



4.1 Determinant and Trace

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 $A \in \mathbb{R}^{n \times n}$, i.e., matrices with the same number of rows and columns. In this book, we write the determinant as $\det(A)$ or sometimes as $|A|$ so that

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

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

©2021 M. P. Deisenroth, A. A. Faisal, C. S. Ong. Published by Cambridge University Press (2020).

100

Matrix Decompositions

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

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

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

Hence, 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 $A \in \mathbb{R}^{2 \times 2}$, i.e.,

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

Example 4.1 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 $A \in \mathbb{R}^{n \times n}$ it holds that A is invertible if and only if $\det(A) \neq 0$.

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

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

For $n = 2$,

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

which we have observed in the preceding example.

For $n = 3$ (known as **Sarrus' rule**),

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

Draft (2021-01-14) of "Mathematics for Machine Learning". Feedback: <https://mml-book.com>.

4.1 Determinant and Trace

101

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 $i > j$, i.e., the matrix is zero below its diagonal. Analogously, we define a **lower-triangular matrix** as a matrix with zeros above its diagonal. For a triangular matrix $T \in \mathbb{R}^{n \times n}$, the determinant is the product of the diagonal elements, i.e.,

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

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 $\det(A)$ is the signed volume of an n -dimensional parallelepiped formed by columns of the matrix 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 b, g that form the columns of a matrix $A = [b, g]$. Then, the absolute value of the determinant of A is the area of the parallelogram with vertices $0, b, g, b + g$. In particular, if b, g are linearly dependent so that $b = \lambda g$ for some $\lambda \in \mathbb{R}$, they no longer form a two-dimensional parallelogram. Therefore, the corresponding area is 0. On the contrary, if b, g are linearly independent and are multiples of the canonical basis vectors e_1, e_2 then they can be written as $b = \begin{bmatrix} b_1 \\ 0 \end{bmatrix}$ and $g = \begin{bmatrix} 0 \\ g_2 \end{bmatrix}$, and the determinant is $\begin{vmatrix} b_1 & 0 \\ 0 & g_2 \end{vmatrix} = b_1 g_2 - 0 = b_1 g_2$.

The sign of the determinant indicates the orientation of the spanning vectors b, g with respect to the standard basis (e_1, e_2) . In our figure, flipping the order to g, b swaps the columns of A and reverses the orientation of the shaded area. This becomes the familiar formula: area = height \times length. This intuition extends to higher dimensions. In \mathbb{R}^3 , we consider three vectors $r, b, g \in \mathbb{R}^3$ 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 $[r, b, g]$ 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 $r, b, g \in \mathbb{R}^3$ given as

$$r = \begin{bmatrix} 2 \\ 0 \\ -8 \end{bmatrix}, \quad g = \begin{bmatrix} 6 \\ 1 \\ 0 \end{bmatrix}, \quad b = \begin{bmatrix} 1 \\ 4 \\ -1 \end{bmatrix}. \quad (4.9)$$

Writing these vectors as the columns of a matrix

$$A = [r, g, b] = \begin{bmatrix} 2 & 6 & 1 \\ 0 & 1 & 4 \\ -8 & 0 & -1 \end{bmatrix} \quad (4.10)$$

allows us to compute the desired volume as

$$V = |\det(A)| = 186. \quad (4.11)$$

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. Theorem 4.2 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 **Laplace expansion** (Theorem 4.2), 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 $A \in \mathbb{R}^{n \times n}$. Then, for all $j = 1, \dots, n$:

1. Expansion along column j

$$\det(A) = \sum_{k=1}^n (-1)^{k+j} a_{kj} \det(A_{k,j}). \quad (4.12)$$

2. Expansion along row j

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

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

Example 4.3 (Laplace Expansion)

Let us compute the determinant of

$$A = \begin{bmatrix} 1 & 2 & 3 \\ 3 & 1 & 2 \\ 0 & 0 & 1 \end{bmatrix} \quad (4.14)$$

using the Laplace expansion along the first row. Applying (4.13) yields

$$\begin{aligned} \det(A) &= \begin{vmatrix} 1 & 2 & 3 \\ 3 & 1 & 2 \\ 0 & 0 & 1 \end{vmatrix} = (-1)^{1+1} \cdot 1 \cdot \begin{vmatrix} 1 & 2 \\ 0 & 1 \end{vmatrix} \\ &\quad + (-1)^{1+2} \cdot 2 \cdot \begin{vmatrix} 3 & 2 \\ 0 & 1 \end{vmatrix} + (-1)^{1+3} \cdot 3 \cdot \begin{vmatrix} 3 & 1 \\ 0 & 0 \end{vmatrix}. \end{aligned} \quad (4.15)$$

Because of the last three properties, we can use **Gaussian elimination** (see Section 2.1) to compute $\det(A)$ by bringing A into row-echelon form. We can **stop** Gaussian elimination when we have 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 $A \in \mathbb{R}^{n \times n}$ has $\det(A) \neq 0$ if and only if $\text{rk}(A) = n$. In other words, 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 $A \in \mathbb{R}^{n \times n}$ is defined as

$$\text{tr}(A) := \sum_{i=1}^n a_{ii}, \quad (4.18)$$

i.e., the **trace is the sum of the diagonal elements of A** .

The trace satisfies the following properties:

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

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.,

$$\text{tr}(AKL) = \text{tr}(KLA) \quad (4.19)$$

for matrices $A \in \mathbb{R}^{n \times k}$, $K \in \mathbb{R}^{k \times l}$, $L \in \mathbb{R}^{l \times n}$. **This property generalizes to products of an arbitrary number of matrices.** As a special case of (4.19), it follows that for **two vectors $x, y \in \mathbb{R}^n$**

$$\text{tr}(xy^T) = \text{tr}(y^T x) = y^T x \in \mathbb{R}. \quad (4.20)$$

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 A . Then the trace of Φ is the trace of A . For a different basis of V , it holds that the corresponding transformation matrix B of Φ can be obtained by a basis change of the form $S^{-1}AS$ for suitable S (see Section 2.7.2). For the corresponding trace of Φ , this means

$$\text{tr}(B) = \text{tr}(S^{-1}AS) \stackrel{(4.19)}{=} \text{tr}(ASS^{-1}) = \text{tr}(A). \quad (4.21)$$

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 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 $A \in \mathbb{R}^{n \times n}$

$$p_A(\lambda) := \det(A - \lambda I) \quad (4.22a)$$

$$= c_0 + c_1\lambda + c_2\lambda^2 + \dots + c_{n-1}\lambda^{n-1} + (-1)^n\lambda^n, \quad (4.22b)$$

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

$$c_0 = \det(A), \quad (4.23)$$

$$c_{n-1} = (-1)^{n-1}\text{tr}(A). \quad (4.24)$$

The characteristic polynomial (4.22a) 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, is transformed by the linear mapping**.

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

$$Ax = \lambda x. \quad (4.25)$$

We call (4.25) 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.

The following statements are equivalent:

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

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 x is an eigenvector of A associated with eigenvalue λ , then for any $c \in \mathbb{R} \setminus \{0\}$ it holds that **cx is an eigenvector of A with the same eigenvalue since**

$$A(cx) = cAx = c\lambda x = \lambda(cx). \quad (4.26)$$

Thus, all vectors **that are collinear to x** are also **eigenvectors of A** .

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

A *diagonal matrix* is a matrix that has value zero on all off-diagonal elements, i.e., they are of the form

$$\mathbf{D} = \begin{bmatrix} c_1 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & c_n \end{bmatrix}. \quad (4.49)$$

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

In this section, we will discuss how to transform matrices into diagonal

- 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

$$\mathbf{A}^k = (\mathbf{P}\mathbf{D}\mathbf{P}^{-1})^k = \mathbf{P}\mathbf{D}^k\mathbf{P}^{-1}. \quad (4.62)$$

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

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

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