21}, \quad l_{31} = \frac{1}{l_{11}}a_{31}, \quad l_{32} = \frac{1}{l_{22}}(a_{32} - l_{31}l_{21}). \quad (4.49)$$

Thus, we constructed the Cholesky decomposition for any symmetric, positive definite 3×3 matrix. The key realization is that we can backward calculate what the components l_{ij} for the \mathbf{L} should be, given the values a_{ij} for \mathbf{A} and previously computed values of l_{ij} .

The Cholesky decomposition is an important tool for the numerical computations underlying machine learning. Here, symmetric positive definite matrices require frequent manipulation, e.g., the covariance matrix of a multivariate Gaussian variable (see Section 6.5) is symmetric, positive definite. The Cholesky factorization of this covariance matrix allows us to generate samples from a Gaussian distribution. It also allows us to perform a linear transformation of random variables, which is heavily exploited when computing gradients in deep stochastic models, such as the variational auto-encoder (Jimenez Rezende et al., 2014; Kingma and Welling, 2014). The Cholesky decomposition also allows us to compute determinants very efficiently. Given the Cholesky decomposition $\mathbf{A} = \mathbf{L}\mathbf{L}^\top$, we know that $\det(\mathbf{A}) = \det(\mathbf{L})\det(\mathbf{L}^\top) = \det(\mathbf{L})^2$. Since \mathbf{L} is a triangular matrix, the determinant is simply the product of its diagonal entries so that $\det(\mathbf{A}) = \prod_i l_{ii}^2$. Thus, many numerical software packages use the Cholesky decomposition to make computations more efficient.

4.4 Eigendecomposition and Diagonalization

diagonal matrix

A *diagonal matrix* is a matrix that have value zero on all off diagonal elements, that is they are of the form

$$\mathbf{D} = \begin{bmatrix} c_1 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & c_n \end{bmatrix}. \quad (4.50)$$

They allow fast computation of determinants, powers and inverses. The determinant is the product of its diagonal entries, a matrix power D^k is given by each diagonal element raised to the power k , and the inverse D^{-1} is the reciprocal of its diagonal elements if all of them are non-zero.

In this section, we will discuss how to transform matrices into diagonal form. This is an important application of the basis change we discussed in Section 2.7.2 and eigenvalues from Section 4.2.

Recall that two matrices A, D are similar (Definition 2.21) if there exists an invertible matrix P , such that $D = P^{-1}AP$. More specifically, we will look at matrices A that are similar to diagonal matrices D that contain the eigenvalues of A on the diagonal.

Definition 4.19 (Diagonalizable). A matrix $A \in \mathbb{R}^{n \times n}$ is *diagonalizable* if it is similar to a diagonal matrix, i.e., if there exists an invertible matrix $P \in \mathbb{R}^{n \times n}$ such that $D = P^{-1}AP$. diagonalizable

In the following, we will see that diagonalizing a matrix $A \in \mathbb{R}^{n \times n}$ is a way of expressing the same linear mapping but in another basis (see Section 2.6.1), which will turn out to be a basis that consists of the eigenvectors of A .

Let $A \in \mathbb{R}^{n \times n}$, let $\lambda_1, \dots, \lambda_n$ be a set of scalars, and let p_1, \dots, p_n be a set of vectors in \mathbb{R}^n . We define $P := [p_1, \dots, p_n]$ and let $D \in \mathbb{R}^{n \times n}$ be a diagonal matrix with diagonal entries $\lambda_1, \dots, \lambda_n$. Then we can show that

$$AP = PD \quad (4.51)$$

if and only if $\lambda_1, \dots, \lambda_n$ are the eigenvalues of A and p_1, \dots, p_n are corresponding eigenvectors of A .

We can see that this statement holds because

$$AP = A[p_1, \dots, p_n] = [Ap_1, \dots, Ap_n], \quad (4.52)$$

$$PD = [p_1, \dots, p_n] \begin{bmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & & \lambda_n \end{bmatrix} = [\lambda_1 p_1, \dots, \lambda_n p_n]. \quad (4.53)$$

Thus, (4.51) implies that

$$Ap_1 = \lambda_1 p_1 \quad (4.54)$$

$$\vdots$$

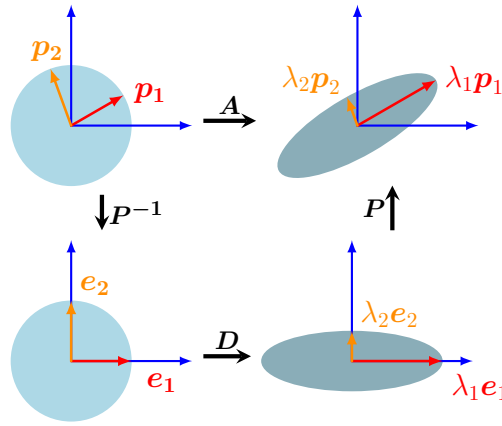
$$Ap_n = \lambda_n p_n. \quad (4.55)$$

Therefore, the columns of P must be eigenvectors of A .

Our definition of diagonalization requires that $P \in \mathbb{R}^{n \times n}$ is invertible, i.e., P has full rank (Theorem 4.3). This requires us to have n linearly independent eigenvectors p_1, \dots, p_n , i.e., the p_i form a basis of \mathbb{R}^n .

Theorem 4.20. (Eigendecomposition/diagonalization) A square matrix $A \in$ diagonalization

Figure 4.7 Intuition behind the eigendecomposition as sequential transformations. Top-left to bottom-left: P^{-1} performs a basis change (here drawn in \mathbb{R}^2 and depicted as a rotation-like operation), mapping the eigenvectors into the standard basis. Bottom-left to bottom-right: D performs a scaling along the remapped orthogonal eigenvectors, depicted here by a circle being stretched to an ellipse. Bottom-right to top-right: P undoes the basis change (depicted as a reverse rotation) and restores the original coordinate frame.



$\mathbb{R}^{n \times n}$ can be factored into

$$A = PDP^{-1}, \quad (4.56)$$

where $P \in \mathbb{R}^{n \times n}$ and D is a diagonal matrix whose diagonal entries are the eigenvalues of A , if and only if the eigenvectors of A form a basis of \mathbb{R}^n .

Theorem 4.20 implies that only non-defective matrices can be diagonalized and that the columns of P are the n eigenvectors of A . For symmetric matrices we can obtain even stronger outcomes for the eigenvalue decomposition.

Theorem 4.21. A symmetric matrix $S \in \mathbb{R}^{n \times n}$ can always be diagonalized.

Theorem 4.21 follows directly from the spectral theorem 4.15. Moreover, the spectral theorem states that we can find an ONB of eigenvectors of \mathbb{R}^n . This makes P an orthogonal matrix so that $D = P^T A P$.

Geometric Intuition for the Eigendecomposition

We can interpret the eigendecomposition of a matrix as follows (see also Figure 4.7): Let A be the transformation matrix of a linear mapping with respect to the standard basis. P^{-1} performs a basis change from the standard basis into the eigenbasis. This identifies the eigenvectors p_i (red and orange arrows in Figure 4.7) onto the standard basis vectors e_i . Then, the diagonal D scales the vectors along these axes by the eigenvalues λ_i . Finally, P transforms these scaled vectors back into the standard/canonical coordinates yielding $\lambda_i p_i$.

Example 4.11 (Eigendecomposition)

Let us compute the eigendecomposition of $A = \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix}$.

The Jordan normal form of a matrix offers a decomposition that works for defective matrices (Lang, 1987) but is beyond the scope of this book.

Step 1: Compute eigenvalues and eigenvectors. The characteristic polynomial of \mathbf{A} is

$$\det(\mathbf{A} - \lambda \mathbf{I}) = \det \left(\begin{bmatrix} 2 - \lambda & 1 \\ 1 & 2 - \lambda \end{bmatrix} \right) \quad (4.57a)$$

$$= (2 - \lambda)^2 - 1 = \lambda^2 - 4\lambda + 3 = (\lambda - 3)(\lambda - 1). \quad (4.57b)$$

Therefore, the eigenvalues of \mathbf{A} are $\lambda_1 = 1$ and $\lambda_2 = 3$ (the roots of the characteristic polynomial), and the associated (normalized) eigenvectors are obtained via

$$\begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix} \mathbf{p}_1 = 1\mathbf{p}_1, \quad \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix} \mathbf{p}_2 = 3\mathbf{p}_2. \quad (4.58)$$

This yields

$$\mathbf{p}_1 = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ -1 \end{bmatrix}, \quad \mathbf{p}_2 = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ 1 \end{bmatrix}. \quad (4.59)$$

Step 2: Check for existence The eigenvectors $\mathbf{p}_1, \mathbf{p}_2$ form a basis of \mathbb{R}^2 . Therefore, \mathbf{A} can be diagonalized.

Step 3: Construct the matrix \mathbf{P} to diagonalize \mathbf{A} We collect the eigenvectors of \mathbf{A} in \mathbf{P} so that

$$\mathbf{P} = [\mathbf{p}_1, \mathbf{p}_2] = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 1 \\ -1 & 1 \end{bmatrix}. \quad (4.60)$$

We then obtain

$$\mathbf{P}^{-1} \mathbf{A} \mathbf{P} = \begin{bmatrix} 1 & 0 \\ 0 & 3 \end{bmatrix} = \mathbf{D}. \quad (4.61)$$

Equivalently, we get (exploiting that $\mathbf{P}^{-1} = \mathbf{P}^\top$ since the eigenvectors \mathbf{p}_1 and \mathbf{p}_2 in this example form an ONB)

$$\underbrace{\begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix}}_{\mathbf{A}} = \underbrace{\frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 1 \\ -1 & 1 \end{bmatrix}}_{\mathbf{P}} \underbrace{\begin{bmatrix} 1 & 0 \\ 0 & 3 \end{bmatrix}}_{\mathbf{D}} \underbrace{\frac{1}{\sqrt{2}} \begin{bmatrix} 1 & -1 \\ 1 & 1 \end{bmatrix}}_{\mathbf{P}^\top}. \quad (4.62)$$

- Diagonal matrices \mathbf{D} can efficiently be raised to a power. Therefore, we can find a matrix power for a matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ via the eigenvalue decomposition (if it exists) so that

$$\mathbf{A}^k = (\mathbf{P} \mathbf{D} \mathbf{P}^{-1})^k = \mathbf{P} \mathbf{D}^k \mathbf{P}^{-1}. \quad (4.63)$$

Computing \mathbf{D}^k is efficient because we apply this operation individually to any diagonal element.

- Assuming the eigendecomposition $\mathbf{A} = \mathbf{P} \mathbf{D} \mathbf{P}^{-1}$ exists. Then,

$$\det(\mathbf{A}) = \det(\mathbf{P} \mathbf{D} \mathbf{P}^{-1}) = \det(\mathbf{P}) \det(\mathbf{D}) \det(\mathbf{P}^{-1}) \quad (4.64a)$$

$$= \det(\mathbf{D}) = \prod_i d_{ii} \quad (4.64b)$$

allows for an efficient computation of the determinant of \mathbf{A} .

The eigenvalue decomposition requires square matrices. It would be useful to perform a decomposition on general matrices. In the next section, we introduce a more general matrix decomposition technique, the singular value decomposition.

4.5 Singular Value Decomposition

The singular value decomposition (SVD) of a matrix is a central matrix decomposition method in linear algebra. It has been referred to as the “fundamental theorem of linear algebra” (Strang, 1993) because it can be applied to all matrices, not only to square matrices, and it always exists. Moreover, as we will explore in the following, the SVD of a matrix \mathbf{A} , which represents a linear mapping $\Phi : V \rightarrow W$, quantifies the change between the underlying geometry of these two vector spaces. We recommend the work by Kalman (1996) and Roy and Banerjee (2014) for a deeper overview of the mathematics of the SVD.

Theorem 4.22 (SVD Theorem). *Let $\mathbf{A}^{m \times n}$ be a rectangular matrix of rank $r \in [0, \min(m, n)]$. The singular value decomposition (SVD) of \mathbf{A} is a decomposition of the form*

$$\begin{matrix} n \\ \boxed{\mathbf{A}} \\ m \end{matrix} = \begin{matrix} m \\ \boxed{\mathbf{U}} \\ m \end{matrix} \begin{matrix} n \\ \boxed{\Sigma} \\ m \end{matrix} \begin{matrix} n \\ \boxed{\mathbf{V}^\top} \\ n \end{matrix} \quad (4.65)$$

with an orthogonal matrix $\mathbf{U} \in \mathbb{R}^{m \times m}$ with column vectors \mathbf{u}_i , $i = 1, \dots, m$, and an orthogonal matrix $\mathbf{V} \in \mathbb{R}^{n \times n}$ with column vectors \mathbf{v}_j , $j = 1, \dots, n$. Moreover, Σ is an $m \times n$ matrix with $\Sigma_{ii} = \sigma_i \geq 0$ and $\Sigma_{ij} = 0$, $i \neq j$.

The diagonal entries σ_i , $i = 1, \dots, r$, of Σ are called the *singular values*, \mathbf{u}_i are called the *left-singular vectors* and \mathbf{v}_j are called the *right-singular vectors*. By convention the singular values are ordered, i.e., $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r \geq 0$.

The *singular value matrix* Σ is unique, but it requires some attention. Observe that the $\Sigma \in \mathbb{R}^{m \times n}$ is rectangular. In particular, Σ is of the same size as \mathbf{A} . This means that Σ has a diagonal submatrix that contains the singular values and needs additional zero padding. Specifically, if $m > n$ then the matrix Σ has diagonal structure up to row n and then consists of

singular value
decomposition
SVD

singular values
left-singular vectors
right-singular
vectors

singular value
matrix

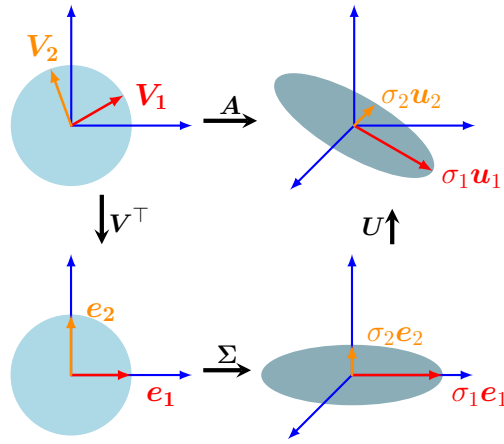


Figure 4.8 Intuition behind the SVD of a matrix $A \in \mathbb{R}^{3 \times 2}$ as sequential transformations. Top-left to bottom-left: V^\top performs a basis change in \mathbb{R}^2 . Bottom-left to bottom right: Σ scales and maps from \mathbb{R}^2 to \mathbb{R}^3 . The ellipse in the bottom-right lives in \mathbb{R}^3 . The third dimension is orthogonal to the surface of the elliptical disk. Bottom-right to top-right: U performs a basis change within \mathbb{R}^3 .

$\mathbf{0}^\top$ row vectors from $n + 1$ to m below so that

$$\Sigma = \begin{bmatrix} \sigma_1 & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & \sigma_n \\ 0 & \dots & 0 \\ \vdots & & \vdots \\ 0 & \dots & 0 \end{bmatrix}. \quad (4.66)$$

If $m < n$ the matrix Σ has a diagonal structure up to column m and columns that consist of $\mathbf{0}$ from $m + 1$ to n :

$$\Sigma = \begin{bmatrix} \sigma_1 & 0 & 0 & 0 & \dots & 0 \\ 0 & \ddots & 0 & 0 & & 0 \\ 0 & 0 & \sigma_n & 0 & \dots & 0 \end{bmatrix}. \quad (4.67)$$

Remark. The SVD exists for any matrix $A \in \mathbb{R}^{m \times n}$. \diamond

4.5.1 Geometric Intuitions for the SVD

The SVD offers geometric intuitions to describe a transformation matrix A . In the following, we will discuss the SVD as sequential linear transformations performed on the bases. In example 4.12, we will then apply transformation matrices of the SVD to a set of vectors in \mathbb{R}^2 , which allows us to visualize the effect of each transformation more clearly.

The SVD of a matrix can be interpreted as a decomposition of a corresponding linear mapping (recall Section 2.7.1) $\Phi : \mathbb{R}^n \rightarrow \mathbb{R}^m$ into three operations (see Figure 4.8). The SVD intuition follows superficially a similar structure to our eigendecomposition intuition, see Figure 4.7: Broadly speaking, the SVD performs a basis change via V^\top followed by a scaling and augmentation (or reduction) in dimensionality via the singular

It is useful to revise basis changes (Section 2.7.2), orthogonal matrices (Definition 3.8) and orthonormal bases (Section 3.5).

value matrix Σ . Finally, it performs a second basis change via U . The SVD entails a number of important details and caveats, which is why we will review our intuition in more detail.

Assume we are given a transformation matrix of a linear mapping $\Phi : \mathbb{R}^n \rightarrow \mathbb{R}^m$ with respect to the standard bases B and C of \mathbb{R}^n and \mathbb{R}^m , respectively. Moreover, assume a second basis \tilde{B} of \mathbb{R}^n and \tilde{C} of \mathbb{R}^m . Then

- 1 The matrix V performs a basis change in the domain \mathbb{R}^n from \tilde{B} (represented by the red and orange vectors v_1 and v_2 in the top-left of Figure 4.8) to the standard basis B . $V^\top = V^{-1}$ performs a basis change from B to \tilde{B} . The red and orange vectors are now aligned with the canonical basis in the bottom left of Figure 4.8.
- 2 Having changed the coordinate system to \tilde{B} , Σ scales the new coordinates by the singular values σ_i (and adds or deletes dimensions), i.e., Σ is the transformation matrix of Φ with respect to \tilde{B} and \tilde{C} , represented by the red and orange vectors being stretched and lying in the e_1 - e_2 plane, which is now embedded in a third dimension in the bottom right of Figure 4.8.
- 3 U performs a basis change in the codomain \mathbb{R}^m from \tilde{C} into the canonical basis of \mathbb{R}^m , represented by a rotation of red and orange vectors out of the e_1 - e_2 plane in the bottom right of Figure 4.8.

The SVD expresses a change of basis in both the domain and codomain. This is in contrast with the eigendecomposition that operates within the same vector space, where the same basis change is applied and then undone. What makes the SVD special is that these two different bases are simultaneously linked by the singular value matrix Σ .

Example 4.12 (Vectors and the SVD)

Consider a mapping of a square grid of vectors $\mathcal{X} \in \mathbb{R}^2$ which fit in a box of size 2×2 centered at the origin. Using the standard basis we map these vectors using

$$A = \begin{bmatrix} 1 & -0.8 \\ 0 & 1 \\ 1 & 0 \end{bmatrix} = U \Sigma V^\top \quad (4.68a)$$

$$= \begin{bmatrix} -0.79 & 0 & -0.62 \\ 0.38 & -0.78 & -0.49 \\ -0.48 & -0.62 & 0.62 \end{bmatrix} \begin{bmatrix} 1.62 & 0 \\ 0 & 1.0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} -0.78 & 0.62 \\ -0.62 & -0.78 \end{bmatrix}. \quad (4.68b)$$

We start with a set of vectors \mathcal{X} (colored dots, see top-left panel of Figure 4.9) arranged in a grid. We then apply $V^\top \in \mathbb{R}^{2 \times 2}$, which rotates \mathcal{X} . The rotated vectors are shown in the bottom-left panel of Figure 4.9. We now map these vectors using the singular value matrix Σ to the codomain

\mathbb{R}^3 (see bottom right panel in Figure 4.9). Note that all vectors lie in the x_1 - x_2 plane. The third coordinate is always 0. The vectors in the x_1 - x_2 plane has been stretched by the singular values.

The direct mapping of the vectors \mathcal{X} by \mathbf{A} to the codomain \mathbb{R}^3 equals the transformation of \mathcal{X} by $\mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top$, where \mathbf{U} performs a rotation within the codomain \mathbb{R}^3 so that the mapped vectors are no longer restricted to the x_1 - x_2 plane; they still are on a plane as shown in the top-right panel of Figure 4.9.

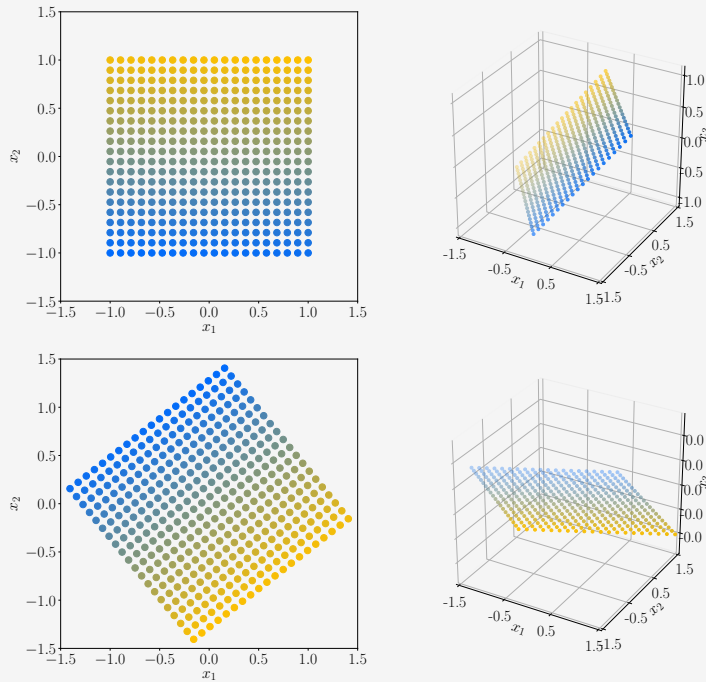


Figure 4.9 SVD and mapping of vectors (represented by discs). The panels follow the same anti-clockwise structure of Figure 4.8.

4.5.2 Construction of the SVD

We will next discuss why the SVD exists and show how to compute it in detail. The SVD of a general matrix shares some similarities with the eigendecomposition of a square matrix.

Remark. Compare the eigendecomposition of an SPD matrix

$$\mathbf{S} = \mathbf{S}^\top = \mathbf{P}\mathbf{D}\mathbf{P}^\top \quad (4.69)$$

with the corresponding SVD

$$\mathbf{S} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top. \quad (4.70)$$

If we set

$$U = P = V, \quad D = \Sigma \quad (4.71)$$

we see that the SVD of SPD matrices is their eigendecomposition. \diamond

In the following, we will explore why Theorem 4.22 holds and how the SVD is constructed. Computing the SVD of $A \in \mathbb{R}^{m \times n}$ is equivalent to finding two sets of orthonormal bases $U = (\mathbf{u}_1, \dots, \mathbf{u}_m)$ and $V = (\mathbf{v}_1, \dots, \mathbf{v}_n)$ of the codomain \mathbb{R}^m and the domain \mathbb{R}^n , respectively. From these ordered bases we will construct the matrices U and V .

Our plan is to start with constructing the orthonormal set of right-singular vectors $\mathbf{v}_1, \dots, \mathbf{v}_n \in \mathbb{R}^n$. We then construct the orthonormal set of left-singular vectors $\mathbf{u}_1, \dots, \mathbf{u}_m \in \mathbb{R}^m$. Thereafter, we will link the two and require that the orthogonality of the \mathbf{v}_i is preserved under the transformation of A . This is important because we know that the images $A\mathbf{v}_i$ form a set of orthogonal vectors. We will then normalize these images by scalar factors, which will turn out to be the singular values.

Let us begin with constructing the right-singular vectors. The spectral theorem (Theorem 4.15) tells us that a symmetric matrix possesses an ONB of eigenvectors, which also means it can be diagonalized. Moreover, from Theorem 4.14 we can always construct a symmetric, positive semi-definite matrix $A^\top A \in \mathbb{R}^{n \times n}$ from any rectangular matrix $A \in \mathbb{R}^{m \times n}$. Thus, we can always diagonalize $A^\top A$ and obtain

$$A^\top A = PDP^\top = P \begin{bmatrix} \lambda_1 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \lambda_n \end{bmatrix} P^\top, \quad (4.72)$$

where P is an orthogonal matrix, which is composed of the orthonormal eigenbasis. The $\lambda_i \geq 0$ are the eigenvalues of $A^\top A$. Let us assume the SVD of A exists and inject (4.65) into (4.72). This yields

$$A^\top A = (U\Sigma V^\top)^\top (U\Sigma V^\top) = V\Sigma^\top U^\top U\Sigma V^\top, \quad (4.73)$$

where U, V are orthogonal matrices. Therefore, with $U^\top U = I$ we obtain

$$A^\top A = V\Sigma^\top \Sigma V^\top = V \begin{bmatrix} \sigma_1^2 & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & \sigma_n^2 \end{bmatrix} V^\top. \quad (4.74)$$

Comparing now (4.72) and (4.74) we identify

$$V^\top = P^\top, \quad (4.75)$$

$$\sigma_i^2 = \lambda_i. \quad (4.76)$$

Therefore, the eigenvectors of $A^\top A$ that compose P are the right-singular vectors V of A (see (4.75)). The eigenvalues of $A^\top A$ are the squared singular values of Σ (see (4.76)).

To obtain the left-singular vectors \mathbf{U} we follow a similar procedure. We start by computing the SVD of the symmetric matrix $\mathbf{A}\mathbf{A}^\top \in \mathbb{R}^{m \times m}$ (instead of the above $\mathbf{A}^\top \mathbf{A} \in \mathbb{R}^{n \times n}$). The SVD of \mathbf{A} yields

$$\mathbf{A}\mathbf{A}^\top = (\mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top)(\mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top)^\top = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top \mathbf{V}\mathbf{\Sigma}^\top \mathbf{U}^\top \quad (4.77a)$$

$$= \mathbf{U} \begin{bmatrix} \sigma_1^2 & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & \sigma_m^2 \end{bmatrix} \mathbf{U}^\top. \quad (4.77b)$$

The spectral theorem tells us that $\mathbf{A}\mathbf{A}^\top = \mathbf{S}\mathbf{D}\mathbf{S}^\top$ can be diagonalized and we can find an ONB of eigenvectors of $\mathbf{A}\mathbf{A}^\top$, which are collected in \mathbf{S} . The orthonormal eigenvectors of $\mathbf{A}\mathbf{A}^\top$ are the left-singular vectors \mathbf{U} and form an orthonormal basis set in the codomain of the SVD.

This leaves the question of the structure of the matrix $\mathbf{\Sigma}$. Since $\mathbf{A}\mathbf{A}^\top$ and $\mathbf{A}^\top \mathbf{A}$ have the same non-zero eigenvalues (see page 109) the non-zero entries of the $\mathbf{\Sigma}$ matrices in the SVD for both cases have to be the same.

The last step is to link up all the parts we touched upon so far. We have an orthonormal set of right-singular vectors in \mathbf{V} . To finish the construction of the SVD we connect them with the orthonormal vectors \mathbf{U} . To reach this goal we use the fact the images of the \mathbf{v}_i under \mathbf{A} have to be orthogonal, too. We can show this by using the results from Section 3.4. We require that the inner product between $\mathbf{A}\mathbf{v}_i$ and $\mathbf{A}\mathbf{v}_j$ must be 0 for $i \neq j$. For any two orthogonal eigenvectors $\mathbf{v}_i, \mathbf{v}_j$, $i \neq j$ it holds that

$$(\mathbf{A}\mathbf{v}_i)^\top (\mathbf{A}\mathbf{v}_j) = \mathbf{v}_i^\top (\mathbf{A}^\top \mathbf{A}) \mathbf{v}_j = \mathbf{v}_i^\top (\lambda_j \mathbf{v}_j) = \lambda_j \mathbf{v}_i^\top \mathbf{v}_j = 0. \quad (4.78)$$

For the case $m \geq r$ it holds that $\{\mathbf{A}\mathbf{v}_1, \dots, \mathbf{A}\mathbf{v}_r\}$ are a basis of an r -dimensional subspace of \mathbb{R}^m .

To complete the SVD construction we need left-singular vectors that are *orthonormal*: we normalize the images of the right-singular vectors $\mathbf{A}\mathbf{v}_i$ and obtain

$$\mathbf{u}_i := \frac{\mathbf{A}\mathbf{v}_i}{\|\mathbf{A}\mathbf{v}_i\|} = \frac{1}{\sqrt{\lambda_i}} \mathbf{A}\mathbf{v}_i = \frac{1}{\sigma_i} \mathbf{A}\mathbf{v}_i, \quad (4.79)$$

where the last equality was obtained from (4.76) and (4.77b) showing us that the eigenvalues of $\mathbf{A}\mathbf{A}^\top$ are such that $\sigma_i^2 = \lambda_i$.

Therefore, the eigenvectors of $\mathbf{A}^\top \mathbf{A}$, which we know are the right-singular vectors \mathbf{v}_i , and their normalized images under \mathbf{A} , the left-singular vectors \mathbf{u}_i , form two self-consistent ONBs that are connected through the singular value matrix $\mathbf{\Sigma}$.

Let us rearrange (4.79) to obtain the *singular value equation*

$$\mathbf{A}\mathbf{v}_i = \sigma_i \mathbf{u}_i, \quad i = 1, \dots, r. \quad (4.80)$$

singular value
equation

This equation closely resembles the eigenvalue equation (4.26), but the vectors on the left and the right-hand sides are not the same.

For $n > m$ (4.80) holds only for $i \leq m$ and (4.80) says nothing about the \mathbf{u}_i for $i > m$. However, we know by construction that they are orthonormal. Conversely, for $m > n$, (4.80) holds only for $i \leq n$. For $i > n$ we have $\mathbf{A}\mathbf{v}_i = \mathbf{0}$ and we still know that the \mathbf{v}_i form an orthonormal set. This means that the SVD also supplies an orthonormal basis of the kernel (null space) of \mathbf{A} , the set of vectors \mathbf{x} with $\mathbf{A}\mathbf{x} = \mathbf{0}$ (see Section 2.7.3).

Moreover, concatenating the \mathbf{v}_i as the columns of \mathbf{V} and the \mathbf{u}_i as the columns of \mathbf{U} yields

$$\mathbf{A}\mathbf{V} = \mathbf{U}\mathbf{\Sigma}, \quad (4.81)$$

where $\mathbf{\Sigma}$ has the same dimensions as \mathbf{A} and a diagonal structure for rows $1, \dots, r$. Hence, right-multiplying with \mathbf{V}^\top yields $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top$, which is the SVD of \mathbf{A} .

Example 4.13 (Computing the SVD)

Let us find the singular value decomposition of

$$\mathbf{A} = \begin{bmatrix} 1 & 0 & 1 \\ -2 & 1 & 0 \end{bmatrix}. \quad (4.82)$$

The SVD requires us to compute the right-singular vectors \mathbf{v}_j , the singular values σ_k and the left-singular values \mathbf{u}_i .

Step 1: Right-singular vectors as the eigenbasis of $\mathbf{A}^\top \mathbf{A}$.

We start by computing

$$\mathbf{A}^\top \mathbf{A} = \begin{bmatrix} 1 & -2 \\ 0 & 1 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} 1 & 0 & 1 \\ -2 & 1 & 0 \end{bmatrix} = \begin{bmatrix} 5 & -2 & 1 \\ -2 & 1 & 0 \\ 1 & 0 & 1 \end{bmatrix}. \quad (4.83)$$

We compute the singular values and right-singular vectors \mathbf{v}_j through the eigenvalue decomposition of $\mathbf{A}^\top \mathbf{A}$, which is given as

$$\mathbf{A}^\top \mathbf{A} = \begin{bmatrix} \frac{5}{\sqrt{30}} & 0 & \frac{-1}{\sqrt{6}} \\ \frac{-2}{\sqrt{30}} & \frac{1}{\sqrt{5}} & \frac{-2}{\sqrt{6}} \\ \frac{1}{\sqrt{30}} & \frac{2}{\sqrt{5}} & \frac{1}{\sqrt{6}} \end{bmatrix} \begin{bmatrix} 6 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} \frac{5}{\sqrt{30}} & \frac{-2}{\sqrt{30}} & \frac{1}{\sqrt{30}} \\ 0 & \frac{1}{\sqrt{5}} & \frac{2}{\sqrt{5}} \\ \frac{-1}{\sqrt{6}} & \frac{-2}{\sqrt{6}} & \frac{1}{\sqrt{6}} \end{bmatrix} = \mathbf{P}\mathbf{D}\mathbf{P}^\top, \quad (4.84)$$

and we obtain the right-singular vectors as the columns of \mathbf{P} so that

$$\mathbf{V} = \mathbf{P} = \begin{bmatrix} \frac{5}{\sqrt{30}} & 0 & \frac{-1}{\sqrt{6}} \\ \frac{-2}{\sqrt{30}} & \frac{1}{\sqrt{5}} & \frac{-2}{\sqrt{6}} \\ \frac{1}{\sqrt{30}} & \frac{2}{\sqrt{5}} & \frac{1}{\sqrt{6}} \end{bmatrix}. \quad (4.85)$$

Step 2: Singular-value matrix.

As the singular values σ_i are the square-roots of the eigenvalues of $\mathbf{A}^\top \mathbf{A}$ we obtain them straight from \mathbf{D} . Since $\text{rk}(\mathbf{A}) = 2$ there are only two non-zero singular values: $\sigma_1 = \sqrt{6}$ and $\sigma_2 = 1$. The singular value matrix must

be the same size as \mathbf{A} , and we obtain

$$\Sigma = \begin{bmatrix} \sqrt{6} & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix}. \quad (4.86)$$

Step 3: Left-singular vectors as the normalized image of the right-singular vectors.

We find the left-singular vectors by computing the image of the right-singular vectors under \mathbf{A} and normalizing them by dividing them by their corresponding singular value. We obtain

$$\mathbf{u}_1 = \frac{1}{\sigma_1} \mathbf{A} \mathbf{v}_1 = \frac{1}{\sqrt{6}} \begin{bmatrix} 1 & 0 & 1 \\ -2 & 1 & 0 \end{bmatrix} \begin{bmatrix} \frac{5}{\sqrt{30}} \\ \frac{-2}{\sqrt{30}} \\ \frac{1}{\sqrt{30}} \end{bmatrix} = \begin{bmatrix} \frac{1}{\sqrt{5}} \\ -\frac{2}{\sqrt{5}} \end{bmatrix}, \quad (4.87)$$

$$\mathbf{u}_2 = \frac{1}{\sigma_2} \mathbf{A} \mathbf{v}_2 = \frac{1}{1} \begin{bmatrix} 1 & 0 & 1 \\ -2 & 1 & 0 \end{bmatrix} \begin{bmatrix} 0 \\ \frac{1}{\sqrt{5}} \\ \frac{2}{\sqrt{5}} \end{bmatrix} = \begin{bmatrix} \frac{2}{\sqrt{5}} \\ \frac{1}{\sqrt{5}} \end{bmatrix}, \quad (4.88)$$

$$\mathbf{U} = [\mathbf{u}_1, \mathbf{u}_2] = \frac{1}{\sqrt{5}} \begin{bmatrix} 1 & 2 \\ -2 & 1 \end{bmatrix}. \quad (4.89)$$

Note that on a computer the approach illustrated here has poor numerical behavior, and the SVD of \mathbf{A} is normally computed without resorting to the eigenvalue decomposition of $\mathbf{A}^\top \mathbf{A}$.

4.5.3 Eigenvalue Decomposition vs Singular Value Decomposition

Let us consider the eigendecomposition $\mathbf{A} = \mathbf{P} \mathbf{D} \mathbf{P}^{-1}$ and the SVD $\mathbf{A} = \mathbf{U} \Sigma \mathbf{V}^\top$ and review the core elements of the past sections.

- The SVD always exists for any matrix $\mathbb{R}^{m \times n}$. The eigendecomposition is only defined for square matrices $\mathbb{R}^{n \times n}$ and only exists if we can find a basis of eigenvectors of \mathbb{R}^n .
- The vectors in the eigendecomposition matrix \mathbf{P} are not necessarily orthogonal, i.e., the change of basis is not a simple rotation and scaling. On the other hand, the vectors in the matrices \mathbf{U} and \mathbf{V} in the SVD are orthonormal, so they do represent rotations.
- Both the eigendecomposition and the SVD are compositions of three linear mappings:
 - 1 Change of basis in the domain
 - 2 Independent scaling of each new basis vector and mapping from domain to codomain
 - 3 Change of basis in the codomain

A key difference between the eigendecomposition and the SVD is that

Figure 4.10 Movie ratings of three people for four movies and its SVD decomposition.

$$\begin{array}{c}
 \text{Star Wars} \\
 \text{Blade Runner} \\
 \text{Amelie} \\
 \text{Delicatessen}
 \end{array}
 \begin{array}{c}
 \text{Ali} \\
 \text{Beatrix} \\
 \text{Chandra}
 \end{array}
 \begin{bmatrix}
 5 & 4 & 1 \\
 5 & 5 & 0 \\
 0 & 0 & 5 \\
 1 & 0 & 4
 \end{bmatrix}
 =
 \begin{bmatrix}
 -0.6710 & 0.0236 & 0.4647 & -0.5774 \\
 -0.7197 & 0.2054 & -0.4759 & 0.4619 \\
 -0.0939 & -0.7705 & -0.5268 & -0.3464 \\
 -0.1515 & -0.6030 & 0.5293 & -0.5774
 \end{bmatrix}
 \begin{bmatrix}
 9.6438 & 0 & 0 \\
 0 & 6.3639 & 0 \\
 0 & 0 & 0.7056 \\
 0 & 0 & 0
 \end{bmatrix}
 \begin{bmatrix}
 -0.7367 & -0.6515 & -0.1811 \\
 0.0852 & 0.1762 & -0.9807 \\
 0.6708 & -0.7379 & -0.0743
 \end{bmatrix}$$

in the SVD, domain and codomain can be vector spaces of different dimensions.

- In the SVD, the left and right singular vector matrices U and V are generally not inverse of each other (they perform basis changes in different vector spaces). In the eigendecomposition, the basis change matrices P and P^{-1} are inverses of each other.
- In the SVD, the entries in the diagonal matrix Σ are all real and non-negative, which is not generally true for the diagonal matrix in the eigendecomposition.
- The SVD and the eigendecomposition are closely related through their projections
 - The left-singular vectors of A are eigenvectors of AA^T
 - The right-singular vectors of A are eigenvectors of $A^T A$.
 - The non-zero singular values of A are the square-roots of the non-zero eigenvalues of AA^T and are equal to the non-zero eigenvalues of $A^T A$.
- For symmetric matrices $A \in \mathbb{R}^{n \times n}$ the eigenvalue decomposition and the SVD are one and the same, which follows from the spectral theorem 4.15.

Example 4.14 (Finding Structure in Movie Ratings and Consumers)

Let us add a practical interpretation of the SVD by analyzing data on people and their preferred movies. Consider 3 viewers (Ali, Beatrix, Chandra) rating four different movies (Star Wars, Blade Runner, Amelie, Delicatessen). Their ratings are values between 0 (worst) and 5 (best) and encoded in a data matrix $A \in \mathbb{R}^{4 \times 3}$ as shown in Figure 4.10. Each row represents a movie and each column a user. Thus, the column vectors of movie ratings, one for each viewer, are \mathbf{x}_{Ali} , $\mathbf{x}_{\text{Beatrix}}$, $\mathbf{x}_{\text{Chandra}}$.

Factoring \mathbf{A} using the SVD offers us a way to capture the relationships of how people rate movies, and especially if there is a structure linking which people like which movies. Applying the SVD to our data matrix \mathbf{A} makes a number of assumptions:

- 1 All viewers rate movies consistently using the same linear mapping.
- 2 There are no errors or noise in the ratings.
- 3 We interpret the left-singular vectors \mathbf{u}_i as stereotypical movies and the right-singular vectors \mathbf{v}_j as stereotypical viewers.

We then make the assumption that any viewer's specific movie preferences can be expressed as a linear combination of the \mathbf{v}_j . Similarly, any movie's like-ability can be expressed as a linear combination of the \mathbf{u}_i . Therefore, a vector in the domain of the SVD can be interpreted as a viewer in the "space" of stereotypical viewers, and a vector in the co-domain of the SVD correspondingly as a movie in the "space" of stereotypical movies. Let us inspect the SVD of our movie-user matrix. The first left-singular vector \mathbf{u}_1 has large absolute values for the two science fiction movies and a large first singular value (red shading in Figure 4.10). Thus, this groups a type of users with a specific set of movies (science fiction theme). Similarly, the first right-singular \mathbf{v}_1 shows large absolute values for Ali and Beatrix, who give high ratings to science fiction movies (green shading in Figure 4.10). This suggests that \mathbf{v}_1 reflects the notion of a science fiction lover.

Similarly, \mathbf{u}_2 , seems to capture a French art house film theme, and \mathbf{v}_2 indicates that Chandra is close to an idealized lover of such movies. An idealized science fiction lover is a purist and only loves science fiction movies, so a science fiction lover \mathbf{v}_1 gives a rating of zero to everything but science fiction themed – this logic is implied the diagonal substructure for the singular value matrix Σ . A specific movie is therefore represented by how it decomposes (linearly) into its stereotypical movies. Likewise a person would be represented by how they decompose (via linear combination) into movie themes.

These two "spaces" are only meaningfully spanned by the respective viewer and movie data if the data itself covers a sufficient diversity of viewers and movies.

It is worth to briefly discuss SVD terminology and conventions as there are different versions used in the literature. The mathematics remains invariant to these differences, but these differences can be confudig.

- For convenience in notation and abstraction we use an SVD notation where the SVD is described as having two square left and right-singular vector matrices, but a non-square singular value matrix. Our definition (4.65) for the SVD is sometimes called the *full SVD*.
- Some authors define the SVD a bit differently and focus on square singular matrices. Then, for $\mathbf{A} \in \mathbb{R}^{m \times n}$ and $m \geq n$

full SVD

$$\mathbf{A}_{m \times n} = \mathbf{U}_{m \times n} \Sigma_{n \times n} \mathbf{V}_{n \times n}^T. \quad (4.90)$$

reduced SVD

Sometimes this formulation is called the *reduced SVD* (e.g., Datta (2010)) or *the SVD* (e.g., Press et al. (2007)). This alternative format changes merely how the matrices are constructed but leaves the mathematical structure of the SVD unchanged. The convenience of this alternative formulation is that Σ is diagonal, as in the eigenvalue decomposition.

truncated SVD

- In Section 4.6, we will learn about matrix approximation techniques using the SVD, which is also called the *truncated SVD*.
- It is possible to define the SVD of a rank- r matrix A so that U is an $m \times r$ matrix, Σ a diagonal matrix $r \times r$, and V an $r \times n$ matrix. This construction is very similar to our definition, and ensures that the diagonal matrix Σ has only non-zero entries along the diagonal. The main convenience of this alternative notation is that Σ is diagonal, as in the eigenvalue decomposition.
- A restriction that the SVD for A only applies to $m \times n$ matrices with $m > n$ is practically unnecessary. When $m < n$ the SVD decomposition will yield Σ with more zero columns than rows and, consequently, the singular values $\sigma_{m+1}, \dots, \sigma_n$ are 0.

The SVD is used in a variety of applications in machine learning from least squares problems in curve fitting to solving systems of linear equations. These applications harness various important properties of the SVD, its relation to the rank of a matrix and its ability to approximate matrices of a given rank with lower-rank matrices. Substituting a matrix with its SVD has often the advantage of making calculation more robust to numerical rounding errors. As we will explore in the next section the SVD's ability to approximate matrices with "simpler" matrices in a principled manner opens up machine learning applications ranging from dimensionality reduction and topic modeling to data compression and clustering.

4.6 Matrix Approximation

We considered the SVD as a way to factorize $A = U\Sigma V^\top \in \mathbb{R}^{m \times n}$ into the product of three matrices, where $U \in \mathbb{R}^{m \times m}$ and $V \in \mathbb{R}^{n \times n}$ are orthogonal and Σ contains the singular values on its main diagonal. Instead of doing the full SVD factorization, we will now investigate how the SVD allows us to represent a matrix A as a sum of simpler (low-rank) matrices A_i , which lends itself to a matrix approximation scheme that is cheaper to compute than the full SVD.

We construct a rank-1 matrix $A_i \in \mathbb{R}^{m \times n}$ as

$$A_i := u_i v_i^\top, \quad (4.91)$$

which is formed by the outer product of the i th orthogonal column vector of U and V . Figure 4.11 shows an image of Stonehenge, which can be represented by a matrix $A \in \mathbb{R}^{1432 \times 1910}$, and some outer products A_i , as defined in (4.91).

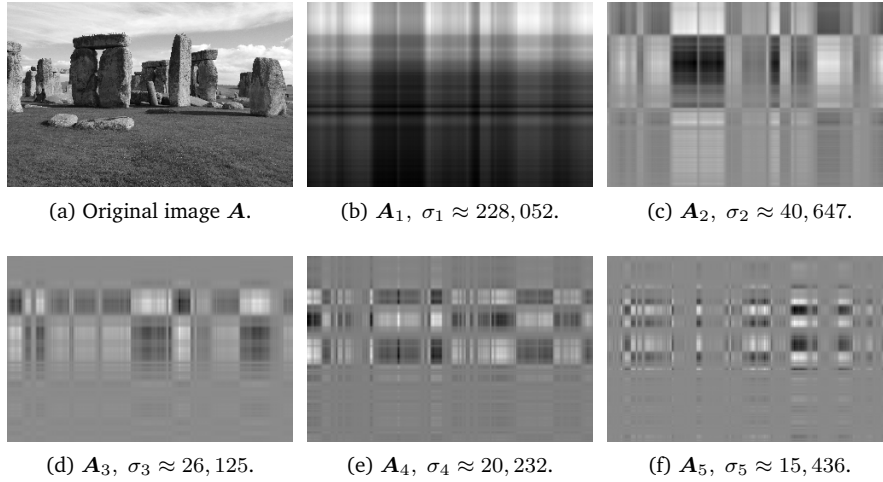


Figure 4.11 Image processing with the SVD. (a) The original grayscale image is a $1,432 \times 1,910$ matrix of values between 0 (black) and 1 (white). (b)–(f) Rank-1 matrices $\mathbf{A}_1, \dots, \mathbf{A}_5$ and their corresponding singular values $\sigma_1, \dots, \sigma_5$. The grid-like structure of each rank-1 matrix is imposed by the outer-product of the left and right-singular vectors.

A matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ of rank r can be written as a sum of rank-1 matrices \mathbf{A}_i so that

$$\mathbf{A} = \sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^\top = \sum_{i=1}^r \sigma_i \mathbf{A}_i, \quad (4.92)$$

where the outer-product matrices \mathbf{A}_i are weighted by the i th singular value σ_i . We can see why (4.92) holds: the diagonal structure of the singular value matrix $\mathbf{\Sigma}$ multiplies only matching left and right-singular vectors $\mathbf{u}_i \mathbf{v}_i^\top$ and scales them by the corresponding singular value σ_i . All terms $\Sigma_{ij} \mathbf{u}_i \mathbf{v}_j^\top$ vanish for $i \neq j$ because $\mathbf{\Sigma}$ is a diagonal matrix. Any terms $i > r$ vanish because the corresponding singular values are 0.

In (4.91), we introduced rank-1 matrices \mathbf{A}_i . We summed up the r individual rank-1 matrices to obtain a rank- r matrix \mathbf{A} , see (4.92). If the sum does not run over all matrices \mathbf{A}_i , $i = 1, \dots, r$, but only up to an intermediate value $k < r$, we obtain a *rank- k approximation*

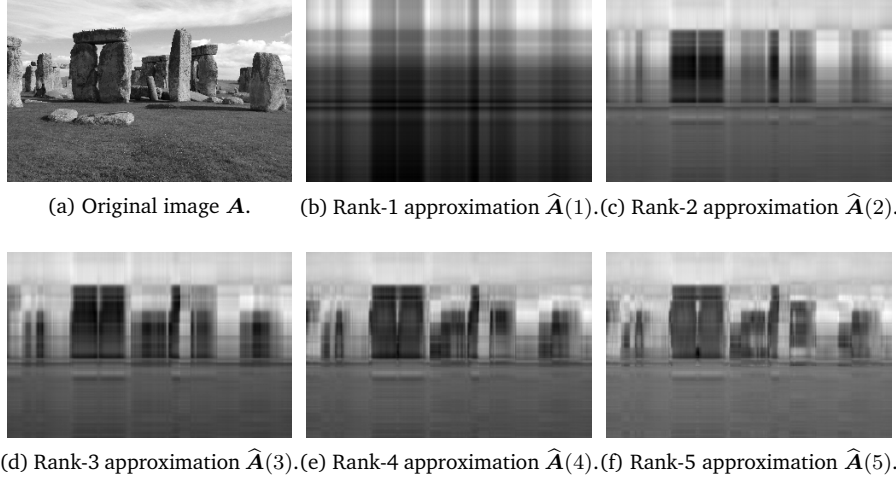
$$\hat{\mathbf{A}}(k) := \sum_{i=1}^k \sigma_i \mathbf{u}_i \mathbf{v}_i^\top = \sum_{i=1}^k \sigma_i \mathbf{A}_i \quad (4.93)$$

rank- k
approximation

of \mathbf{A} with $\text{rk}(\hat{\mathbf{A}}(k)) = k$. Figure 4.12 shows low-rank approximations $\hat{\mathbf{A}}(k)$ of an original image \mathbf{A} of Stonehenge. The shape of the rocks becomes increasingly visible and clearly recognizable in the rank-5 approximation. While the original image requires $1,432 \cdot 1,910 = 2,735,120$ numbers, the rank-5 approximation requires us only to store the five singular values and the five left and right singular vectors (1,432 and 1,910-dimensional each) for a total of $5 \cdot (1,432 + 1,910 + 1) = 16,715$ numbers – just above 0.6% of the original.

To measure the difference (error) between \mathbf{A} and its rank- k approximation $\hat{\mathbf{A}}(k)$ we need the notion of a norm. In Section 3.1, we already used

Figure 4.12 Image reconstruction with the SVD. (a) Original image. (b)–(f) Image reconstruction using the low-rank approximation of the SVD, where the rank- k approximation is given by $\hat{\mathbf{A}}(k) = \sum_{i=1}^k \sigma_i \mathbf{A}_i$.



norms on vectors that measure the length of a vector. By analogy we can also define norms on matrices.

spectral norm

Definition 4.23 (Spectral Norm of a Matrix). For $\mathbf{x} \in \mathbb{R}^n \setminus \{\mathbf{0}\}$ the *spectral norm* of a matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ is defined as

$$\|\mathbf{A}\|_2 := \max_{\mathbf{x}} \frac{\|\mathbf{A}\mathbf{x}\|_2}{\|\mathbf{x}\|_2}. \quad (4.94)$$

We introduce the notation of a subscript in the matrix norm (left-hand side), similar to the Euclidean norm for vectors (right-hand side), which has subscript 2. The spectral norm (4.94) determines how long any vector \mathbf{x} can at most become when multiplied by \mathbf{A} .

Theorem 4.24. *The spectral norm of \mathbf{A} is its largest singular value σ_1 .*

We leave the proof of this theorem as an exercise.

Theorem 4.25 (Eckart-Young Theorem (Eckart and Young, 1936)). *Consider a $\mathbf{A} \in \mathbb{R}^{m \times n}$ of rank r and let $\mathbf{B} \in \mathbb{R}^{m \times n}$ be a matrix of rank k . For any $k \leq r$ with $\hat{\mathbf{A}}(k) = \sum_{i=1}^k \sigma_i \mathbf{u}_i \mathbf{v}_i^\top$ it holds that*

$$\hat{\mathbf{A}}(k) = \operatorname{argmin}_{\operatorname{rk}(\mathbf{B}) \leq k} \|\mathbf{A} - \mathbf{B}\|_2, \quad (4.95)$$

$$\|\mathbf{A} - \hat{\mathbf{A}}(k)\|_2 = \sigma_{k+1}. \quad (4.96)$$

The Eckart-Young theorem states explicitly how much error we introduce by approximating \mathbf{A} using a rank- k approximation. We can interpret the rank- k approximation obtained with the SVD as a projection of the full-rank matrix \mathbf{A} onto a lower-dimensional space of rank-at-most- k matrices. Of all possible projections the SVD minimizes the error (with respect to the spectral norm) between \mathbf{A} and any rank- k approximation.

We can retrace some of the steps to understand why (4.96) should hold.

We observe that the difference between $\mathbf{A} - \hat{\mathbf{A}}(k)$ is a matrix containing the sum of the remaining rank-1 matrices

$$\mathbf{A} - \hat{\mathbf{A}}(k) = \sum_{i=k+1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^\top. \quad (4.97)$$

By Theorem 4.24, we immediately obtain σ_{k+1} as the spectral norm of the difference matrix. Let us have a closer look at (4.95). If we assume that there is another matrix \mathbf{B} with $\text{rk}(\mathbf{B}) \leq k$ such that

$$\|\mathbf{A} - \mathbf{B}\|_2 < \|\mathbf{A} - \hat{\mathbf{A}}(k)\|_2 \quad (4.98)$$

then there exists an $(n - k)$ -dimensional nullspace $Z \subseteq \mathbb{R}^n$ such that $\mathbf{x} \in Z$ implies that $\mathbf{B}\mathbf{x} = \mathbf{0}$. Then it follows that

$$\|\mathbf{A}\mathbf{x}\|_2 = \|(\mathbf{A} - \mathbf{B})\mathbf{x}\|_2, \quad (4.99)$$

and by using a version of the Cauchy-Schwartz inequality (3.17) that encompasses norms of matrices we obtain

$$\|\mathbf{A}\mathbf{x}\|_2 \leq \|\mathbf{A} - \mathbf{B}\|_2 \|\mathbf{x}\|_2 < \sigma_{k+1} \|\mathbf{x}\|_2. \quad (4.100)$$

However, there exists a $(k + 1)$ -dimensional subspace where $\|\mathbf{A}\mathbf{x}\|_2 \geq \sigma_{k+1} \|\mathbf{x}\|_2$, which is spanned by the right-singular vectors $\mathbf{v}_j, j \leq k + 1$ of \mathbf{A} . Adding up dimensions of these two spaces yields a number greater n , as there must be a non-zero vector in both spaces. This is a contradiction of the Rank-Nullity Theorem 2.23 in Section 2.7.3.

The Eckart-Young theorem implies that we can use SVD to reduce a rank- r matrix \mathbf{A} to a rank- k matrix $\hat{\mathbf{A}}$ in a principled, optimal (in the spectral norm sense) manner. We can interpret the approximation of \mathbf{A} by a rank- k matrix as a form of lossy compression. Therefore, the low-rank approximation of a matrix appears in many machine learning applications, e.g., image processing, noise filtering and regularization of ill-posed problems. Furthermore, it plays a key role in dimensionality reduction and principal component analysis as we will see in Chapter 10.

Example 4.15 (Finding Structure in Movie Ratings and Consumers (continued))

Coming back to our movie-rating example we can now apply the concept of low-rank approximations to approximate the original data matrix. Recall that our first singular value captures the notion of science fiction theme in movies and science fiction lovers. Thus, by using only the first singular value term in a rank-1 decomposition of the movie-rating matrix we obtain the predicted ratings

$$\mathbf{A}_1 = \mathbf{u}_1 \mathbf{v}_1^\top = \begin{bmatrix} -0.6710 \\ -0.7197 \\ -0.0939 \\ -0.1515 \end{bmatrix} \begin{bmatrix} -0.7367 & -0.6515 & -0.1811 \end{bmatrix} \quad (4.101a)$$

$$= \begin{bmatrix} 0.4943 & -0.0572 & -0.4500 \\ 0.5302 & -0.0613 & -0.4828 \\ 0.0692 & -0.0080 & -0.0630 \\ 0.1116 & -0.0129 & -0.1016 \end{bmatrix}. \quad (4.101b)$$

This first rank-1 approximation \mathbf{A}_1 is insightful: it tells us that Ali and Beatrix like science fiction movies, such as Star Wars and Bladerunner (entries have values > 4), but fails to capture the ratings of the other movies by Chandra. This is not surprising as Chandra's type of movies are not captured by the first singular value. The second singular value gives us a better rank-1 approximation for those movie-theme lovers:

$$\mathbf{A}_2 = \mathbf{u}_2 \mathbf{v}_2^\top = \begin{bmatrix} 0.0236 \\ 0.2054 \\ -0.7705 \\ -0.6030 \end{bmatrix} [0.0852 \quad 0.1762 \quad -0.9807] \quad (4.102a)$$

$$= \begin{bmatrix} -0.0154 & 0.0042 & -0.0174 \\ -0.1338 & 0.0362 & -0.1516 \\ 0.5019 & -0.1358 & 0.5686 \\ 0.3928 & -0.1063 & 0.445 \end{bmatrix} \quad (4.102b)$$

In this second rank-1 approximation, \mathbf{A}_2 we capture Chandra's ratings and movie types well, but for the science fiction movies. This leads us to consider the rank-2 approximation $\hat{\mathbf{A}}(2)$, where we combine the first two rank-1 approximations

$$\hat{\mathbf{A}}(2) = \sigma_1 \mathbf{A}_1 + \sigma_2 \mathbf{A}_2 = \begin{bmatrix} 4.7801 & 4.2419 & 1.0244 \\ 5.2252 & 4.7522 & -0.0250 \\ 0.2493 & -0.2743 & 4.9724 \\ 0.7495 & 0.2756 & 4.0278 \end{bmatrix}. \quad (4.103)$$

$\hat{\mathbf{A}}(2)$ is similar to the original movie ratings table

$$\mathbf{A} = \begin{bmatrix} 5 & 4 & 1 \\ 5 & 5 & 0 \\ 0 & 0 & 5 \\ 1 & 0 & 4 \end{bmatrix}, \quad (4.104)$$

and this suggests that we can ignore the contribution of \mathbf{A}_3 . We can interpret this so that in the data table there is no evidence of a third movie-theme/movie-lovers category. This also means that the entire space of movie-themes/movie-lovers in our example is a two-dimensional space spanned by science fiction and French art house movies and lovers.

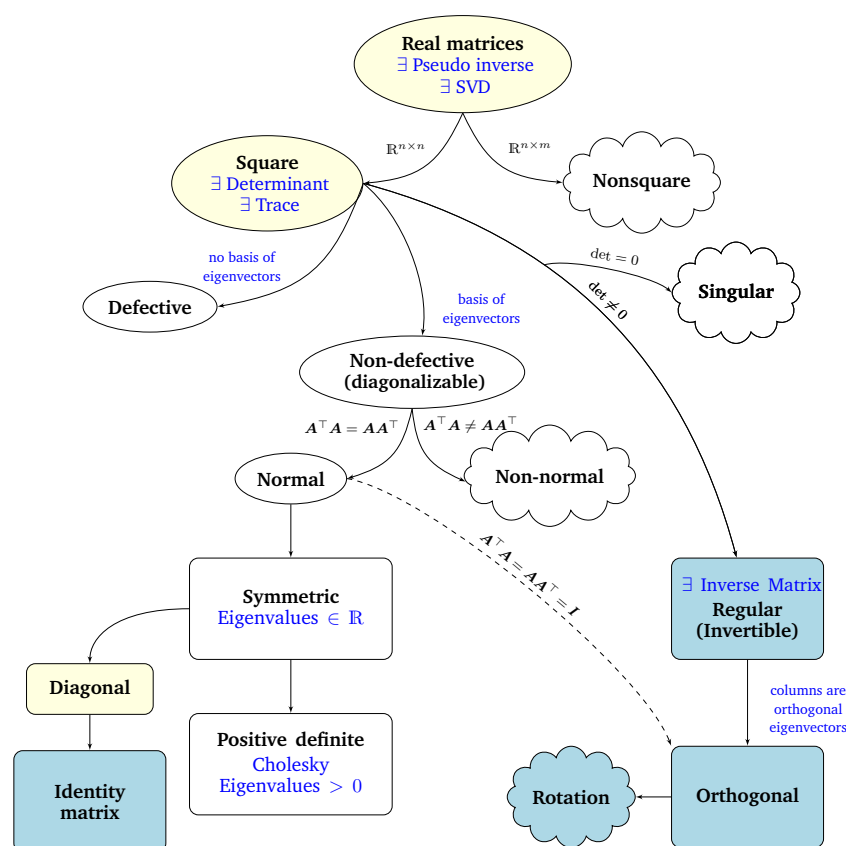


Figure 4.13 A functional phylogeny of matrices encountered in machine learning.

4.7 Matrix Phylogeny

In Chapters 2 and 3, we covered the basics of linear algebra and analytic geometry. In this chapter, we looked at fundamental characteristics of matrices and linear mappings. Figure 4.13 depicts the phylogenetic tree of relationships between different types of matrices (black arrows indicating “is a subset of”) and the covered operations we can perform on them (in blue). We consider all *real matrices* $A \in \mathbb{R}^{n \times m}$. For non-square matrices (where $n \neq m$), the SVD always exists as we saw in this chapter. Focusing on *square matrices* $A \in \mathbb{R}^{n \times n}$ the *determinant* informs us whether a square matrix possesses an *inverse matrix*, i.e., whether it belongs to the class of regular, invertible matrices. If the square $n \times n$ matrix possesses n linearly independent eigenvectors then the matrix is *non-defective* and an *eigendecomposition* exists (Theorem 4.12). We know that repeated eigenvalues may result in defective matrices, which cannot be diagonalized.

Non-singular and non-defective matrices are not the same. For example, a rotation matrix will be invertible (determinant is non-zero) but not diagonalizable in the real numbers (eigenvalues are not guaranteed to be real numbers).

The word “phylogenetic” describes how we capture the relationships among individuals or groups and derived from the Greek words for “tribe” and “source”.

We dive further into the branch of non-defective square $n \times n$ matrices. \mathbf{A} is *normal* if the condition $\mathbf{A}^\top \mathbf{A} = \mathbf{A} \mathbf{A}^\top$ holds. Moreover, if the more restrictive condition holds that $\mathbf{A}^\top \mathbf{A} = \mathbf{A} \mathbf{A}^\top = \mathbf{I}$, then \mathbf{A} is called *orthogonal* (see Definition 3.8). The set of orthogonal matrices is a subset of the regular (invertible) matrices and satisfies $\mathbf{A}^\top = \mathbf{A}^{-1}$.

Normal matrices have a frequently encountered subset, the symmetric matrices $\mathbf{S} \in \mathbb{R}^{n \times n}$ which satisfy $\mathbf{S} = \mathbf{S}^\top$. Symmetric matrices have only real eigenvalues. A subset of the symmetric matrices are the positive definite matrices \mathbf{P} that satisfy the condition of $\mathbf{x}^\top \mathbf{P} \mathbf{x} > 0$ for all $\mathbf{x} \in \mathbb{R}^n \setminus \{\mathbf{0}\}$. In this case, a unique *Cholesky decomposition* exists (Theorem 4.18). Positive definite matrices have only positive eigenvalues and are always invertible (i.e., have a non-zero determinant).

Another subset of symmetric matrices are the *diagonal matrices* \mathbf{D} . Diagonal matrices are closed under multiplication and addition, but do not necessarily form a group (this is only the case if all diagonal entries are non-zero so that the matrix is invertible). A special diagonal matrix is the identity matrix \mathbf{I} .

4.8 Further Reading

Most of the content in this chapter establishes underlying mathematics and connects them to methods for studying mappings, many of which are at the heart of machine learning at the level of underpinning software solutions and building blocks for almost all machine learning theory. Matrix characterization using determinants, eigenspectra and eigenspaces are fundamental features and conditions for categorizing and analyzing matrices. This extends to all forms of representations of data and mappings involving data, as well as judging the numerical stability of computational operations on such matrices (Press et al., 2007).

Determinants are fundamental tools in order to invert matrices and compute eigenvalues “by hand”. However, for almost all but the smallest instances numerical computation by Gaussian elimination outperforms determinants (Press et al., 2007). Determinants remain nevertheless a powerful theoretical concept, e.g., to gain intuition about the orientation of a basis based on the sign of the determinant. Eigenvectors can be used to perform basis changes to transform data into the coordinates of meaningful orthogonal, feature vectors. Similarly, matrix decomposition methods, such as the Cholesky decomposition, reappear often when we compute or simulate random events (Rubinstein and Kroese, 2016). Therefore, the Cholesky decomposition enables us to compute the *reparametrization trick* where we want to perform continuous differentiation over random variables, e.g., in variational autoencoders (Kingma and Ba, 2014; Jimenez Rezende et al., 2014).

Eigendecomposition is fundamental in enabling us to extract meaningful and interpretable information that characterizes linear mappings.

Therefore, the eigendecomposition underlies a general class of machine learning algorithms called *spectral methods* that perform eigendecomposition of a positive-definite kernel. These spectral decomposition methods encompass classical approaches to statistical data analysis, such as:

- *Principal component analysis* (PCA (Pearson, 1901), see also Chapter 10), in which a low-dimensional subspace, which explains most of the variability in the data, is sought.
- *Fisher discriminant analysis*, which aims to determine a separating hyperplane for data classification (Mika et al., 1999).
- *Multidimensional scaling* (MDS) (Carroll and Chang, 1970).

principal component
analysis

Fisher discriminant
analysis

multidimensional
scaling

The computational efficiency of these methods typically comes from finding the best rank- k approximation to a symmetric, positive semi-definite matrix. More contemporary examples of spectral methods have different origins, but each of them requires the computation of the eigenvectors and eigenvalues of a positive-definite kernel, such as *Isomap* (Tenenbaum et al., 2000), *Laplacian eigenmaps* (Belkin and Niyogi, 2003), *Hessian eigenmaps* (Donoho and Grimes, 2003), and *spectral clustering* (Shi and Malik, 2000). The core computations of these are generally underpinned by low-rank matrix approximation techniques (Belabbas and Wolfe, 2009) as we encountered here via the SVD.

Isomap
Laplacian
eigenmaps
Hessian eigenmaps
spectral clustering

The SVD allows us to discover some of the same kind of information as the eigendecomposition. However, the SVD is more generally applicable to non-square matrices and data tables. These matrix factorization methods become relevant whenever we want to identify heterogeneity in data when we want to perform data compression by approximation, e.g., instead of storing $n \times m$ values just storing $(n+m)k$ values, or when we want to perform data pre-processing, e.g., to decorrelate predictor variables of a design matrix (Ormonet et al., 2001). The SVD operates on matrices, which we can interpret as rectangular arrays with two indices (rows and columns). The extension of matrix-like structure to higher-dimensional arrays are called tensors. It turns out that the SVD is the special case of a more general family of decompositions that operate on such tensors (Kolda and Bader, 2009). SVD-like operations and low-rank approximations on tensors are for example the *Tucker decomposition* (Tucker, 1966) or the *CP decomposition* (Carroll and Chang, 1970).

Tucker
decomposition
CP decomposition

The SVD low-rank approximation is frequently used in machine learning for computational efficiency reasons. This is because it reduces the amount of memory and operations with non-zero multiplications we need to perform on potentially very large matrices of data (Trefethen and Bau III, 1997). Moreover, low-rank approximations are used to operate on matrices that may contain missing values as well as for purposes of lossy compression and dimensionality reduction (Moonen and De Moor, 1995; Markovsky, 2011).

Exercises

- 4.1 Compute the determinant using the Laplace expansion (using the first row) and the Sarrus Rule for

$$\mathbf{A} = \begin{bmatrix} 1 & 3 & 5 \\ 2 & 4 & 6 \\ 0 & 2 & 4 \end{bmatrix}. \quad (4.105)$$

- 4.2 Compute the following determinant efficiently:

$$\begin{bmatrix} 2 & 0 & 1 & 2 & 0 \\ 2 & -1 & 0 & 1 & 1 \\ 0 & 1 & 2 & 1 & 2 \\ -2 & 0 & 2 & -1 & 2 \\ 2 & 0 & 0 & 1 & 1 \end{bmatrix}. \quad (4.106)$$

- 4.3 Compute the eigenspaces of $\begin{bmatrix} 1 & 0 \\ 1 & 1 \end{bmatrix}$, $\begin{bmatrix} -2 & 2 \\ 2 & 1 \end{bmatrix}$.

- 4.4 Compute the eigenspaces of

$$\mathbf{A} = \begin{bmatrix} 0 & -1 & 1 & 1 \\ -1 & 1 & -2 & 3 \\ 2 & -1 & 0 & 0 \\ 1 & -1 & 1 & 0 \end{bmatrix}. \quad (4.107)$$

- 4.5 Diagonalizability of a matrix is unrelated to its invertibility. Determine for the following four matrices if it is diagonalizable and/or invertible

$$\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, \quad \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}, \quad \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix}, \quad \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix}. \quad (4.108)$$

- 4.6 Find the SVD of the matrix

$$\mathbf{A} = \begin{bmatrix} 3 & 2 & 2 \\ 2 & 3 & -2 \end{bmatrix}. \quad (4.109)$$

- 4.7 Find the singular value decomposition of

$$\mathbf{A} = \begin{bmatrix} 2 & 2 \\ -1 & 1 \end{bmatrix}. \quad (4.110)$$

- 4.8 Find the best rank-1 approximation of

$$\mathbf{A} = \begin{bmatrix} 3 & 2 & 2 \\ 2 & 3 & -2 \end{bmatrix}. \quad (4.111)$$

- 4.9 Show that for any $\mathbf{A} \in \mathbb{R}^{m \times n}$ the matrices $\mathbf{A}^\top \mathbf{A}$ and $\mathbf{A} \mathbf{A}^\top$ possess the same non-zero eigenvalues.

- 4.10 Show that for $\mathbf{x} \neq \mathbf{0}$ Theorem 4.24 holds, i.e., show that

$$\max_{\mathbf{x}} \frac{\|\mathbf{A}\mathbf{x}\|_2}{\|\mathbf{x}\|_2} = \sigma_1, \quad (4.112)$$

where σ_1 is the largest singular value of $\mathbf{A} \in \mathbb{R}^{m \times n}$.