

# MATHEMATICS FOR MACHINE LEARNING

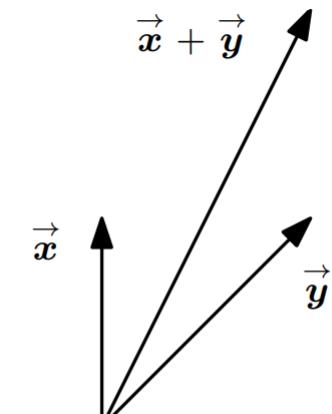
---

TUYEN NGOC LE

# Vector

---

1. Geometric vectors. This example of a vector may be familiar from high school mathematics and physics. Geometric vectors – see Figure 2.1(a) – are directed segments, which can be drawn (at least in two dimensions). Two geometric vectors  $\vec{x}$ ,  $\vec{y}$  can be added, such that  $\vec{x} + \vec{y} = \vec{z}$  is another geometric vector. Furthermore, multiplication by a scalar  $\lambda \vec{x}$ ,  $\lambda \in \mathbb{R}$ , is also a geometric vector. In fact, it is the original vector scaled by  $\lambda$ . Therefore, geometric vectors are instances of the vector concepts introduced previously. Interpreting vectors as geometric vectors enables us to use our intuitions about direction and magnitude to reason about mathematical operations.

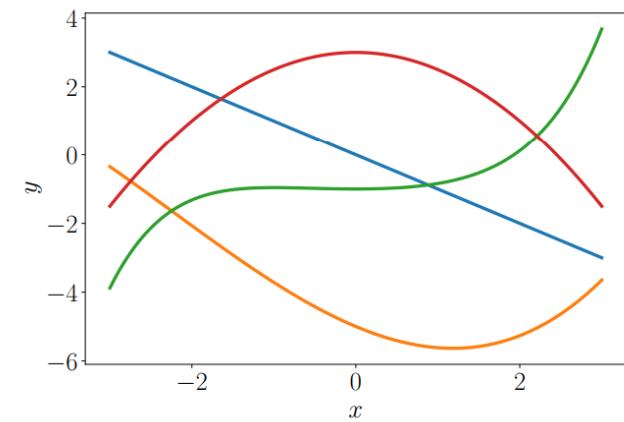


(a) Geometric vectors.

# Vector

---

Polynomials are also vectors; see Figure 2.1(b): Two polynomials can be added together, which results in another polynomial; and they can be multiplied by a scalar  $\lambda \in \mathbb{R}$ , and the result is a polynomial as well. Therefore, polynomials are (rather unusual) instances of vectors. Note that polynomials are very different from geometric vectors. While geometric vectors are concrete “drawings”, polynomials are abstract concepts. However, they are both vectors in the sense previously described.

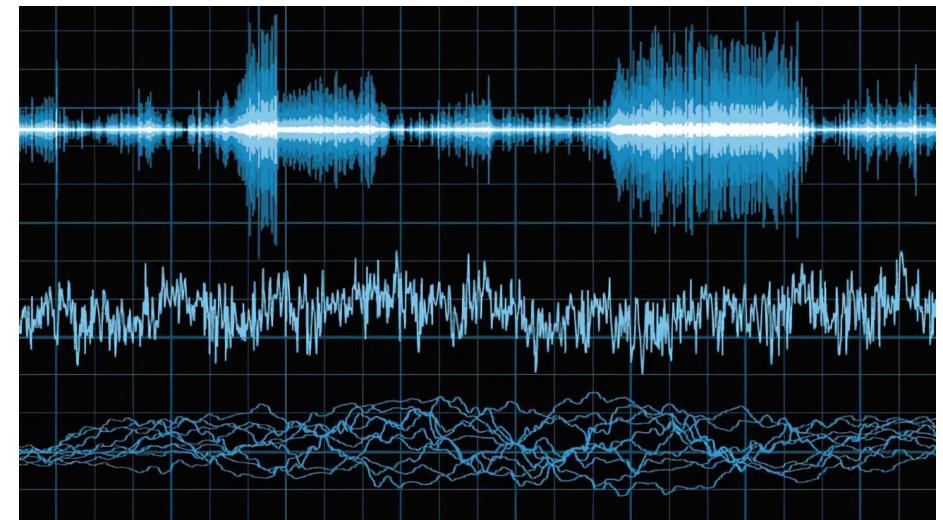


(b) Polynomials.

# Vector

---

Audio signals are vectors. Audio signals are represented as a series of numbers. We can add audio signals together, and their sum is a new audio signal. If we scale an audio signal, we also obtain an audio signal. Therefore, audio signals are a type of vector, too.



# Vector

---

Elements of  $\mathbb{R}^n$  (tuples of  $n$  real numbers) are vectors.  $\mathbb{R}^n$  is more abstract than polynomials, and it is the concept we focus on in this book. For instance,

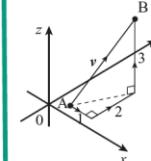
$$\mathbf{a} = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} \in \mathbb{R}^3 \quad (2.1)$$

is an example of a triplet of numbers. Adding two vectors  $\mathbf{a}, \mathbf{b} \in \mathbb{R}^n$  component-wise results in another vector:  $\mathbf{a} + \mathbf{b} = \mathbf{c} \in \mathbb{R}^n$ . Moreover, multiplying  $\mathbf{a} \in \mathbb{R}^n$  by  $\lambda \in \mathbb{R}$  results in a scaled vector  $\lambda\mathbf{a} \in \mathbb{R}^n$ . Considering vectors as elements of  $\mathbb{R}^n$  has an additional benefit that it loosely corresponds to arrays of real numbers on a computer. Many programming languages support array operations, which allow for convenient implementation of algorithms that involve vector operations.

## What is a 3-dimensional vector?

3D vectors are described using components parallel to the  $x$ -,  $y$ - and  $z$ -axes. In the diagram, vector  $\mathbf{v}$ , represented by  $\overline{AB}$ , has components

$$\begin{pmatrix} 1 \\ 2 \\ 3 \end{pmatrix}.$$



$$\mathbf{v} = \begin{pmatrix} 1 \\ 2 \\ 3 \end{pmatrix} \leftarrow \begin{array}{l} x\text{-component} \\ \leftarrow y\text{-component} \\ \leftarrow z\text{-component} \end{array}$$

3D vectors can be added and subtracted by adding and subtracting their corresponding components. For example if

$$\mathbf{a} = \begin{pmatrix} -1 \\ 1 \\ 3 \end{pmatrix} \text{ and } \mathbf{b} = \begin{pmatrix} -1 \\ -2 \\ 5 \end{pmatrix}$$

then

$$\mathbf{a} + \mathbf{b} = \begin{pmatrix} -1 \\ 1 \\ 3 \end{pmatrix} + \begin{pmatrix} -1 \\ -2 \\ 5 \end{pmatrix} = \begin{pmatrix} -1 + (-1) \\ 1 + (-2) \\ 3 + 5 \end{pmatrix} = \begin{pmatrix} -2 \\ -1 \\ 8 \end{pmatrix}$$

$$\mathbf{a} - \mathbf{b} = \begin{pmatrix} -1 \\ 1 \\ 3 \end{pmatrix} - \begin{pmatrix} -1 \\ -2 \\ 5 \end{pmatrix} = \begin{pmatrix} -1 - (-1) \\ 1 - (-2) \\ 3 - 5 \end{pmatrix} = \begin{pmatrix} 0 \\ 3 \\ -2 \end{pmatrix}$$

# Outer product

---

Given two vectors of size  $m \times 1$  and  $n \times 1$  respectively

$$\mathbf{u} = \begin{bmatrix} u_1 \\ u_2 \\ \vdots \\ u_m \end{bmatrix}, \quad \mathbf{v} = \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_n \end{bmatrix}$$

their outer product, denoted  $\mathbf{u} \otimes \mathbf{v}$ , is defined as the  $m \times n$  matrix  $\mathbf{A}$  obtained by multiplying each element of  $\mathbf{u}$  by each element of  $\mathbf{v}$ :<sup>[1]</sup>

$$\mathbf{u} \otimes \mathbf{v} = \mathbf{A} = \begin{bmatrix} u_1 v_1 & u_1 v_2 & \dots & u_1 v_n \\ u_2 v_1 & u_2 v_2 & \dots & u_2 v_n \\ \vdots & \vdots & \ddots & \vdots \\ u_m v_1 & u_m v_2 & \dots & u_m v_n \end{bmatrix}$$

Or in index notation:

$$(\mathbf{u} \otimes \mathbf{v})_{ij} = u_i v_j$$

# Matrix

---

**Definition 2.1** (Matrix). With  $m, n \in \mathbb{N}$  a real-valued  $(m, n)$  matrix  $A$  is an  $m \cdot n$ -tuple of elements  $a_{ij}$ ,  $i = 1, \dots, m$ ,  $j = 1, \dots, n$ , which is ordered according to a rectangular scheme consisting of  $m$  rows and  $n$  columns:

$$A = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mn} \end{bmatrix}, \quad a_{ij} \in \mathbb{R}. \quad (2.11)$$

By convention  $(1, n)$ -matrices are called *rows* and  $(m, 1)$ -matrices are called *columns*. These special matrices are also called *row/column vectors*.

$\mathbb{R}^{m \times n}$  is the set of all real-valued  $(m, n)$ -matrices.  $A \in \mathbb{R}^{m \times n}$  can be equivalently represented as  $a \in \mathbb{R}^{mn}$  by stacking all  $n$  columns of the matrix into a long vector; see Figure 2.4.

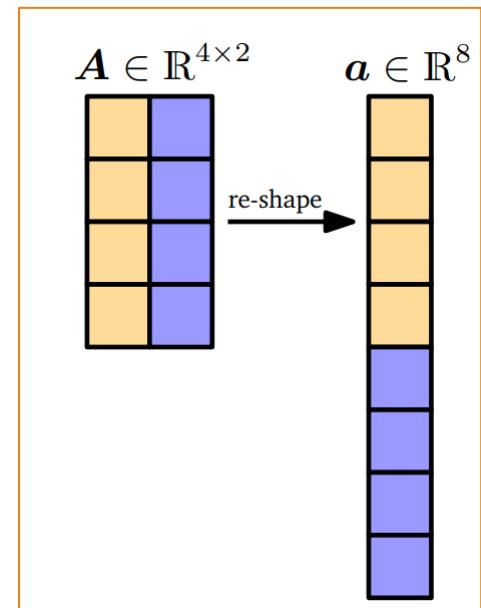


Figure 2.4 By stacking its columns, a matrix  $A$  can be represented as a long vector  $a$ .

# Matrix Addition and Multiplication

---

The sum of two matrices  $A \in \mathbb{R}^{m \times n}$ ,  $B \in \mathbb{R}^{m \times n}$  is defined as the element-wise sum, i.e.,

$$A + B := \begin{bmatrix} a_{11} + b_{11} & \cdots & a_{1n} + b_{1n} \\ \vdots & & \vdots \\ a_{m1} + b_{m1} & \cdots & a_{mn} + b_{mn} \end{bmatrix} \in \mathbb{R}^{m \times n}. \quad (2.12)$$

# Matrix Addition and Multiplication

---

For matrices  $\mathbf{A} \in \mathbb{R}^{m \times n}$ ,  $\mathbf{B} \in \mathbb{R}^{n \times k}$ , the elements  $c_{ij}$  of the product  $\mathbf{C} = \mathbf{AB} \in \mathbb{R}^{m \times k}$  are computed as

$$c_{ij} = \sum_{l=1}^n a_{il} b_{lj}, \quad i = 1, \dots, m, \quad j = 1, \dots, k. \quad (2.13)$$

# Matrix Addition and Multiplication

---

For  $A = \begin{bmatrix} 1 & 2 & 3 \\ 3 & 2 & 1 \end{bmatrix} \in \mathbb{R}^{2 \times 3}$ ,  $B = \begin{bmatrix} 0 & 2 \\ 1 & -1 \\ 0 & 1 \end{bmatrix} \in \mathbb{R}^{3 \times 2}$ , we obtain

$$AB = \begin{bmatrix} 1 & 2 & 3 \\ 3 & 2 & 1 \end{bmatrix} \begin{bmatrix} 0 & 2 \\ 1 & -1 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} 2 & 3 \\ 2 & 5 \end{bmatrix} \in \mathbb{R}^{2 \times 2},$$



$$BA = \begin{bmatrix} 0 & 2 \\ 1 & -1 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 2 & 3 \\ 3 & 2 & 1 \end{bmatrix} = \begin{bmatrix} 6 & 4 & 2 \\ -2 & 0 & 2 \\ 3 & 2 & 1 \end{bmatrix} \in \mathbb{R}^{3 \times 3}.$$

# Matrix Addition and Multiplication

---

For  $A = \begin{bmatrix} 1 & 2 & 3 \\ 3 & 2 & 1 \end{bmatrix} \in \mathbb{R}^{2 \times 3}$ ,  $B = \begin{bmatrix} 0 & 2 \\ 1 & -1 \\ 0 & 1 \end{bmatrix} \in \mathbb{R}^{3 \times 2}$ , we obtain

$$AB = \begin{bmatrix} 1 & 2 & 3 \\ 3 & 2 & 1 \end{bmatrix} \begin{bmatrix} 0 & 2 \\ 1 & -1 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} 2 & 3 \\ 2 & 5 \end{bmatrix} \in \mathbb{R}^{2 \times 2},$$



$$BA = \begin{bmatrix} 0 & 2 \\ 1 & -1 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 2 & 3 \\ 3 & 2 & 1 \end{bmatrix} = \begin{bmatrix} 6 & 4 & 2 \\ -2 & 0 & 2 \\ 3 & 2 & 1 \end{bmatrix} \in \mathbb{R}^{3 \times 3}.$$

# Identity Matrix

---

$$\boldsymbol{I}_n := \begin{bmatrix} 1 & 0 & \cdots & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 1 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 0 & \cdots & 1 \end{bmatrix} \in \mathbb{R}^{n \times n}$$

# Properties of matrices

---

- *Associativity:*

$$\forall \mathbf{A} \in \mathbb{R}^{m \times n}, \mathbf{B} \in \mathbb{R}^{n \times p}, \mathbf{C} \in \mathbb{R}^{p \times q} : (\mathbf{AB})\mathbf{C} = \mathbf{A}(\mathbf{BC})$$

- *Distributivity:*

$$\begin{aligned}\forall \mathbf{A}, \mathbf{B} \in \mathbb{R}^{m \times n}, \mathbf{C}, \mathbf{D} \in \mathbb{R}^{n \times p} : & (\mathbf{A} + \mathbf{B})\mathbf{C} = \mathbf{AC} + \mathbf{BC} \\ & \mathbf{A}(\mathbf{C} + \mathbf{D}) = \mathbf{AC} + \mathbf{AD}\end{aligned}$$

- Multiplication with the identity matrix:

$$\forall \mathbf{A} \in \mathbb{R}^{m \times n} : \mathbf{I}_m \mathbf{A} = \mathbf{A} \mathbf{I}_n = \mathbf{A}$$

Note that  $\mathbf{I}_m \neq \mathbf{I}_n$  for  $m \neq n$ .

# Inverse and Transpose

---

**Definition 2.3** (Inverse). Consider a square matrix  $A \in \mathbb{R}^{n \times n}$ . Let matrix  $B \in \mathbb{R}^{n \times n}$  have the property that  $AB = I_n = BA$ .  $B$  is called the *inverse* of  $A$  and denoted by  $A^{-1}$ .

## Example 2.4 (Inverse Matrix)

The matrices

$$A = \begin{bmatrix} 1 & 2 & 1 \\ 4 & 4 & 5 \\ 6 & 7 & 7 \end{bmatrix}, \quad B = \begin{bmatrix} -7 & -7 & 6 \\ 2 & 1 & -1 \\ 4 & 5 & -4 \end{bmatrix}$$

are inverse to each other since  $AB = I = BA$ .

# Inverse and Transpose

---

**Definition 2.4** (Transpose). For  $A \in \mathbb{R}^{m \times n}$  the matrix  $B \in \mathbb{R}^{n \times m}$  with  $b_{ij} = a_{ji}$  is called the *transpose* of  $A$ . We write  $B = A^\top$ .

$$\begin{bmatrix} 1 & 2 & 3 \\ 0 & -6 & 7 \end{bmatrix}^\top = \begin{bmatrix} 1 & 0 \\ 2 & -6 \\ 3 & 7 \end{bmatrix}$$

# Properties of inverses and transposes

---

$$\mathbf{A}\mathbf{A}^{-1} = \mathbf{I} = \mathbf{A}^{-1}\mathbf{A}$$

$$(\mathbf{AB})^{-1} = \mathbf{B}^{-1}\mathbf{A}^{-1}$$

$$(\mathbf{A} + \mathbf{B})^{-1} \neq \mathbf{A}^{-1} + \mathbf{B}^{-1}$$

$$(\mathbf{A}^\top)^\top = \mathbf{A}$$

$$(\mathbf{A} + \mathbf{B})^\top = \mathbf{A}^\top + \mathbf{B}^\top$$

$$(\mathbf{AB})^\top = \mathbf{B}^\top \mathbf{A}^\top$$

# Multiplication by a Scalar

---

Let us look at what happens to matrices when they are multiplied by a scalar  $\lambda \in \mathbb{R}$ . Let  $A \in \mathbb{R}^{m \times n}$  and  $\lambda \in \mathbb{R}$ . Then  $\lambda A = K$ ,  $K_{ij} = \lambda a_{ij}$ . Practically,  $\lambda$  scales each element of  $A$ . For  $\lambda, \psi \in \mathbb{R}$ , the following holds:

- *Associativity:*

$$(\lambda\psi)C = \lambda(\psi C), \quad C \in \mathbb{R}^{m \times n}$$

- $\lambda(BC) = (\lambda B)C = B(\lambda C) = (BC)\lambda, \quad B \in \mathbb{R}^{m \times n}, C \in \mathbb{R}^{n \times k}$ .

Note that this allows us to move scalar values around.

- $(\lambda C)^\top = C^\top \lambda^\top = C^\top \lambda = \lambda C^\top$  since  $\lambda = \lambda^\top$  for all  $\lambda \in \mathbb{R}$ .

- *Distributivity:*

$$(\lambda + \psi)C = \lambda C + \psi C, \quad C \in \mathbb{R}^{m \times n}$$

$$\lambda(B + C) = \lambda B + \lambda C, \quad B, C \in \mathbb{R}^{m \times n}$$

# Compact Representations of Systems of Linear Equations

If we consider the system of linear equations

$$\begin{aligned} 2x_1 + 3x_2 + 5x_3 &= 1 \\ 4x_1 - 2x_2 - 7x_3 &= 8 \\ 9x_1 + 5x_2 - 3x_3 &= 2 \end{aligned} \tag{2.35}$$

and use the rules for matrix multiplication, we can write this equation system in a more compact form as

$$\begin{bmatrix} 2 & 3 & 5 \\ 4 & -2 & -7 \\ 9 & 5 & -3 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 1 \\ 8 \\ 2 \end{bmatrix}. \tag{2.36}$$

$$Ax = b \iff x = A^{-1}b \quad (\text{if } A \text{ is a square matrix and invertible})$$

$$Ax = b \iff A^\top Ax = A^\top b \iff x = (A^\top A)^{-1}A^\top b$$

# Rank

---

The number of linearly independent columns of a matrix  $A \in \mathbb{R}^{m \times n}$  equals the number of linearly independent rows and is called the *rank* of  $A$  and is denoted by  $\text{rk}(A)$ .

# Rank's properties

---

- $\text{rk}(\mathbf{A}) = \text{rk}(\mathbf{A}^\top)$ , i.e., the column rank equals the row rank.
- The columns of  $\mathbf{A} \in \mathbb{R}^{m \times n}$  span a subspace  $U \subseteq \mathbb{R}^m$  with  $\dim(U) = \text{rk}(\mathbf{A})$ . Later we will call this subspace the *image* or *range*. A basis of  $U$  can be found by applying Gaussian elimination to  $\mathbf{A}$  to identify the pivot columns.
- The rows of  $\mathbf{A} \in \mathbb{R}^{m \times n}$  span a subspace  $W \subseteq \mathbb{R}^n$  with  $\dim(W) = \text{rk}(\mathbf{A})$ . A basis of  $W$  can be found by applying Gaussian elimination to  $\mathbf{A}^\top$ .
- For all  $\mathbf{A} \in \mathbb{R}^{n \times n}$  it holds that  $\mathbf{A}$  is regular (invertible) if and only if  $\text{rk}(\mathbf{A}) = n$ .

# Rank's properties

---

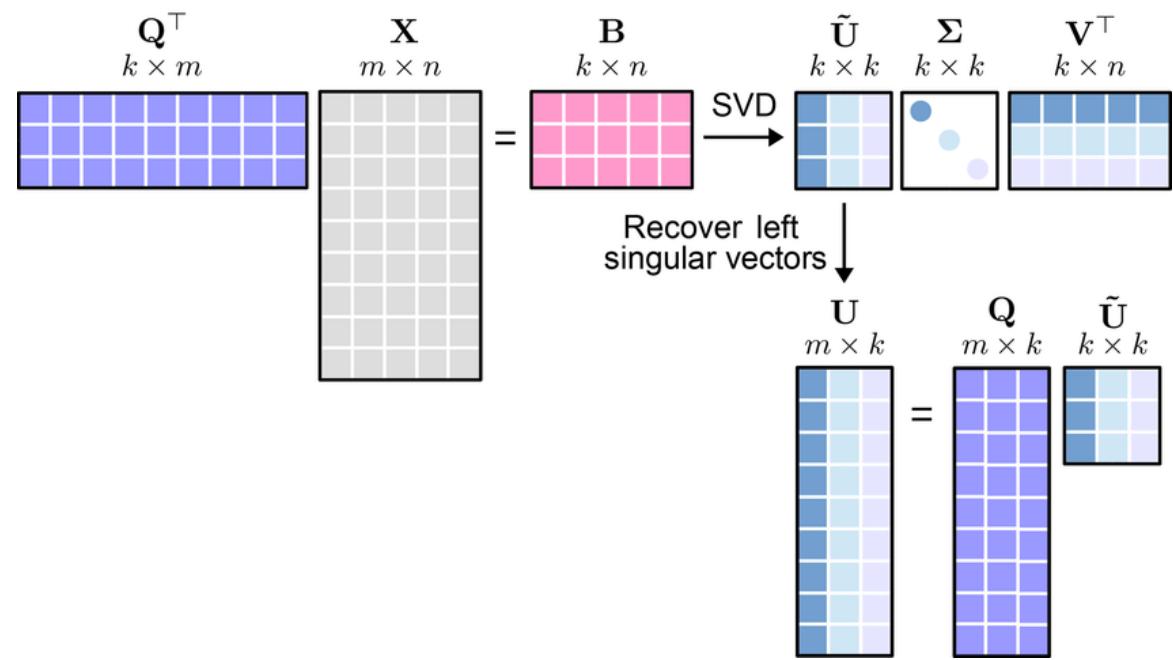
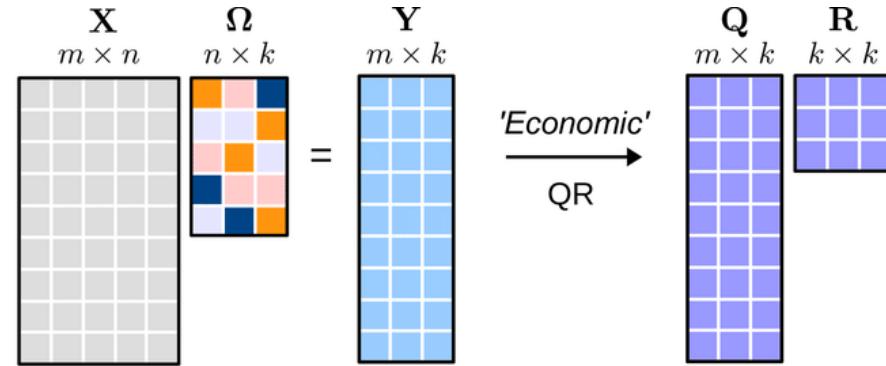
- For all  $A \in \mathbb{R}^{m \times n}$  and all  $b \in \mathbb{R}^m$  it holds that the linear equation system  $Ax = b$  can be solved if and only if  $\text{rk}(A) = \text{rk}(A|b)$ , where  $A|b$  denotes the augmented system.
- For  $A \in \mathbb{R}^{m \times n}$  the subspace of solutions for  $Ax = 0$  possesses dimension  $n - \text{rk}(A)$ . Later, we will call this subspace the *kernel* or the *null space*.
- A matrix  $A \in \mathbb{R}^{m \times n}$  has *full rank* if its rank equals the largest possible rank for a matrix of the same dimensions. This means that the rank of a full-rank matrix is the lesser of the number of rows and columns, i.e.,  $\text{rk}(A) = \min(m, n)$ . A matrix is said to be *rank deficient* if it does not have full rank.

# Rank's properties

---

- For all  $A \in \mathbb{R}^{m \times n}$  and all  $b \in \mathbb{R}^m$  it holds that the linear equation system  $Ax = b$  can be solved if and only if  $\text{rk}(A) = \text{rk}(A|b)$ , where  $A|b$  denotes the augmented system.
- For  $A \in \mathbb{R}^{m \times n}$  the subspace of solutions for  $Ax = 0$  possesses dimension  $n - \text{rk}(A)$ . Later, we will call this subspace the *kernel* or the *null space*.
- A matrix  $A \in \mathbb{R}^{m \times n}$  has *full rank* if its rank equals the largest possible rank for a matrix of the same dimensions. This means that the rank of a full-rank matrix is the lesser of the number of rows and columns, i.e.,  $\text{rk}(A) = \min(m, n)$ . A matrix is said to be *rank deficient* if it does not have full rank.

# Matrix Decompositions



# Determinant

---

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

# Determinant

---

For  $n = 1$ ,  $\det(\mathbf{A}) = \det(a_{11}) = a_{11}$

For  $n = 2$ ,

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

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

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

# Determinant

---

**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  $j$*

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

2. *Expansion along row  $j$*

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

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  $j$ .

# Determinant

---

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 2.22) possess the same determinant. Therefore, for a linear mapping  $\Phi : V \rightarrow V$  all transformation matrices  $\mathbf{A}_\Phi$  of  $\Phi$  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  $\lambda$ . In particular,  $\det(\lambda\mathbf{A}) = \lambda^n \det(\mathbf{A})$ .
- Swapping two rows/columns changes the sign of  $\det(\mathbf{A})$ .

# Determinant

---

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

**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.*

# Trace

---

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

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

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  $A \in \mathbb{R}^{n \times k}$ ,  $B \in \mathbb{R}^{k \times n}$

# Eigenvalues and Eigenvectors

---

**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  $x \in \mathbb{R}^n \setminus \{\mathbf{0}\}$  is the corresponding *eigenvector* of  $A$  if

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

We call (4.25) the *eigenvalue equation*.

The following statements are equivalent:

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

# Eigenvalues and Eigenvectors

---

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

$$p_{\mathbf{A}}(\lambda) := \det(\mathbf{A} - \lambda \mathbf{I}) \tag{4.22a}$$

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

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

$$\begin{aligned} c_0 &= \det(\mathbf{A}), \\ c_{n-1} &= (-1)^{n-1} \text{tr}(\mathbf{A}). \end{aligned}$$

# Computing Eigenvalues, Eigenvectors

---

Let us find the eigenvalues and eigenvectors of the  $2 \times 2$  matrix

$$A = \begin{bmatrix} 4 & 2 \\ 1 & 3 \end{bmatrix}. \quad (4.28)$$

**Step 1: Characteristic Polynomial.** From our definition of the eigenvector  $x \neq \mathbf{0}$  and eigenvalue  $\lambda$  of  $A$ , there will be a vector such that  $Ax = \lambda x$ , i.e.,  $(A - \lambda I)x = \mathbf{0}$ . Since  $x \neq \mathbf{0}$ , this requires that the kernel (null space) of  $A - \lambda I$  contains more elements than just  $\mathbf{0}$ . This means that  $A - \lambda I$  is not invertible and therefore  $\det(A - \lambda I) = 0$ . Hence, we need to compute the roots of the characteristic polynomial (4.22a) to find the eigenvalues.

# Computing Eigenvalues, Eigenvectors

---

**Step 2: Eigenvalues.** The characteristic polynomial is

$$p_A(\lambda) = \det(A - \lambda I) \quad (4.29a)$$

$$= \det \left( \begin{bmatrix} 4 & 2 \\ 1 & 3 \end{bmatrix} - \begin{bmatrix} \lambda & 0 \\ 0 & \lambda \end{bmatrix} \right) = \begin{vmatrix} 4 - \lambda & 2 \\ 1 & 3 - \lambda \end{vmatrix} \quad (4.29b)$$

$$= (4 - \lambda)(3 - \lambda) - 2 \cdot 1. \quad (4.29c)$$

We factorize the characteristic polynomial and obtain

$$p(\lambda) = (4 - \lambda)(3 - \lambda) - 2 \cdot 1 = 10 - 7\lambda + \lambda^2 = (2 - \lambda)(5 - \lambda) \quad (4.30)$$

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

# Computing Eigenvalues, Eigenvectors

---

**Step 3: Eigenvectors and Eigenspaces.** We find the eigenvectors that correspond to these eigenvalues by looking at vectors  $x$  such that

$$\begin{bmatrix} 4 - \lambda & 2 \\ 1 & 3 - \lambda \end{bmatrix} x = \mathbf{0}. \quad (4.31)$$

For  $\lambda = 5$  we obtain

$$\begin{bmatrix} 4 - 5 & 2 \\ 1 & 3 - 5 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} -1 & 2 \\ 1 & -2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \mathbf{0}. \quad (4.32)$$

We solve this homogeneous system and obtain a solution space

$$E_5 = \text{span} \left[ \begin{bmatrix} 2 \\ 1 \end{bmatrix} \right]. \quad (4.33)$$

# Computing Eigenvalues, Eigenvectors

---

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

$$\begin{bmatrix} 4 - 2 & 2 \\ 1 & 3 - 2 \end{bmatrix} \mathbf{x} = \begin{bmatrix} 2 & 2 \\ 1 & 1 \end{bmatrix} \mathbf{x} = \mathbf{0}. \quad (4.34)$$

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

$$E_2 = \text{span}\left[\begin{bmatrix} 1 \\ -1 \end{bmatrix}\right]. \quad (4.35)$$

# Eigenvalues, Eigenvectors

---

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

$$\det(A) = \prod_{i=1}^n \lambda_i , \quad (4.42)$$

*where  $\lambda_i \in \mathbb{C}$  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.,*

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

*where  $\lambda_i \in \mathbb{C}$  are (possibly repeated) eigenvalues of  $A$ .*

# Singular Value Decomposition

---

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.

# Singular Value Decomposition

---

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

$$\begin{matrix} n \\ \approx \\ m \end{matrix} \boxed{\mathbf{A}} = \begin{matrix} m \\ \approx \\ m \end{matrix} \boxed{\mathbf{U}} \begin{matrix} n \\ \approx \\ m \end{matrix} \boxed{\Sigma} \begin{matrix} n \\ \approx \\ m \end{matrix} \boxed{\mathbf{V}^\top} \approx \quad (4.64)$$

*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,  $\Sigma$  is an  $m \times n$  matrix with  $\Sigma_{ii} = \sigma_i \geq 0$  and  $\Sigma_{ij} = 0$ ,  $i \neq j$ .*

The diagonal entries  $\sigma_i$ ,  $i = 1, \dots, r$ , of  $\Sigma$  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.,  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r \geq 0$ .

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

# Singular Value Decomposition

---

if  $m > n$

$$\Sigma = \begin{bmatrix} \sigma_1 & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & \sigma_n \\ 0 & \dots & 0 \\ \vdots & & \vdots \\ 0 & \dots & 0 \end{bmatrix}$$

If  $m < n$

$$\Sigma = \begin{bmatrix} \sigma_1 & 0 & 0 & 0 & \dots & 0 \\ 0 & \ddots & 0 & \vdots & & \vdots \\ 0 & 0 & \sigma_m & 0 & \dots & 0 \end{bmatrix}$$

# Matrix Approximation

---

We considered the SVD as a way to factorize  $\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^\top \in \mathbb{R}^{m \times n}$  into the product of three matrices, where  $\mathbf{U} \in \mathbb{R}^{m \times m}$  and  $\mathbf{V} \in \mathbb{R}^{n \times n}$  are orthogonal and  $\Sigma$  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  $\mathbf{A}$  as a sum of simpler (low-rank) matrices  $\mathbf{A}_i$ , which lends itself to a matrix approximation scheme that is cheaper to compute than the full SVD.

# Matrix Approximation

---

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

$$\mathbf{A}_i := \mathbf{u}_i \mathbf{v}_i^\top, \quad (4.90)$$

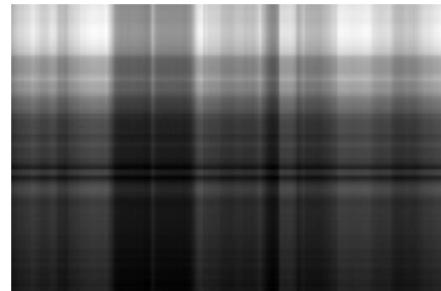
which is formed by the outer product of the  $i$ th orthogonal column vector of  $\mathbf{U}$  and  $\mathbf{V}$ . Figure 4.11 shows an image of Stonehenge, which can be represented by a matrix  $\mathbf{A} \in \mathbb{R}^{1432 \times 1910}$ , and some outer products  $\mathbf{A}_i$ , as defined in (4.90).

# Matrix Approximation

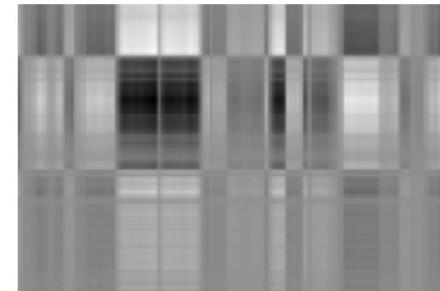
---



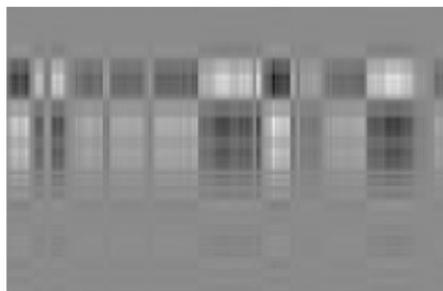
(a) Original image  $\mathbf{A}$ .



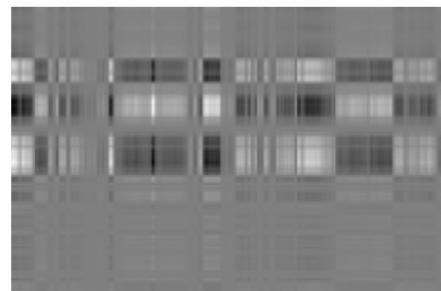
(b)  $\mathbf{A}_1$ ,  $\sigma_1 \approx 228,052$ .



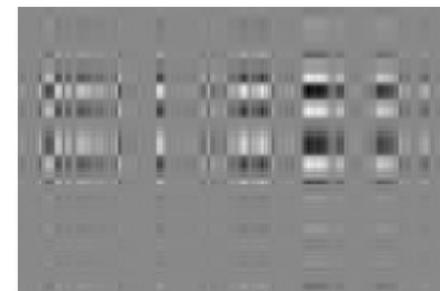
(c)  $\mathbf{A}_2$ ,  $\sigma_2 \approx 40,647$ .



(d)  $\mathbf{A}_3$ ,  $\sigma_3 \approx 26,125$ .



(e)  $\mathbf{A}_4$ ,  $\sigma_4 \approx 20,232$ .



(f)  $\mathbf{A}_5$ ,  $\sigma_5 \approx 15,436$ .

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.

# Matrix Approximation

---

Given two vectors of size  $m \times 1$  and  $n \times 1$  respectively

$$\mathbf{u} = \begin{bmatrix} u_1 \\ u_2 \\ \vdots \\ u_m \end{bmatrix}, \quad \mathbf{v} = \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_n \end{bmatrix}$$

their outer product, denoted  $\mathbf{u} \otimes \mathbf{v}$ , is defined as the  $m \times n$  matrix  $\mathbf{A}$  obtained by multiplying each element of  $\mathbf{u}$  by each element of  $\mathbf{v}$ :<sup>[1]</sup>

$$\mathbf{u} \otimes \mathbf{v} = \mathbf{A} = \begin{bmatrix} u_1 v_1 & u_1 v_2 & \dots & u_1 v_n \\ u_2 v_1 & u_2 v_2 & \dots & u_2 v_n \\ \vdots & \vdots & \ddots & \vdots \\ u_m v_1 & u_m v_2 & \dots & u_m v_n \end{bmatrix}$$

Or in index notation:

$$(\mathbf{u} \otimes \mathbf{v})_{ij} = u_i v_j$$

# Matrix Approximation

---

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

$$\mathbf{A} = \sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^\top = \sum_{i=1}^r \sigma_i \mathbf{A}_i, \quad (4.91)$$

where the outer-product matrices  $\mathbf{A}_i$  are weighted by the  $i$ th singular value  $\sigma_i$ . We can see why (4.91) holds: The diagonal structure of the singular value matrix  $\Sigma$  multiplies only matching left- and right-singular vectors  $\mathbf{u}_i \mathbf{v}_i^\top$  and scales them by the corresponding singular value  $\sigma_i$ . All terms  $\Sigma_{ij} \mathbf{u}_i \mathbf{v}_j^\top$  vanish for  $i \neq j$  because  $\Sigma$  is a diagonal matrix. Any terms  $i > r$  vanish because the corresponding singular values are 0.

# Matrix Approximation

---

In (4.90), 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 (4.91). 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*

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

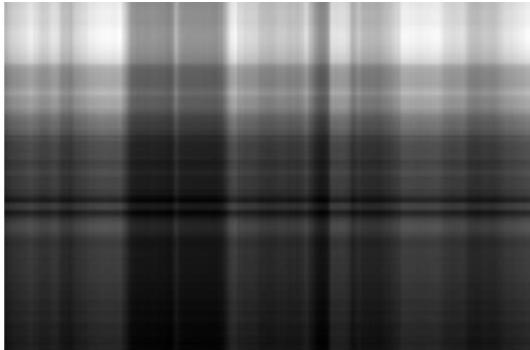
of  $\mathbf{A}$  with  $\text{rk}(\widehat{\mathbf{A}}(k)) = k$ . Figure 4.12 shows low-rank approximations  $\widehat{\mathbf{A}}(k)$  of an original image  $\mathbf{A}$  of Stonehenge. The shape of the rocks becomes increasingly visible and clearly recognizable in the rank-5 approximation. While the 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.

# Matrix Approximation

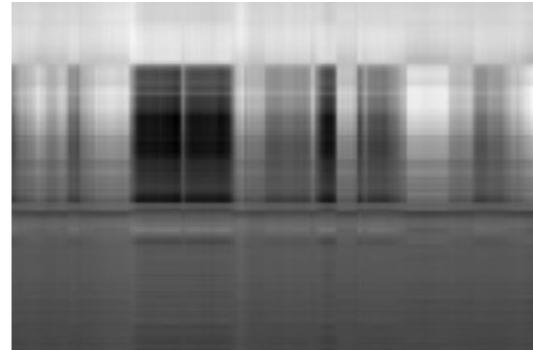
---



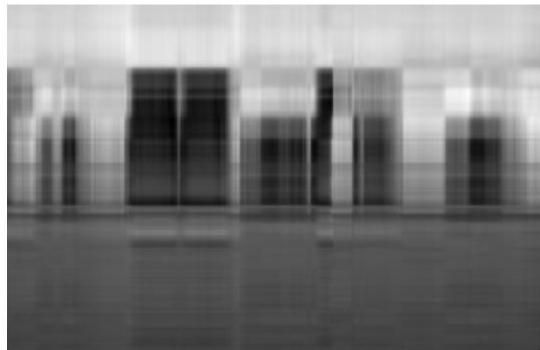
(a) Original image  $\mathbf{A}$ .



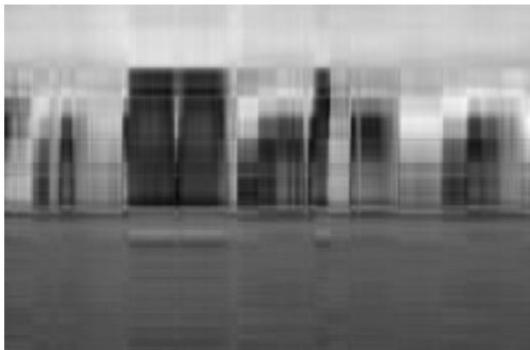
(b) Rank-1 approximation  $\hat{\mathbf{A}}(1)$ .



(c) Rank-2 approximation  $\hat{\mathbf{A}}(2)$ .



(d) Rank-3 approximation  $\hat{\mathbf{A}}(3)$ .



(e) Rank-4 approximation  $\hat{\mathbf{A}}(4)$ .

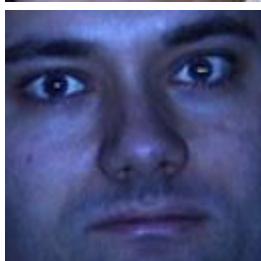


(f) Rank-5 approximation  $\hat{\mathbf{A}}(5)$ .

**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  $\hat{\mathbf{A}}(k) = \sum_{i=1}^k \sigma_i \mathbf{A}_i$ .

# Matrix Approximation

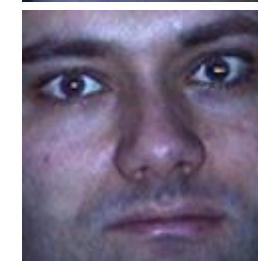
$$W_A = U_{W,A} \Sigma_{W,A} V_{W,A}^T$$



$$W_A = U_{W,A} V_{W,A}^T$$



$$W_A = U_{W,A} \Sigma_{M,A} V_{W,A}^T$$



$$M_A = U_{M,A} \Sigma_{M,A} V_{M,A}^T$$

(a)

$$M_A = U_{M,A} V_{M,A}^T$$

(b)

$$M_A = U_{M,A} \Sigma_{W,A} V_{M,A}^T$$

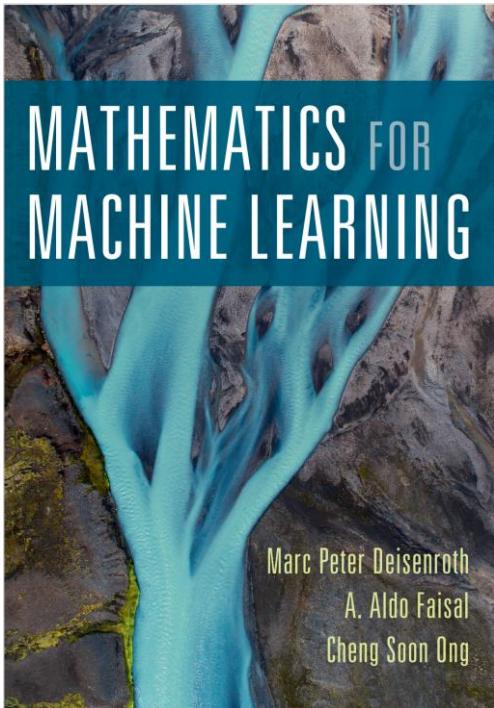
(c)

From top to bottom: (a) Two color face images from the CMU-PIE database. (b) The reconstructed images of each images in (a) when let all SVs equal to 1 and (c) The reconstructed images by swaping singular value matrix between 2 face images.

# References

---

[Mathematics for Machine Learning \(mml-book.github.io\)](https://mml-book.github.io)



[math-deep.pdf \(upenn.edu\)](http://math-deep.pdf.upenn.edu)

Algebra, Topology, Differential Calculus, and  
Optimization Theory  
For Computer Science and Machine Learning

Jean Gallier and Jocelyn Quaintance  
Department of Computer and Information Science  
University of Pennsylvania  
Philadelphia, PA 19104, USA  
e-mail: [jean@seas.upenn.edu](mailto:jean@seas.upenn.edu)

© Jean Gallier  
March 18, 2022