# Mathematics Refresher

Daniel Wehner, M.Sc.

SAP SE / DHBW Mannheim

Summer term 2025

Find all slides on `GitHub` (DaWe1992/Applied_ML_Fundamentals)
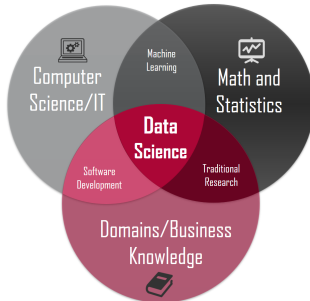
# Agenda for this Unit

1. Introduction

2. Linear Algebra

3. Probability Theory and Statistics

4. Wrap-Up

**Section:**

**Introduction**

Introduction
Math is important!

Introduction
Linear Algebra
Probability Theory and Statistics
Wrap-Up

Introduction
Math is important!

# Introduction

**Mathematics play a major role in data science and machine learning!**
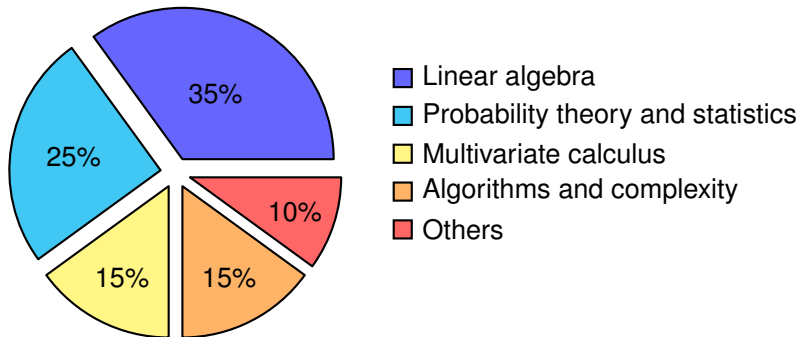


You will need it to understand:

- **Statistical** machine learning
- How **optimization** is used in learning and empirical risk minimization
- How linear algebra, calculus, and statistics are used to make learning and inference more efficient

Introduction
Linear Algebra
Probability Theory and Statistics
Wrap-Up

Introduction
Math is important!

# Math is important!

Rough importance of mathematical disciplines in data science and machine learning:



- Linear algebra
- Probability theory and statistics
- Multivariate calculus
- Algorithms and complexity
- Others

**Section:**

**Linear Algebra**

Vectors and Vector Operations
Matrix Operations
Determinants and Inverses
Eigenvectors and Eigenvalues
Symmetric Matrices and Definiteness

Introduction
Linear Algebra
Probability Theory and Statistics
Wrap-Up

Vectors and Vector Operations
Matrix Operations
Determinants and Inverses
Eigenvectors and Eigenvalues
Symmetric Matrices and Definiteness

## What is a Vector?

General:

$$\boldsymbol{x} = \begin{pmatrix} x_1 \\ \vdots \\ x_D \end{pmatrix} \in \mathbb{R}^D$$

Example: ($D = 2$)

$$\boldsymbol{v} = \begin{pmatrix} 3 \\ 1 \end{pmatrix} \qquad \boldsymbol{w} = \begin{pmatrix} 2 \\ 4 \end{pmatrix}$$

Introduction
Linear Algebra
Probability Theory and Statistics
Wrap-Up

Vectors and Vector Operations
Matrix Operations
Determinants and Inverses
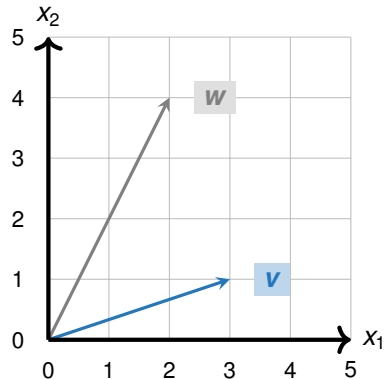Eigenvectors and Eigenvalues
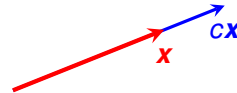Symmetric Matrices and Definiteness

## Multiplication of Vectors by Scalars

General: Let $c \in \mathbb{R}$ and $\boldsymbol{x} \in \mathbb{R}^D$:

$$c\boldsymbol{x} = c \begin{pmatrix} x_1 \\ \vdots \\ x_D \end{pmatrix} = \begin{pmatrix} cx_1 \\ \vdots \\ cx_D \end{pmatrix}$$

Example: $(D = 2)$

$$2\boldsymbol{v} = 2 \begin{pmatrix} 3 \\ 1 \end{pmatrix} = \begin{pmatrix} 6 \\ 2 \end{pmatrix}$$

Introduction
**Linear Algebra**
Probability Theory and Statistics
Wrap-Up

**Vectors and Vector Operations**
Matrix Operations
Determinants and Inverses
Eigenvectors and Eigenvalues
Symmetric Matrices and Definiteness

## Addition of Vectors

General: Let $\boldsymbol{x}, \boldsymbol{y} \in \mathbb{R}^D$

$$\boldsymbol{x} + \boldsymbol{y} = \begin{pmatrix} x_1 \\ \vdots \\ x_D \end{pmatrix} + \begin{pmatrix} y_1 \\ \vdots \\ y_D \end{pmatrix} = \begin{pmatrix} x_1 + y_1 \\ \vdots \\ x_D + y_D \end{pmatrix}$$

Example: $(D = 2)$

$$\boldsymbol{v} + \boldsymbol{w} = \begin{pmatrix} 3 \\ 1 \end{pmatrix} + \begin{pmatrix} 2 \\ 4 \end{pmatrix} = \begin{pmatrix} 5 \\ 5 \end{pmatrix}$$

Introduction
Linear Algebra
Probability Theory and Statistics
Wrap-Up

Vectors and Vector Operations
Matrix Operations
Determinants and Inverses
Eigenvectors and Eigenvalues
Symmetric Matrices and Definiteness

## Linear Combination of Vectors and Span

- Let $c_1, \ldots, c_N \in \mathbb{R}$ and $\boldsymbol{x}^1, \ldots, \boldsymbol{x}^N \in \mathbb{R}^D$

- A **linear combination** of these vectors is given by $\boldsymbol{u} \in \mathbb{R}^D$:

$$\boldsymbol{u} := \sum_{n=1}^{N} c_n \boldsymbol{x}^n \tag{1}$$

- The **span** (German: *lineare Hülle*) of $\boldsymbol{x}^1, \ldots, \boldsymbol{x}^N \in \mathbb{R}^D$ is defined by:

$$\text{span}(\boldsymbol{x}^1, \ldots, \boldsymbol{x}^N) := \left\{ \boldsymbol{u} \in \mathbb{R}^D : \exists c_1, \ldots, c_N : \boldsymbol{u} = \sum_{n=1}^{N} c_n \boldsymbol{x}^n \right\} \tag{2}$$

Introduction
Linear Algebra
Probability Theory and Statistics
Wrap-Up

Vectors and Vector Operations
Matrix Operations
Determinants and Inverses
Eigenvectors and Eigenvalues
Symmetric Matrices and Definiteness

# Linear Combination of Vectors and Span (Ctd.)

Introduction
**Linear Algebra**
Probability Theory and Statistics
Wrap-Up

**Vectors and Vector Operations**
Matrix Operations
Determinants and Inverses
Eigenvectors and Eigenvalues
Symmetric Matrices and Definiteness

## Vector Transpose, Inner and Outer Product

- Let $\boldsymbol{x}, \boldsymbol{y} \in \mathbb{R}^D$ be given

- **Transposition:**

$$\boldsymbol{x} = \begin{pmatrix} x_1 \\ \vdots \\ x_D \end{pmatrix} \qquad \boldsymbol{x}^\top = \begin{pmatrix} x_1 & \dots & x_D \end{pmatrix} \tag{3}$$

- **Inner product** (also referred to as **dot product** or **scalar product**):

$$\boldsymbol{x}^\top \boldsymbol{y} \equiv \langle \boldsymbol{x}, \boldsymbol{y} \rangle := \begin{pmatrix} x_1 & \dots & x_D \end{pmatrix} \begin{pmatrix} y_1 \\ \vdots \\ y_D \end{pmatrix} = \sum_{d=1}^{D} x_d y_d \tag{4}$$

Introduction
**Linear Algebra**
Probability Theory and Statistics
Wrap-Up

Vectors and Vector Operations
Matrix Operations
Determinants and Inverses
Eigenvectors and Eigenvalues
Symmetric Matrices and Definiteness

## Vector Transpose, Inner and Outer Product (Ctd.)

- **Outer product:**

$$\boldsymbol{x}\boldsymbol{y}^\top = \begin{pmatrix} x_1 \\ \vdots \\ x_D \end{pmatrix} \begin{pmatrix} y_1 & \dots & y_D \end{pmatrix} = \begin{pmatrix} x_1y_1 & x_1y_2 & \dots & x_1y_D \\ x_2y_1 & x_2y_2 & \dots & x_2y_D \\ \vdots & \vdots & \ddots & \vdots \\ x_Dy_1 & x_Dy_2 & \dots & x_Dy_D \end{pmatrix} \tag{5}$$

**Remember:** The inner product yields a scalar value; The result of an outer product is a matrix!

Introduction
Linear Algebra
Probability Theory and Statistics
Wrap-Up

Vectors and Vector Operations
Matrix Operations
Determinants and Inverses
Eigenvectors and Eigenvalues
Symmetric Matrices and Definiteness

# Example: Vector Transpose, Inner and Outer Product

- Let $\boldsymbol{v} = \begin{pmatrix} 3 \\ 1 \end{pmatrix} \in \mathbb{R}^2$ and $\boldsymbol{w} = \begin{pmatrix} 2 \\ 4 \end{pmatrix} \in \mathbb{R}^2$

- **Transposition:**

$$\boldsymbol{v}^\top = \begin{pmatrix} 3 & 1 \end{pmatrix}$$

- **Inner product:**

$$\boldsymbol{v}^\top \boldsymbol{w} = \begin{pmatrix} 3 & 1 \end{pmatrix} \begin{pmatrix} 2 \\ 4 \end{pmatrix} = 3 \cdot 2 + 1 \cdot 4 = 10$$

Introduction
Linear Algebra
Probability Theory and Statistics
Wrap-Up

Vectors and Vector Operations
Matrix Operations
Determinants and Inverses
Eigenvectors and Eigenvalues
Symmetric Matrices and Definiteness

# Example: Vector Transpose, Inner and Outer Product (Ctd.)

- **Outer product:**

$$\boldsymbol{v}\boldsymbol{w}^\top = \begin{pmatrix} 3 \\ 1 \end{pmatrix} \begin{pmatrix} 2 & 4 \end{pmatrix} = \begin{pmatrix} 6 & 12 \\ 2 & 4 \end{pmatrix}$$

Introduction
**Linear Algebra**
Probability Theory and Statistics
Wrap-Up

Vectors and Vector Operations
Matrix Operations
Determinants and Inverses
Eigenvectors and Eigenvalues
Symmetric Matrices and Definiteness

## Length of a Vector and Norms

- Length of a vector **(EUCLIDean norm)**: Let $\boldsymbol{x}, \boldsymbol{y} \in \mathbb{R}^D$ and $c \in \mathbb{R}$

$$\|\boldsymbol{x}\| := \sqrt{\boldsymbol{x}^\top \boldsymbol{x}} \tag{6}$$

$$\|c\boldsymbol{x}\| = |c| \cdot \|\boldsymbol{x}\| \tag{7}$$

$$\|\boldsymbol{x} + \boldsymbol{y}\| \leqslant \|\boldsymbol{x}\| + \|\boldsymbol{y}\| \tag{8}$$

- Example: Let $\boldsymbol{v} = \begin{pmatrix} 3 \\ 1 \end{pmatrix} \in \mathbb{R}^2$

$$\|\boldsymbol{v}\| = \sqrt{3^2 + 1^2} = \sqrt{10}$$

Introduction
Linear Algebra
Probability Theory and Statistics
Wrap-Up

Vectors and Vector Operations
Matrix Operations
Determinants and Inverses
Eigenvectors and Eigenvalues
Symmetric Matrices and Definiteness

## Angle between Vectors

- The **angle** between two vectors $\boldsymbol{x}, \boldsymbol{y} \in \mathbb{R}^D$ is given by:

$$\cos \measuredangle(\boldsymbol{x}, \boldsymbol{y}) = \frac{\boldsymbol{x}^\top \boldsymbol{y}}{\|\boldsymbol{x}\| \cdot \|\boldsymbol{y}\|} = \frac{\sum_{d=1}^{D} x_d y_d}{\sqrt{\sum_{d=1}^{D} x_d^2} \cdot \sqrt{\sum_{d=1}^{D} y_d^2}} \qquad (9)$$

$$\cos \measuredangle(\boldsymbol{v}, \boldsymbol{w}) = \frac{\boldsymbol{v}^\top \boldsymbol{w}}{\|\boldsymbol{v}\| \cdot \|\boldsymbol{w}\|} = \frac{10}{\sqrt{10} \cdot \sqrt{20}} \approx 0.71$$

- Inner product: $\boldsymbol{x}^\top \boldsymbol{y} = \|\boldsymbol{x}\| \cdot \|\boldsymbol{y}\| \cdot \cos \measuredangle(\boldsymbol{x}, \boldsymbol{y})$

Introduction
Linear Algebra
Probability Theory and Statistics
Wrap-Up

Vectors and Vector Operations
Matrix Operations
Determinants and Inverses
Eigenvectors and Eigenvalues
Symmetric Matrices and Definiteness

## (Orthogonal) Projection of Vectors

- How is the projection of $\boldsymbol{y}$ onto $\boldsymbol{x}$ defined?

- Formally, we have:

$$
\begin{aligned}
p &= \|\boldsymbol{y}\| \cos \measuredangle(\boldsymbol{x}, \boldsymbol{y}) \\
&= \|\boldsymbol{y}\| \frac{\boldsymbol{x}^\top \boldsymbol{y}}{\|\boldsymbol{x}\| \cdot \|\boldsymbol{y}\|} \\
&= \frac{\boldsymbol{x}^\top \boldsymbol{y}}{\|\boldsymbol{x}\|}
\end{aligned}
\tag{10}
$$

- Note that $p$ is **not a vector!**

Introduction
Linear Algebra
Probability Theory and Statistics
Wrap-Up

Vectors and Vector Operations
Matrix Operations
Determinants and Inverses
Eigenvectors and Eigenvalues
Symmetric Matrices and Definiteness

## What is a Matrix?

- **General case:**

$$\boldsymbol{X} = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1M} \\ x_{21} & x_{22} & \dots & x_{2M} \\ \vdots & \vdots & \ddots & \vdots \\ x_{N1} & x_{N2} & \dots & x_{NM} \end{pmatrix} \in \mathbb{R}^{N \times M}$$

- $x_{nm}$ is the entry in row $n$ and column $m$

**Remember: Z**eilen **z**uerst, **S**palten **s**päter (German)

Introduction
Linear Algebra
Probability Theory and Statistics
Wrap-Up

Vectors and Vector Operations
Matrix Operations
Determinants and Inverses
Eigenvectors and Eigenvalues
Symmetric Matrices and Definiteness

# Example: Matrices

$$A = \begin{pmatrix} 3 & 4 & 5 \\ 1 & 0 & 1 \end{pmatrix} \in \mathbb{R}^{2\times3}$$

$A$ has three rows, but only two columns.

$$B = \begin{pmatrix} 10 & 1 \\ 11 & 2 \end{pmatrix} \in \mathbb{R}^{2\times2}$$

$B$ is a **square matrix** as it has the same number of rows and columns.

Square matrices play a special role in mathematics, e.g. matrix inversion and determinants are only defined for square matrices.

$$C = \begin{pmatrix} 3 & 0 \\ 0 & 7 \end{pmatrix} \in \mathbb{R}^{2\times2}$$

$C$ is a square matrix, but also a **diagonal matrix** because $c_{ij} = 0$ for $i \neq j$. We often write $C := \mathrm{diag}(3, 7)$.

Introduction
Linear Algebra
Probability Theory and Statistics
Wrap-Up

Vectors and Vector Operations
Matrix Operations
Determinants and Inverses
Eigenvectors and Eigenvalues
Symmetric Matrices and Definiteness

## Special Matrices

**Identity matrix:**

$$I_N := \begin{pmatrix} 1 & 0 & 0 & \dots & 0 \\ 0 & 1 & 0 & \dots & 0 \\ 0 & 0 & 1 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & 1 \end{pmatrix} \in \mathbb{R}^{N \times N}$$

**Zero matrix:**

$$\mathbf{0} := \begin{pmatrix} 0 & 0 & 0 & \dots & 0 \\ 0 & 0 & 0 & \dots & 0 \\ 0 & 0 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & 0 \end{pmatrix} \in \mathbb{R}^{N \times N}$$

Introduction
Linear Algebra
Probability Theory and Statistics
Wrap-Up

Vectors and Vector Operations
Matrix Operations
Determinants and Inverses
Eigenvectors and Eigenvalues
Symmetric Matrices and Definiteness

## Matrix Transpose

- **Matrix transposition:**

$$\boldsymbol{X}^\top = \begin{pmatrix} x_{11} & x_{21} & \dots & x_{N1} \\ x_{12} & x_{22} & \dots & x_{N2} \\ \vdots & \vdots & \ddots & \vdots \\ x_{1M} & x_{2M} & \dots & x_{NM} \end{pmatrix} \in \mathbb{R}^{M \times N} \tag{11}$$

- **Please note:** $\boldsymbol{X} \in \mathbb{R}^{N \times M}$, but $\boldsymbol{X}^\top \in \mathbb{R}^{M \times N}$

Introduction
**Linear Algebra**
Probability Theory and Statistics
Wrap-Up

Vectors and Vector Operations
Matrix Operations
Determinants and Inverses
Eigenvectors and Eigenvalues
Symmetric Matrices and Definiteness

## Matrix Addition

- **Addition of matrices:** Let $X$, $Y \in \mathbb{R}^{N \times M}$

$$X + Y = \begin{pmatrix} x_{11} & \ldots & x_{1M} \\ \vdots & \ddots & \vdots \\ x_{N1} & \ldots & x_{NM} \end{pmatrix} + \begin{pmatrix} y_{11} & \ldots & y_{1M} \\ \vdots & \ddots & \vdots \\ y_{N1} & \ldots & y_{NM} \end{pmatrix}$$

$$= \begin{pmatrix} x_{11} + y_{11} & \ldots & x_{1M} + y_{1M} \\ \vdots & \ddots & \vdots \\ x_{N1} + y_{N1} & \ldots & x_{NM} + y_{NM} \end{pmatrix} \in \mathbb{R}^{N \times M} \qquad (12)$$

- **Please note:** $X$ and $Y$ must be of the same size!

Introduction
Linear Algebra
Probability Theory and Statistics
Wrap-Up

Vectors and Vector Operations
Matrix Operations
Determinants and Inverses
Eigenvectors and Eigenvalues
Symmetric Matrices and Definiteness

## Multiplication of Matrices by Scalars

- **Multiplication by scalars:** Let $\boldsymbol{X} \in \mathbb{R}^{N \times M}$ and $c \in \mathbb{R}$

$$c\boldsymbol{X} = c \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1M} \\ x_{21} & x_{22} & \dots & x_{2M} \\ \vdots & \vdots & \ddots & \vdots \\ x_{N1} & x_{N2} & \dots & x_{NM} \end{pmatrix} = \begin{pmatrix} cx_{11} & cx_{12} & \dots & cx_{1M} \\ cx_{21} & cx_{22} & \dots & cx_{2M} \\ \vdots & \vdots & \ddots & \vdots \\ cx_{N1} & cx_{N2} & \dots & cx_{NM} \end{pmatrix} \in \mathbb{R}^{N \times M} \quad (13)$$

- This is defined for all matrices

Introduction
**Linear Algebra**
Probability Theory and Statistics
Wrap-Up

Vectors and Vector Operations
Matrix Operations
Determinants and Inverses
Eigenvectors and Eigenvalues
Symmetric Matrices and Definiteness

## Multiplication of Matrices by Vectors

- **Matrix-vector multiplication:** Let $\boldsymbol{X} \in \mathbb{R}^{N \times M}$ and $\boldsymbol{y} \in \mathbb{R}^M$

$$\boldsymbol{Xy} = \begin{pmatrix} x_{11} & x_{12} & \ldots & x_{1M} \\ x_{21} & x_{22} & \ldots & x_{2M} \\ \vdots & \vdots & \ddots & \vdots \\ x_{N1} & x_{N2} & \ldots & x_{NM} \end{pmatrix} \begin{pmatrix} y_1 \\ v_2 \\ \vdots \\ y_M \end{pmatrix} = \begin{pmatrix} \sum_{m=1}^{M} x_{1m} y_m \\ \sum_{m=1}^{M} x_{2m} y_m \\ \vdots \\ \sum_{m=1}^{M} x_{Nm} y_m \end{pmatrix} \in \mathbb{R}^N \quad (14)$$

- **Please note:** The number of columns of $\boldsymbol{X}$ and the number of rows of $\boldsymbol{y}$ must be equal in order for the matrix-vector product to exist!

- The order is important: $\boldsymbol{Xy}$ is defined, but $\boldsymbol{yX}$ is not, if $n > 1$

Introduction
Linear Algebra
Probability Theory and Statistics
Wrap-Up

Vectors and Vector Operations
Matrix Operations
Determinants and Inverses
Eigenvectors and Eigenvalues
Symmetric Matrices and Definiteness

## Matrix Multiplication (Ctd.)

- **Matrix-matrix multiplication:** Let $X \in \mathbb{R}^{L \times M}$ and $Y \in \mathbb{R}^{M \times N}$

$$XY = \begin{pmatrix} x_{11} & \dots & x_{1M} \\ \vdots & \ddots & \vdots \\ x_{L1} & \dots & x_{LM} \end{pmatrix} \begin{pmatrix} y_{11} & \dots & y_{1N} \\ \vdots & \ddots & \vdots \\ y_{M1} & \dots & y_{MN} \end{pmatrix} = \begin{pmatrix} z_{11} & \dots & z_{1N} \\ \vdots & \ddots & \vdots \\ z_{L1} & \dots & z_{LN} \end{pmatrix} \quad (15)$$

where:

$$z_{\ell n} = \sum_{m=1}^{M} x_{\ell m} y_{nk} \quad (16)$$

- **Please note:** The number of columns of $X$ and the number of rows of $Y$ must match!

Introduction
**Linear Algebra**
Probability Theory and Statistics
Wrap-Up

Vectors and Vector Operations
Matrix Operations
Determinants and Inverses
Eigenvectors and Eigenvalues
Symmetric Matrices and Definiteness

## Determinants of Square Matrices

- **Determinants are defined for square matrices only!**

- The **determinant** of a $(2 \times 2)$-matrix is given by:

$$\det \begin{pmatrix} a & b \\ c & d \end{pmatrix} := ad - bc \tag{17}$$

- The determinant of a $(3 \times 3)$-matrix is given by **(rule of SARRUS)**:

$$\det \begin{pmatrix} a & b & c \\ d & e & f \\ g & h & i \end{pmatrix} := aei + bfg + cdh - gec - hfa - idb \tag{18}$$

Introduction
**Linear Algebra**
Probability Theory and Statistics
Wrap-Up

Vectors and Vector Operations
Matrix Operations
**Determinants and Inverses**
Eigenvectors and Eigenvalues
Symmetric Matrices and Definiteness

# LAPLACE Expansion

Use the **LAPLACE expansion** for $(N \times N)$-matrices if $N > 3$:

**LAPLACE expansion:** Let $\boldsymbol{X} \in \mathbb{R}^{N \times N}$ be given. Then

$$\det(\boldsymbol{X}) = \sum_{n=1}^{N} x_{nm} \cdot (-1)^{n+m} \cdot \det(\boldsymbol{X}_{nm}), \qquad (19)$$

where $\boldsymbol{X}_{nm}$ is the matrix obtained by removing row $n$ and column $m$ from $\boldsymbol{X}$.

Introduction
Linear Algebra
Probability Theory and Statistics
Wrap-Up

Vectors and Vector Operations
Matrix Operations
**Determinants and Inverses**
Eigenvectors and Eigenvalues
Symmetric Matrices and Definiteness

## Matrix Inversion

- **Matrix inversion is defined for square matrices only!**

- $\boldsymbol{X} \in \mathbb{R}^{N \times N}$ multiplied by its inverse $\boldsymbol{X}^{-1} \in \mathbb{R}^{N \times N}$ gives the identity matrix $\boldsymbol{I}_N$:

$$\boldsymbol{X}\boldsymbol{X}^{-1} = \boldsymbol{I}_N \tag{20}$$

- Also, the order is not important, i.e.:

$$\boldsymbol{X}^{-1}\boldsymbol{X} = \boldsymbol{I}_N \tag{21}$$

- We call $\boldsymbol{X}$ **non-singular** or **invertible**, if $\boldsymbol{X}^{-1}$ exists

Introduction
Linear Algebra
Probability Theory and Statistics
Wrap-Up

Vectors and Vector Operations
Matrix Operations
**Determinants and Inverses**
Eigenvectors and Eigenvalues
Symmetric Matrices and Definiteness

## Matrix Inversion (Ctd.)

Let $\boldsymbol{X} \in \mathbb{R}^{N \times N}$ be a square matrix. The following statements are equivalent:

$\boldsymbol{X}$ is invertible $\iff \boldsymbol{X}$ is non-singular

$\iff \det(\boldsymbol{X}) \neq 0$

$\iff \boldsymbol{X}$ has rank $N$ (full rank)

$\iff \boldsymbol{X}$ does not have eigenvalue 0

$\iff$ The **reduced row echelon form** of $\boldsymbol{X}$ is the identity matrix $\boldsymbol{I}_N$

Introduction
Linear Algebra
Probability Theory and Statistics
Wrap-Up

Vectors and Vector Operations
Matrix Operations
**Determinants and Inverses**
Eigenvectors and Eigenvalues
Symmetric Matrices and Definiteness

## Matrix Inversion (Ctd.)

- The inverse of a matrix can be computed using the **GAUSS-JORDAN algorithm**

- **Special case:** Do not use the GAUSS-JORDAN algorithm for $(2 \times 2)$-matrices!

You can be more efficient using the identity

$$\boldsymbol{X}\mathrm{adj}(\boldsymbol{X}) = \det(\boldsymbol{X})\boldsymbol{I}_N \quad \overset{\det(\boldsymbol{X}) \neq 0}{\Longleftrightarrow} \quad \boldsymbol{X}\overbrace{\frac{1}{\det(\boldsymbol{X})}\mathrm{adj}(\boldsymbol{X})}^{=: \ \boldsymbol{X}^{-1}} = \boldsymbol{I}_N, \tag{22}$$

where $\mathrm{adj}(\boldsymbol{X})$ is the **adjugate matrix** of $\boldsymbol{X}$.

Introduction
Linear Algebra
Probability Theory and Statistics
Wrap-Up

Vectors and Vector Operations
Matrix Operations
Determinants and Inverses
Eigenvectors and Eigenvalues
Symmetric Matrices and Definiteness

## Matrix Inversion (Ctd.)

- In general, $\text{adj}(\boldsymbol{X})$ is hard to compute, but it is easy for $(2 \times 2)$-matrices

- Let

$$\boldsymbol{X} = \begin{pmatrix} a & b \\ c & d \end{pmatrix} \in \mathbb{R}^{2 \times 2}$$

- The adjugate matrix $\text{adj}(\boldsymbol{X})$ of $\boldsymbol{X}$ is then given by:

$$\text{adj}(\boldsymbol{X}) := \begin{pmatrix} d & -b \\ -c & a \end{pmatrix} \in \mathbb{R}^{2 \times 2} \tag{23}$$

Introduction
**Linear Algebra**
Probability Theory and Statistics
Wrap-Up

Vectors and Vector Operations
Matrix Operations
Determinants and Inverses
**Eigenvectors and Eigenvalues**
Symmetric Matrices and Definiteness

# Eigenvectors and Eigenvalues

- Let $\boldsymbol{X} \in \mathbb{R}^{N \times N}$ be a square matrix

- Some vectors $\boldsymbol{v} \in \mathbb{R}^N$ only change their length (but not their direction) when multiplied by $\boldsymbol{X}$

- Such vectors are called **eigenvectors** of $\boldsymbol{X}$ and the scaling factors are known as **eigenvalues** of $\boldsymbol{X}$

**Eigenvectors and eigenvalues satisfy the equation:**

$$\boldsymbol{X}\boldsymbol{v} = \lambda \boldsymbol{v} \tag{24}$$

Introduction
Linear Algebra
Probability Theory and Statistics
Wrap-Up

Vectors and Vector Operations
Matrix Operations
Determinants and Inverses
Eigenvectors and Eigenvalues
Symmetric Matrices and Definiteness

## Example: Eigenvectors and Eigenvalues

- Let $\boldsymbol{X} := \begin{pmatrix} 4 & -1 \\ 2 & 1 \end{pmatrix} \in \mathbb{R}^{2 \times 2}$ be given

- We have *(please verify for yourself)*:

$$\begin{pmatrix} 4 & -1 \\ 2 & 1 \end{pmatrix} \begin{pmatrix} 1 \\ 2 \end{pmatrix} = \begin{pmatrix} 2 \\ 4 \end{pmatrix} = 2 \begin{pmatrix} 1 \\ 2 \end{pmatrix}$$

- Thus, $\boldsymbol{v} = \begin{pmatrix} 1 \\ 2 \end{pmatrix}$ is an eigenvector of $\boldsymbol{X}$ and $\lambda = 2$ is the corresponding eigenvalue

Introduction
**Linear Algebra**
Probability Theory and Statistics
Wrap-Up

Vectors and Vector Operations
Matrix Operations
Determinants and Inverses
**Eigenvectors and Eigenvalues**
Symmetric Matrices and Definiteness

# Eigenvectors form a Basis (Eigenbasis)

- Let $\boldsymbol{v}^1, \boldsymbol{v}^2, \ldots, \boldsymbol{v}^N$ be $N$ eigenvectors of $\boldsymbol{X} \in \mathbb{R}^{N \times N}$ with corresponding eigenvalues $\lambda_1, \lambda_2, \ldots, \lambda_N$

- **Theorem (*):**
  - If $\boldsymbol{v}^1, \boldsymbol{v}^2, \ldots, \boldsymbol{v}^N$ correspond to **distinct eigenvalues** $\lambda_1, \lambda_2, \ldots, \lambda_N$, then the system $\left( \boldsymbol{v}^1, \boldsymbol{v}^2, \ldots, \boldsymbol{v}^N \right)$ represents a basis of $\mathbb{R}^N$ **(eigenbasis)**
  - Hence, $\boldsymbol{v}^1, \ldots, \boldsymbol{v}^N$ are **linearly independent** and any vector $\boldsymbol{u} \in \mathbb{R}^N$ can be **uniquely** expressed as a **linear combination** of the eigenvectors of $\boldsymbol{X}$, i.e. $\exists c_1, \ldots, c_N \in \mathbb{R}$ such that

$$\boldsymbol{u} = \sum_{n=1}^{N} c_n \boldsymbol{v}^n \qquad \forall \boldsymbol{u} \in \mathbb{R}^N$$

Introduction
**Linear Algebra**
Probability Theory and Statistics
Wrap-Up

Vectors and Vector Operations
Matrix Operations
Determinants and Inverses
**Eigenvectors and Eigenvalues**
Symmetric Matrices and Definiteness

## How to compute Eigenvalues and Eigenvectors

- Let $\boldsymbol{X} \in \mathbb{R}^{N \times N}$ be a square matrix

- The eigenvalues are the roots (German: *Nullstellen*) of the **characteristic polynomial** defined by

$$\chi_{\boldsymbol{X}}(\lambda) := \det(\lambda \boldsymbol{I}_N - \boldsymbol{X}) \tag{25}$$

- For each eigenvalue $\lambda_n$ we have to solve the **homogeneous system of linear equations**

$$(\boldsymbol{X} - \lambda_n \boldsymbol{I}_N)\boldsymbol{v} = \boldsymbol{0} \tag{26}$$

to obtain the respective eigenvectors (see $\Rightarrow$ here)

Introduction
Linear Algebra
Probability Theory and Statistics
Wrap-Up

Vectors and Vector Operations
Matrix Operations
Determinants and Inverses
Eigenvectors and Eigenvalues
Symmetric Matrices and Definiteness

# Example: Computation of Eigenvalues and Eigenvectors

**Computation of the eigenvalues:**

- Let $A := \begin{pmatrix} 2 & -1 \\ -1 & 2 \end{pmatrix} \in \mathbb{R}^{2 \times 2}$

- The eigenvalues are the roots of (25) which is given by

$$\chi_A := \det \begin{pmatrix} \lambda - 2 & -1 \\ -1 & \lambda - 2 \end{pmatrix} = (\lambda - 2)^2 - 1 = (\lambda - 1)(\lambda - 3)$$

- We directly see: $\lambda_1 = 1, \lambda_2 = 3$

Introduction
Linear Algebra
Probability Theory and Statistics
Wrap-Up

Vectors and Vector Operations
Matrix Operations
Determinants and Inverses
Eigenvectors and Eigenvalues
Symmetric Matrices and Definiteness

## Example: Computation of Eigenvalues and Eigenvectors (Ctd.)

**Computation of the eigenspaces:**

- We start with the eigenvalue $\lambda_1 = 1$ and solve (26):

$$(\boldsymbol{A} - \boldsymbol{I}_2) = \begin{pmatrix} 1 & -1 \\ -1 & 1 \end{pmatrix} \xrightarrow{I+II} \begin{pmatrix} 1 & -1 \\ 0 & 0 \end{pmatrix}$$

- The eigenspace for eigenvalue 1 is therefore given by
$\mathcal{E}(1) = \left\{ t \cdot \left(-1, -1\right)^{\top} : t \in \mathbb{R} \right\}$

- Similarly, we can show that $\mathcal{E}(3) = \left\{ t \cdot \left(1, -1\right)^{\top} : t \in \mathbb{R} \right\}$

Introduction
Linear Algebra
Probability Theory and Statistics
Wrap-Up

Vectors and Vector Operations
Matrix Operations
Determinants and Inverses
Eigenvectors and Eigenvalues
Symmetric Matrices and Definiteness

## Diagonalizable Matrices

- Let $A \in \mathbb{R}^{N \times N}$ be a square matrix

- If the conditions of theorem (*) are met, then we can find a non-singular matrix $S \in \mathbb{R}^{N \times N}$ such that:

$$S^{-1}AS = B \in \mathbb{R}^{N \times N} \tag{27}$$

- The columns of $S$ are given by the eigenvectors of $A$, and $B := \mathrm{diag}(\lambda_1, \ldots, \lambda_N)$ is a diagonal matrix containing the eigenvalues of $A$

- We say $A$ is a **diagonalizable matrix**

- $A$ and $B$ are called **similar matrices**

Introduction
**Linear Algebra**
Probability Theory and Statistics
Wrap-Up

Vectors and Vector Operations
Matrix Operations
Determinants and Inverses
Eigenvectors and Eigenvalues
Symmetric Matrices and Definiteness

## Symmetric Matrices

- A square $(N \times N)$-matrix $\boldsymbol{X}$ is called **symmetric**, if and only if

$$\boldsymbol{X} = \boldsymbol{X}^\top \tag{28}$$

- **Some properties:**
    - The inverse $\boldsymbol{X}^{-1}$ of $\boldsymbol{X}$ is also a symmetric matrix
    - **Eigen-decomposition:** Let $\boldsymbol{X}$ be a symmetric matrix. In this case the conditions of theorem (*) are met and we can find an **orthogonal matrix** $\boldsymbol{Q}$ (i. e. $\boldsymbol{Q}^{-1} = \boldsymbol{Q}^\top$) such that $\boldsymbol{Q}^\top \boldsymbol{X} \boldsymbol{Q} = \boldsymbol{D}$. The columns of $\boldsymbol{Q}$ are given by the normalized eigenvectors of $\boldsymbol{X}$, and $\boldsymbol{D}$ is a diagonal matrix whose entries are the corresponding eigenvalues

Introduction
**Linear Algebra**
Probability Theory and Statistics
Wrap-Up

Vectors and Vector Operations
Matrix Operations
Determinants and Inverses
Eigenvectors and Eigenvalues
**Symmetric Matrices and Definiteness**

# Example: Eigen-Decomposition

- Consider $\boldsymbol{A} := \begin{pmatrix} 2 & -1 \\ -1 & 2 \end{pmatrix} \in \mathbb{R}^{2 \times 2}$

- We choose one eigenvector for each eigenvalue and divide them by their lengths:

  Eigenvalue 1: $\left( 1/\sqrt{2}, \, 1/\sqrt{2} \right)^{\top}$    Eigenvalue 3: $\left( 1/\sqrt{2}, \, -1/\sqrt{2} \right)^{\top}$

  Thus, the eigen-decomposition of $\boldsymbol{A}$ is given by:

  $$\boldsymbol{D} = \begin{pmatrix} 1 & 0 \\ 0 & 3 \end{pmatrix} = \begin{pmatrix} 1/\sqrt{2} & 1/\sqrt{2} \\ 1/\sqrt{2} & -1/\sqrt{2} \end{pmatrix} \begin{pmatrix} 2 & -1 \\ -1 & 2 \end{pmatrix} \begin{pmatrix} 1/\sqrt{2} & 1/\sqrt{2} \\ 1/\sqrt{2} & -1/\sqrt{2} \end{pmatrix} = \boldsymbol{Q}^{\top} \boldsymbol{A} \boldsymbol{Q}$$

Introduction
**Linear Algebra**
Probability Theory and Statistics
Wrap-Up

Vectors and Vector Operations
Matrix Operations
Determinants and Inverses
Eigenvectors and Eigenvalues
Symmetric Matrices and Definiteness

## Positive (semi-)definite Matrices

- A symmetric matrix $X \in \mathbb{R}^{N \times N}$ is called **positive definite** (notation: $X \succ 0$), if

$$z^\top X z > 0 \quad \forall z \in \mathbb{R}^N \setminus \{0\} \tag{29}$$

- Or **positive semi-definite** (notation: $X \succeq 0$), if

$$z^\top X z \geqslant 0 \quad \forall z \in \mathbb{R}^N \tag{30}$$

**Such matrices are important in machine learning.** For instance, the covariance matrices $\Sigma$ are always positive semi-definite.

**Section:**

**Probability Theory and Statistics**

Random Variables and Common Distributions
Basic Rules of Probability
Expectation and Variance
Kullback-Leibler Divergence

Introduction
Linear Algebra
Probability Theory and Statistics
Wrap-Up

Random Variables and Common Distributions
Basic Rules of Probability
Expectation and Variance
Kullback-Leibler Divergence

## Random Variables

- What is a **random variable**?

Introduction
Linear Algebra
Probability Theory and Statistics
Wrap-Up

Random Variables and Common Distributions
Basic Rules of Probability
Expectation and Variance
Kullback-Leibler Divergence

## Random Variables

- What is a **random variable**?
  - It's a random number determined by chance (according to an underlying distribution).
    To be precise: A random variable $\mathcal{X}$ is a **measurable function**

    $$\mathcal{X} : \Omega \to \mathbb{R} \qquad \text{where } \Omega \text{ is the } \textbf{sample space}$$

  - Examples of random variables in machine learning: Input data, output data, noise
- What is a **probability distribution**?

Introduction
Linear Algebra
Probability Theory and Statistics
Wrap-Up

Random Variables and Common Distributions
Basic Rules of Probability
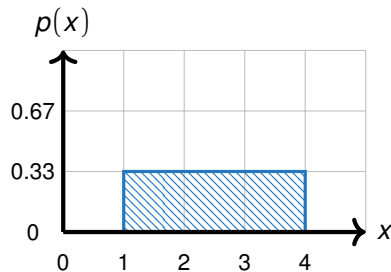Expectation and Variance
Kullback-Leibler Divergence

# Random Variables

- What is a **random variable**?
  - It's a random number determined by chance (according to an underlying distribution).
    To be precise: A random variable $\mathcal{X}$ is a **measurable function**

    $$\mathcal{X} : \Omega \to \mathbb{R} \qquad \text{where } \Omega \text{ is the } \textbf{sample space}$$

  - Examples of random variables in machine learning: Input data, output data, noise
- What is a **probability distribution**?
  - It describes the probability that a random variable is equal to a certain value
  - It can be given by the physics of an experiment (e. g. throwing dice)
  - **Discrete** vs. **continuous** distributions

Introduction
Linear Algebra
Probability Theory and Statistics
Wrap-Up

Random Variables and Common Distributions
Basic Rules of Probability
Expectation and Variance
Kullback-Leibler Divergence

## Uniform Distribution



Every outcome is equally probable within a bounded region $\mathcal{R} := [a, b]$

$$p(\mathcal{X} = x) := \frac{1}{b - a}$$

Introduction
Linear Algebra
Probability Theory and Statistics
Wrap-Up

Random Variables and Common Distributions
Basic Rules of Probability
Expectation and Variance
Kullback-Leibler Divergence

## Discrete Distributions

> A **discrete random variable** takes on discrete values.
>
> **Please note:** Discrete does not mean finite!

**Examples:**

- When throwing a die, the possible values are given by the finite set:

$$\mathcal{X} \in \big\{ 1, 2, 3, 4, 5, 6 \big\}$$

- The number of sand grains at the beach (countably infinite set):

$$\mathcal{X} \in \mathbb{N}$$

Introduction
Linear Algebra
Probability Theory and Statistics
Wrap-Up

Random Variables and Common Distributions
Basic Rules of Probability
Expectation and Variance
Kullback-Leibler Divergence

## Discrete Distributions (Ctd.)

- All probabilities sum up to 1, i. e.:

$$\sum_{x \in \mathcal{X}(\Omega)} p(\mathcal{X} = x) = 1$$

- Discrete distributions are particularly important in classification

- A discrete distribution is described by a **probability mass function** (also called **frequency function**)

Introduction
Linear Algebra
Probability Theory and Statistics
Wrap-Up

Random Variables and Common Distributions
Basic Rules of Probability
Expectation and Variance
Kullback-Leibler Divergence

# BERNOULLI Distribution

- A **BERNOULLI random variable** only takes on two values (e. g. 0 and 1):

$$\mathcal{X} \in \{0, 1\} \tag{31}$$

$$p(\mathcal{X} = 1; \mu) = \mu \tag{32}$$

$$p(\mathcal{X} = 0; \mu) = 1 - \mu \tag{33}$$

$$\mathbb{E}\{\mathcal{X}\} = \mu \tag{34}$$

$$\mathbb{V}\{\mathcal{X}\} = \mu(1 - \mu) \tag{35}$$

- The BERNOULLI distribution is governed only by the parameter $\mu$, the **probability of success**

Introduction
Linear Algebra
Probability Theory and Statistics
Wrap-Up

Random Variables and Common Distributions
Basic Rules of Probability
Expectation and Variance
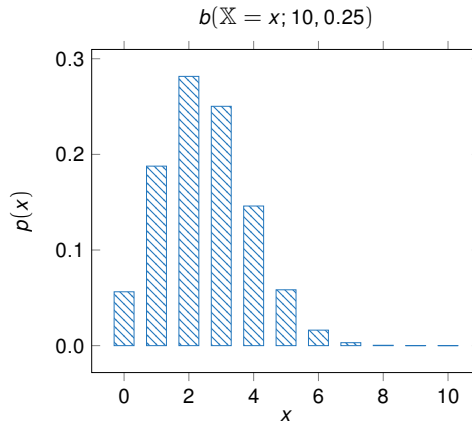Kullback-Leibler Divergence

## Binomial Distribution

- Repeating a BERNOULLI experiment $N$ times leads to the **binomial distribution**
- **Example:** What is the probability of getting $n \in \mathbb{N}$ heads in $N$ trials?

$$b(\mathcal{X} = n; N, \mu) := \binom{N}{n} \mu^n (1 - \mu)^{N-n} \tag{36}$$

$$\mathbb{E}\{\mathcal{X}\} = N\mu \tag{37}$$

$$\mathbb{V}\{\mathcal{X}\} = N\mu(1 - \mu) \tag{38}$$

Introduction
Linear Algebra
Probability Theory and Statistics
Wrap-Up

Random Variables and Common Distributions
Basic Rules of Probability
Expectation and Variance
Kullback-Leibler Divergence

# Binomial Distribution (Ctd.)



$b(\mathbb{X} = x; 10, 0.25)$

Introduction
Linear Algebra
Probability Theory and Statistics
Wrap-Up

Random Variables and Common Distributions
Basic Rules of Probability
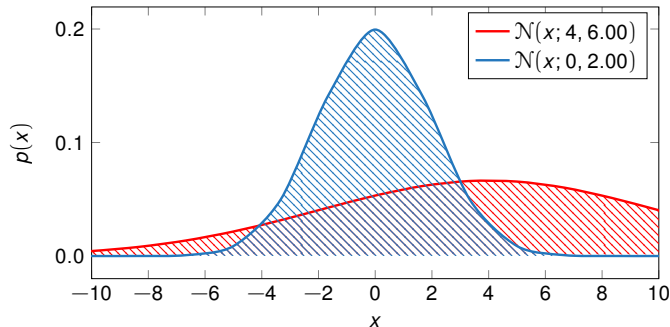Expectation and Variance
Kullback-Leibler Divergence

## Continuous Distributions

**Continuous random variables** take on continuous values

- Continuous distributions are discrete distributions where the **number of discrete values goes to infinity**, while the **probability of each value goes to zero**

- A continuous random variable $\mathcal{X}$ is described by a **probability density function** which integrates to 1, i. e.:

$$\int_{-\infty}^{\infty} p(\mathcal{X} = x)\, \mathrm{d}x = 1$$

Introduction
Linear Algebra
Probability Theory and Statistics
Wrap-Up

Random Variables and Common Distributions
Basic Rules of Probability
Expectation and Variance
Kullback-Leibler Divergence

## GAUSSIAN Distribution



$$\mathcal{N}(\mathcal{X} = x; \mu, \sigma^2) := \frac{1}{\sqrt{2\pi\sigma^2}} \cdot \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right) \qquad (39)$$

Introduction
Linear Algebra
Probability Theory and Statistics
Wrap-Up

Random Variables and Common Distributions
Basic Rules of Probability
Expectation and Variance
Kullback-Leibler Divergence
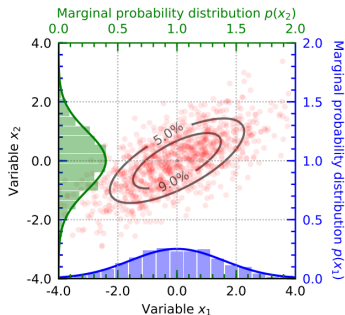
## Central Limit Theorem

**Central Limit Theorem:**

The distribution of the sum of $N$ i.i.d. (independent and identically distributed) random variables **becomes increasingly GAUSSIAN as $N$ increases**

- The GAUSSIAN distribution is one of the most important distributions
- GAUSSIAN distributions often are a good model (due to the central limit theorem)
- Working with GAUSSIANS leads to **analytical solutions for complex operations**

Introduction
Linear Algebra
Probability Theory and Statistics
Wrap-Up

Random Variables and Common Distributions
Basic Rules of Probability
Expectation and Variance
Kullback-Leibler Divergence

## Multivariate GAUSSian Distribution

$$\mathcal{N}_D(\boldsymbol{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) := \frac{1}{\sqrt{(2\pi)^D \det(\boldsymbol{\Sigma})}} \cdot \exp\left(-\frac{1}{2}(\boldsymbol{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\boldsymbol{x} - \boldsymbol{\mu})\right) \tag{40}$$



**Please note:** $\boldsymbol{x}$ and $\boldsymbol{\mu}$ are vectors, while $\boldsymbol{\Sigma}$ is a matrix.

The probability given by $\mathcal{N}_D(\boldsymbol{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$ is still a scalar value!

Introduction
Linear Algebra
Probability Theory and Statistics
Wrap-Up

Random Variables and Common Distributions
Basic Rules of Probability
Expectation and Variance
Kullback-Leibler Divergence

## Basic Rules of Probability

- **Joint distribution:**

$$p(\mathcal{X} \cap \mathcal{Y}) \tag{41}$$

- **Marginal distribution:**

$$p(\mathcal{Y}) = \int_{\mathcal{X}} p(\mathcal{X} \cap \mathcal{Y}) \, \mathrm{d}\mathcal{X} \tag{42}$$

- **Conditional distribution:**

$$p(\mathcal{Y}|\mathcal{X}) = \frac{p(\mathcal{X} \cap \mathcal{Y})}{p(\mathcal{X})} \tag{43}$$

Introduction
Linear Algebra
Probability Theory and Statistics
Wrap-Up

Random Variables and Common Distributions
Basic Rules of Probability
Expectation and Variance
Kullback-Leibler Divergence

## Basic Rules of Probability (Ctd.)

- **Probabilistic independence:**

$$p(\mathcal{X} \cap \mathcal{Y}) = p(\mathcal{X})p(\mathcal{Y}) \tag{44}$$

- **Chain rule of probabilities:**

$$p(\mathcal{X}_1 \cap \ldots \cap \mathcal{X}_N) = p(\mathcal{X}_1 | \mathcal{X}_2 \cap \ldots \cap \mathcal{X}_N)p(\mathcal{X}_2 \cap \ldots \cap \mathcal{X}_N)$$
$$= p(\mathcal{X}_1 | \mathcal{X}_2 \cap \ldots \cap \mathcal{X}_N) \ldots p(\mathcal{X}_{N-1} | \mathcal{X}_N)p(\mathcal{X}_N) \tag{45}$$

- **BAYES' rule:**

$$p(\mathcal{Y} | \mathcal{X}) = \frac{p(\mathcal{X} | \mathcal{Y})p(\mathcal{Y})}{p(\mathcal{X})} \tag{46}$$

Introduction
Linear Algebra
Probability Theory and Statistics
Wrap-Up

Random Variables and Common Distributions
Basic Rules of Probability
Expectation and Variance
Kullback-Leibler Divergence

## Expectation

- The **expectation** $\mathbb{E}$ of a random variable $\mathcal{X}$ is defined by

$$\mathbb{E}\{\mathcal{X}\} := \sum_{k \in \Omega(\mathcal{X})} k \cdot p(\mathcal{X} = k) \tag{47}$$

- Expectations of functions:

$$\mathbb{E}_x\{f\} := \sum_x p(x)f(x) \tag{48}$$

- **Remark:** In the continuous case we have to replace $\sum$ by $\int$

Introduction
Linear Algebra
Probability Theory and Statistics
Wrap-Up

Random Variables and Common Distributions
Basic Rules of Probability
Expectation and Variance
Kullback-Leibler Divergence

## Expectation (Ctd.)

**Rules of expectations:**

- The expectation is a **linear** operation:

$$\mathbb{E}\{a\mathcal{X} + b\mathcal{Y}\} = a\mathbb{E}\{\mathcal{X}\} + b\mathbb{E}\{\mathcal{Y}\} \tag{49}$$

- More general: $\mathbb{E}\left\{\sum_{n=1}^{N} a_n \mathcal{X}_n\right\} = \sum_{n=1}^{N} a_n \mathbb{E}\{\mathcal{X}_n\}$

- If $\mathcal{X}$ and $\mathcal{Y}$ are independent: $\mathbb{E}\{\mathcal{X}\mathcal{Y}\} = \mathbb{E}\{\mathcal{X}\}\mathbb{E}\{\mathcal{Y}\}$

- The expectation is **monotonous**:

$$\mathcal{X} \leqslant \mathcal{Y} \implies \mathbb{E}\{\mathcal{X}\} \leqslant \mathbb{E}\{\mathcal{Y}\} \tag{50}$$

Introduction
Linear Algebra
Probability Theory and Statistics
Wrap-Up

Random Variables and Common Distributions
Basic Rules of Probability
Expectation and Variance
Kullback-Leibler Divergence

## Variance

- The **variance** $\mathbb{V}$ of a random variable $\mathcal{X}$ is defined by

$$\mathbb{V}\{\mathcal{X}\} := \mathbb{E}\{\mathcal{X} - \mathbb{E}^2\{\mathcal{X}\}\} = \mathbb{E}\{\mathcal{X}^2\} - \mathbb{E}^2\{\mathcal{X}\} \tag{51}$$

- $\mathbb{V}$ is **not** linear:

$$\mathbb{V}\{a + b\mathcal{X}\} = b^2\mathbb{V}\{\mathcal{X}\} \tag{52}$$

$$\mathbb{V}\{\mathcal{X} + \mathcal{Y}\} = \mathbb{V}\{\mathcal{X}\} + \mathbb{V}\{\mathcal{Y}\} + \text{cov}\{\mathcal{X}, \mathcal{Y}\} \tag{53}$$

- **BIENAYMÉ's identity:** If $\mathcal{X}$ and $\mathcal{Y}$ are **uncorrelated**, we get:

$$\mathbb{V}\{\mathcal{X} + \mathcal{Y}\} = \mathbb{V}\{\mathcal{X}\} + \mathbb{V}\{\mathcal{Y}\} \tag{54}$$

Introduction
Linear Algebra
Probability Theory and Statistics
Wrap-Up

Random Variables and Common Distributions
Basic Rules of Probability
Expectation and Variance
Kullback-Leibler Divergence

Important

## Covariance

- **Covariances** give a measure of correlation, i.e. how much variables change together

$$\text{cov}\{\mathcal{X}, \mathcal{Y}\} := \mathbb{E}\left\{ (\mathcal{X} - \mathbb{E}\{\mathcal{X}\})(\mathcal{Y} - \mathbb{E}\{\mathcal{Y}\}) \right\}$$
$$= \mathbb{E}\{\mathcal{X}\mathcal{Y}\} - \mathbb{E}\{\mathcal{X}\}\mathbb{E}\{\mathcal{Y}\} \tag{55}$$

- The variance $\mathbb{V}$ of a random variable $\mathcal{X}$ is a special case:

$$\mathbb{V}\{\mathcal{X}\} = \text{cov}\{\mathcal{X}, \mathcal{X}\} \tag{56}$$

Introduction
Linear Algebra
Probability Theory and Statistics
Wrap-Up

Random Variables and Common Distributions
Basic Rules of Probability
Expectation and Variance
Kullback-Leibler Divergence

# Kullback-Leibler Divergence

- The **Kullback-Leibler ($\mathbb{KL}$) divergence** is a similarity measure between two distributions $p$ and $q$:

$$\mathbb{KL}(p\|q) := \sum_x p(x) \cdot \log \frac{p(x)}{q(x)} \tag{57}$$

- Some properties:
  - It is not symmetric: $\mathbb{KL}(p\|q) \neq \mathbb{KL}(q\|p)$
  - It is non-negative: $\mathbb{KL}(p\|q) \geqslant 0$
  - If $\forall x : p(x) = q(x) \implies \mathbb{KL}(p\|q) = 0$

**Section:**

**Wrap-Up**

Summary
Recommended Literature
Self-Test Questions
Lecture Outlook

Introduction
Linear Algebra
Probability Theory and Statistics
Wrap-Up

Summary
Recommended Literature
Self-Test Questions
Lecture Outlook

## Summary

- **Mathematics play a major role in machine learning!**
- **Linear algebra:**
  - You should know what vectors are and what you can do with them (addition, multiplication, transpose, ...)
  - The same applies to matrices
  - You should know the concept of **determinants** and how to **invert matrices**
  - **Eigenvectors** and **eigenvalues** are important tools in machine learning
  - The **eigen-decomposition** plays an import role in many machine learning applications

Introduction
Linear Algebra
Probability Theory and Statistics
Wrap-Up

Summary
Recommended Literature
Self-Test Questions
Lecture Outlook

# Summary (Ctd.)

- **Probability theory and statistics:**
  - Random variables are numbers **determined by chance**
  - Probability distributions describe a **probability mass** or **probability density**
  - **Discrete distributions:** BERNOULLI, Binomial, Multinomial
  - **Continuous distribution:** GAUSSian distribution
  - GAUSSians are important in machine learning and have appealing properties
  - Terms you should know: Joint-, marginal- and conditional distribution, chain rule, probabilistic independence, Bayes' rule
  - You should know what **expectation** and **variance** of distributions are

Introduction
Linear Algebra
Probability Theory and Statistics
**Wrap-Up**

Summary
Recommended Literature
Self-Test Questions
Lecture Outlook

# Recommended Literature

**❶ Linear algebra:**

- [DEISENROTH.2019], chapter 2
- [DEISENROTH.2019], chapter 3
- [DEISENROTH.2019], chapter 4

**❷ Probability theory and statistics:**

- [DEISENROTH.2019], chapter 6

(For free PDF versions, see list in GitHub readme!)

Introduction
Linear Algebra
Probability Theory and Statistics
Wrap-Up

Summary
Recommended Literature
Self-Test Questions
Lecture Outlook

## Self-Test Questions

1. What is a vector and what is a matrix?

2. What is the result of an inner product / outer product?

3. How can you invert matrices? Is this always possible?

4. What is an eigenvalue problem? How can you compute eigenvectors and eigenvalues?

5. What are random variables and probability distributions?

6. Why is the GAUSSIAN distribution so important?

7. What is BAYES' rule? Explain its components!

Introduction
Linear Algebra
Probability Theory and Statistics
Wrap-Up

Summary
Recommended Literature
Self-Test Questions
Lecture Outlook

## What's next...?

- **I** Machine Learning Introduction
- **II** Optimization Techniques
- **III** Bayesian Decision Theory
- **IV** Non-parametric Density Estimation
- **V** Probabilistic Graphical Models
- **VI** Linear Regression
- **VII** Logistic Regression
- **VIII** Deep Learning

- **IX** Evaluation
- **X** Decision Trees
- **XI** Support Vector Machines
- **XII** Clustering
- **XIII** Principal Component Analysis
- **XIV** Reinforcement Learning
- **XV** Advanced Regression

# Thank you very much for the attention!

**\* \* \* Artificial Intelligence and Machine Learning \* \* \***

**Topic:** Mathematics Refresher

**Term:** Summer term 2025

**Contact:**

Daniel Wehner, M.Sc.

SAP SE / DHBW Mannheim

<u>daniel.wehner@sap.com</u>

## Do you have any questions?