

# Probabilistic Artificial Intelligence

Lecturer(s): Prof. Andreas Krause

Author: Saurav Banka

Semester: HS 2025

*Last edited:* September 23, 2025

# Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
1.1	Motivation . . . . .	3
1.2	Probability Basics . . . . .	3
1.3	Multivariate Gaussians . . . . .	5
<b>2</b>	<b>Bayesian Linear Regression</b>	<b>7</b>
<b>3</b>	<b>Gaussian Processes</b>	<b>8</b>
<b>4</b>	<b>Variational Inference</b>	<b>9</b>
<b>5</b>	<b>Markov Chain Monte Carlo.</b>	<b>10</b>
<b>6</b>	<b>Bayesian Deep Learning</b>	<b>11</b>
<b>7</b>	<b>Active Learning and Bayesian Optimization</b>	<b>12</b>
<b>8</b>	<b>Markov Decision Processes</b>	<b>13</b>
<b>9</b>	<b>Reinforcement Learning</b>	<b>14</b>
<b>A</b>	<b>Math background: Fourier Transforms</b>	<b>15</b>
<b>B</b>	<b>Useful Math Identities</b>	<b>16</b>

# 1 Introduction

## 1.1 Motivation

Uncertainty is all around us. This course discusses how to enable data-driven reasoning and decision-making under uncertainty.

## 1.2 Probability Basics

Probability is formally defined by a probability space  $(\Omega, F, P)$ :

- $\Omega$ : Set of atomic events (e.g., throwing a die).
- $F \subseteq 2^\Omega$  is a  $\sigma$ -algebra. Intuitively, it formalizes which types of events can occur.
  - $\Omega \in F$
  - $A \in F \implies \Omega \setminus A \in F$  (intuition: the complement of an event is also an event).
  - $A_1, A_2, \dots \in F \implies \bigcup_i A_i \in F$  (intuition: the countable union of events must also be an event).
- Probability measure  $P : F \rightarrow [0, 1]$  assigns a value to each event in  $F$ .

**Axioms (Kolmogorov):**

- $P(\Omega) = 1$
- $P(A) \geq 0$  for all  $A \in F$
- For all disjoint  $A_i \in F$ ,

$$P\left(\bigcup_{i=1}^{\infty} A_i\right) = \sum_{i=1}^{\infty} P(A_i).$$

**Random variables.** A random variable (RV) is a mapping  $X : \Omega \rightarrow D$  where  $D$  is some set of interest. Then,

$$P(X = x) = P(\{\omega \in \Omega : X(\omega) = x\}),$$

i.e., the probability that the variable  $X$  takes on a particular value  $x$ .

**Continuous case.** For continuous values, we introduce a probability density:

$$p(t) = \lim_{\delta \rightarrow 0} \frac{P(T \in [t, t + \delta])}{\delta}.$$

**Definition 1.1.** (Expected value) Given a random variable  $X$  on domain  $D$  and a function  $f : D \rightarrow \mathbb{R}$ , the expectation is

$$\mathbb{E}[f(X)] = \int_{x \in D} f(x) p(x) dx.$$

Instead of a random variable, we can specify a random vector  $\mathbf{X} = [X_i(\omega)]_{i=1}^N$ . The joint distribution describes relations among all variables.

**Definition 1.2.** Conditional probability: For events  $A, B$  with  $P(B) > 0$ ,

$$P(A | B) = \frac{P(A, B)}{P(B)}.$$

**Definition 1.3.** (Marginalization / Sum Rule)

$$P(X_{[1:n] \setminus i}) = \sum_{x_i} P(X_{1:i-1}, x_i, X_{i+1:n}).$$

**Definition 1.4.** (Product Rule / Chain Rule)

$$P(X_1, \dots, X_n) = P(X_1) \prod_{i=2}^n P(X_i | X_1, \dots, X_{i-1}).$$

**Theorem 1.1.** (Bayes' Rule)

$$P(X | Y) = \frac{P(X)P(Y | X)}{\sum_x P(X = x) P(Y | X = x)}.$$

Proof: Can be seen by directly applying the product and sum rule to the def of conditional independence.

**Definition 1.5.** (Independence) Random variables  $X_1, \dots, X_n$  are independent if

$$P(X_1 = x_1, \dots, X_n = x_n) = \prod_{i=1}^n P(X_i = x_i).$$

Equivalently,  $P(X_1 | X_2) = P(X_1)$ .

**Definition 1.6.** (Conditional independence) We write  $X \perp\!\!\!\perp Y | Z$  iff for all  $x, y, z$ ,

$$P(X = x, Y = y | Z = z) = P(X = x | Z = z) P(Y = y | Z = z),$$

or equivalently,

$$P(X = x | Y = y, Z = z) = P(X = x | Z = z).$$

**High-dimensional challenges.** Without dependence structure, we face:

- **Representation:** requires  $2^N - 1$  parameters for  $N$  binary RVs.
- **Marginalization:** computing marginals requires summing over  $2^{N-1}$  terms.
- **Inference:** conditional queries, learning, and prediction are all expensive.

### 1.3 Multivariate Gaussians

Multivariate Gaussians (MVGs) are tractable for many of these problems.

**Definition 1.7.** (Covariance Matrix)

For a random vector  $\mathbf{x} \in \mathbb{R}^d$  with mean  $\boldsymbol{\mu} = \mathbb{E}[\mathbf{x}]$ , the covariance matrix is

$$\boldsymbol{\Sigma} = \mathbb{E}[(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^\top].$$

Equivalently,

$$\Sigma_{ij} = \text{Cov}(X_i, X_j) = \mathbb{E}[(X_i - \mu_i)(X_j - \mu_j)].$$

**Definition 1.8.** (Multivariate Gaussian) The density of  $\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$  is

$$p(\mathbf{x}) = \frac{1}{(2\pi)^{d/2} \sqrt{|\boldsymbol{\Sigma}|}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right).$$

**Theorem 1.2.** In the Gaussian case, uncorrelated variables are independent.

#### 1.3.1 Computing marginals

Marginals of an index set  $A$  are obtained by simply selecting the corresponding components of  $\boldsymbol{\mu}$  and  $\boldsymbol{\Sigma}$ . Formally,  $X_A \sim \mathcal{N}(\mu_A, \Sigma_{AA})$ .

#### 1.3.2 Conditional distributions

Partition  $\mathbf{x} = \begin{bmatrix} \mathbf{x}_A \\ \mathbf{x}_B \end{bmatrix}$  with corresponding mean and covariance

$$\boldsymbol{\mu} = \begin{bmatrix} \mu_A \\ \mu_B \end{bmatrix}, \quad \boldsymbol{\Sigma} = \begin{bmatrix} \Sigma_{AA} & \Sigma_{AB} \\ \Sigma_{BA} & \Sigma_{BB} \end{bmatrix}.$$

Then,

$$\mathbf{x}_A \mid \mathbf{x}_B = b \sim \mathcal{N}(\mu_A + \Sigma_{AB}\Sigma_{BB}^{-1}(b - \mu_B), \Sigma_{AA} - \Sigma_{AB}\Sigma_{BB}^{-1}\Sigma_{BA}).$$

**Note 1.1. Interpretation.**

- The updated mean

$$\mu_{A|B} = \mu_A + \Sigma_{AB}\Sigma_{BB}^{-1}(b - \mu_B)$$

is the prior mean  $\mu_A$  shifted by an adjustment term.

- $\Sigma_{AB}$  is the *cross-covariance*, describing how uncertainty in  $A$  co-varies with  $B$ .
- $\Sigma_{BB}^{-1}$  is the *precision matrix* of  $B$ , which corrects for redundancy among components of  $B$ : highly correlated features of  $B$  contribute less uniquely to the update.
- Together,  $\Sigma_{AB}\Sigma_{BB}^{-1}$  acts like a regression coefficient matrix mapping observed deviations

$(b - \mu_B)$  into updates of  $A$ 's mean.

- The updated covariance

$$\Sigma_{A|B} = \Sigma_{AA} - \Sigma_{AB}\Sigma_{BB}^{-1}\Sigma_{BA}$$

is always smaller (in the PSD sense) than  $\Sigma_{AA}$ . Intuitively: conditioning on  $B$  reduces uncertainty about  $A$ .

### 1.3.3 Affine transformations of Gaussians

If  $\mathbf{x} \sim \mathcal{N}(\mu, \Sigma)$  and  $\mathbf{y} = M\mathbf{x} + b$  for matrix  $M$  and vector  $b$ , then

$$\mathbf{y} \sim \mathcal{N}(M\mu + b, M\Sigma M^T).$$

A natural extension is to use 1-hot vectors as rows of  $A$  to express marginals or sums: this shows directly that linear combinations (and in particular sums) of Gaussians are Gaussian. It also allows us to represent *degenerate Gaussians*, where some components have constant mean and zero variance to allow for the covariance matrix to have certain eigenvalues corresponding to 0.

### 1.3.4 Conditional Linear Gaussians

If  $\mathbf{y} | \mathbf{x} \sim \mathcal{N}(A\mathbf{x} + b, \Sigma_y)$  with  $\mathbf{x} \sim \mathcal{N}(\mu_x, \Sigma_x)$ , then the joint distribution is Gaussian:

$$\begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} \sim \mathcal{N}\left(\begin{bmatrix} \mu_x \\ A\mu_x + b \end{bmatrix}, \begin{bmatrix} \Sigma_x & \Sigma_x A^T \\ A\Sigma_x & A\Sigma_x A^T + \Sigma_y \end{bmatrix}\right).$$

This structure is fundamental in Bayesian networks with Gaussian conditional distributions.

#### Note 1.2. Interpretation.

The conditional distribution  $P(\mathbf{x} | \mathbf{y})$  has the form

$$\mathbf{x} | \mathbf{y} \sim \mathcal{N}(M\mathbf{y} + b, \Sigma_{x|y}),$$

for some matrix  $M$ , vector  $b$ , and covariance  $\Sigma_{x|y}$ .

- The term  $M\mathbf{y} + b$  shows that the conditional mean is an *affine function* of the observed variable  $\mathbf{y}$ . This means that, once we observe  $\mathbf{y}$ , our best prediction of  $\mathbf{x}$  is a linear regression on  $\mathbf{y}$ .
- The covariance  $\Sigma_{x|y}$  represents the residual uncertainty that cannot be explained by  $\mathbf{y}$ . It acts like *independent Gaussian noise* added to the linear prediction.
- Equivalently: the conditional Gaussian says

$$\mathbf{x} = M\mathbf{y} + b + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \Sigma_{x|y}),$$

which makes explicit the regression + noise interpretation.

This viewpoint underlies Gaussian graphical models and Bayesian linear regression.

## 2 Bayesian Linear Regression

### 3 Gaussian Processes



## 4 Variational Inference

## 5 Markov Chain Monte Carlo.

## 6 Bayesian Deep Learning

## 7 Active Learning and Bayesian Optimization

## 8 Markov Decision Processes

## 9 Reinforcement Learning

## **A Math background: Fourier Transforms**

## **B Useful Math Identities**