

Chapter 4

Multivariate Random Variables

In this chapter, we will review some topics related to random vectors, which will be of use in the following chapters.

4.1 Gaussian Random Vectors (Multivariate Normal Distribution)

Recall that a random variable X is Gaussian (normal) with mean μ and variance $\sigma^2 > 0$ if the pdf of X is given by

$$p_X(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp -\frac{(x - \mu)^2}{2\sigma^2}. \quad (4.1)$$

Definition 4.1. A collection of random variables is **jointly Gaussian** if any linear combination of these variables is Gaussian. A **Gaussian random vector**, also known as a *multivariate normal vector*, is a vector whose elements are jointly Gaussian. A collection of random vectors is jointly Gaussian if the vector obtained by concatenating them is jointly Gaussian.

Example 4.2. If $\begin{pmatrix} X \\ Y \end{pmatrix}$ is a Gaussian vector, then $Z = 2X + 3Y$ is Gaussian. Furthermore,

$$\mathbb{E}[Z] = 2\mathbb{E}[X] + 3\mathbb{E}[Y], \quad (4.2)$$

$$\text{Var}(Z) = \text{Cov}(2X + 3Y, 2X + 3Y) = 4\text{Cov}(X, X) + 12\text{Cov}(X, Y) + 9\text{Cov}(Y, Y) \quad (4.3)$$

$$= 4\text{Var}(X) + 12\text{Cov}(X, Y) + 9\text{Var}(Y), \quad (4.4)$$

which completely characterizes the distribution of Z as $Z \sim \mathcal{N}(\mathbb{E}[Z], \text{Var}(Z))$. \triangle

For a Gaussian random vector \mathbf{X} of dimension d , with mean $\mathbb{E}[\mathbf{X}] = \boldsymbol{\mu}$ and covariance matrix $\mathbf{K} = \text{Cov}(\mathbf{X}) = \mathbb{E}[(\mathbf{X} - \boldsymbol{\mu})(\mathbf{X} - \boldsymbol{\mu})^T]$, we have

$$p_{\mathbf{X}}(\mathbf{x}) = \frac{1}{(2\pi)^{d/2} |\mathbf{K}|^{1/2}} \exp \left(-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \mathbf{K}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right), \quad (4.5)$$

provided that the covariance matrix is invertible.

The elements of \mathbf{X} are **independent** if and only if the covariance matrix is diagonal.

4.1.1 Maximum likelihood estimation

Consider a d -dimensional random vector $\mathbf{X} = (X_1, \dots, X_d)$ with distribution $\mathcal{N}(\boldsymbol{\theta}^*, \mathbf{K}^*)$ given in (4.5), where $\boldsymbol{\theta}^*, \mathbf{K}^*$ are unknown. Suppose we are interested in the relationship between X_d and X_1, \dots, X_{d-1} . For

example, for $\mathbf{X}^T = (X_1, X_2, X_3)$, X_1 and X_2 could indicate the heights of the parents and X_3 could be the height of the child. We may, for example, be interested in finding $\mathbb{E}[X_d | X_1, \dots, X_{d-1}]$, thus estimating X_d based on X_1, \dots, X_{d-1} . If we find the distribution, in other words, $\boldsymbol{\theta}^*, \mathbf{K}^*$, we can do so. Furthermore, the matrix \mathbf{K}^* can indicate which dimensions are more strongly correlated.

Consider a set of n iid samples $\mathcal{D} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$, where each \mathbf{x}_i is a sample of \mathbf{X} . We denote the elements of \mathbf{x}_i as $\mathbf{x}_i = (x_{i1}, \dots, x_{id})$.

To estimate $\boldsymbol{\theta}^*$ and \mathbf{K}^* , we write

$$\ell(\boldsymbol{\theta}, \mathbf{K}) = \ln p(\mathcal{D}; \boldsymbol{\theta}, \mathbf{K}) = \sum_{i=1}^n \ln p(\mathbf{x}_i; \boldsymbol{\theta}, \mathbf{K}) \quad (4.6)$$

$$\doteq \frac{n}{2} \ln |\mathbf{K}^{-1}| - \frac{1}{2} \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\theta})^T \mathbf{K}^{-1} (\mathbf{x}_i - \boldsymbol{\theta}), \quad (4.7)$$

where we have used the fact that $|\mathbf{K}^{-1}| = \frac{1}{|\mathbf{K}|}$.

As seen in the appendix (last chapter), for a symmetric matrix \mathbf{A} , we have $\frac{d}{d\mathbf{v}}(\mathbf{y}^T \mathbf{A} \mathbf{y}) = 2\mathbf{y}^T \mathbf{A} \frac{d\mathbf{y}}{d\mathbf{v}}$. Hence,

$$\frac{\partial \ell}{\partial \boldsymbol{\theta}} = -\frac{1}{2} \sum_{i=1}^n 2(\mathbf{x}_i - \boldsymbol{\theta})^T \mathbf{K}^{-1} (-\mathbf{I}) = \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\theta})^T \mathbf{K}^{-1}. \quad (4.8)$$

Setting this equal to zero yields

$$\hat{\boldsymbol{\theta}}_{ML} = \bar{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i. \quad (4.9)$$

Exercise 4.3. Using the facts

$$\frac{\partial}{\partial \mathbf{A}} \mathbf{x}^T \mathbf{A} \mathbf{x} = \mathbf{x} \mathbf{x}^T, \quad \frac{\partial}{\partial \mathbf{A}} \ln |\mathbf{A}| = \mathbf{A}^{-T} \quad (4.10)$$

prove that

$$\hat{\mathbf{K}}_{ML} = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T \quad (4.11)$$

△

4.1.2 Bayesian estimation

We now solve the same problem using Bayesian estimation, with the following likelihood

$$\mathbf{X} | \boldsymbol{\Theta} \sim \mathcal{N}(\boldsymbol{\Theta}, \mathbf{K}), \quad (4.12)$$

$$p(\mathbf{x}_1^n | \boldsymbol{\theta}) \propto \exp \left(-\frac{1}{2} \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\theta})^T \mathbf{K}^{-1} (\mathbf{x}_i - \boldsymbol{\theta}) \right), \quad (4.13)$$

where, for simplicity, we assume \mathbf{K} is known and we only need to estimate $\boldsymbol{\Theta}$. As the prior, we choose

$$\boldsymbol{\Theta} \sim \mathcal{N}(\boldsymbol{\mu}_0, \mathbf{S}_0) \quad (4.14)$$

$$p(\boldsymbol{\theta}) \propto \exp \left(-\frac{1}{2} (\boldsymbol{\theta} - \boldsymbol{\mu}_0)^T \mathbf{S}_0^{-1} (\boldsymbol{\theta} - \boldsymbol{\mu}_0) \right). \quad (4.15)$$

Hence,

$$p(\boldsymbol{\theta} | \mathbf{x}_1^n) \propto \exp \left(-\frac{1}{2} (\boldsymbol{\theta} - \boldsymbol{\mu}_0)^T \mathbf{S}_0^{-1} (\boldsymbol{\theta} - \boldsymbol{\mu}_0) - \frac{1}{2} \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\theta})^T \mathbf{K}^{-1} (\mathbf{x}_i - \boldsymbol{\theta}) \right). \quad (4.16)$$

The exponent in the posterior is quadratic in $\boldsymbol{\theta}$, indicating that $\boldsymbol{\Theta}$ has a Gaussian distribution. So $\boldsymbol{\Theta}|\mathbf{x}_1^n \sim \mathcal{N}(\hat{\boldsymbol{\theta}}_n, \mathbf{S}_n)$, for appropriate choices of $\hat{\boldsymbol{\theta}}_n$ and \mathbf{S}_n ,

$$p(\boldsymbol{\theta}|\mathbf{x}_1^n) \propto \exp\left(-\frac{1}{2}(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}_n)^T \mathbf{S}_n^{-1}(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}_n)\right). \quad (4.17)$$

To find $\hat{\boldsymbol{\theta}}_n$ and \mathbf{S}_n , we equate (4.16) and (4.17), ignoring constant multiplicative factors, which leads to

$$(\boldsymbol{\theta} - \boldsymbol{\mu}_0)^T \mathbf{S}_0^{-1}(\boldsymbol{\theta} - \boldsymbol{\mu}_0) + \sum_{i=1}^n (\boldsymbol{\theta} - \mathbf{x}_i)^T \mathbf{K}^{-1}(\boldsymbol{\theta} - \mathbf{x}_i) \doteq (\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}_n)^T \mathbf{S}_n^{-1}(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}_n), \quad (4.18)$$

$$\boldsymbol{\theta}^T \mathbf{S}_0^{-1} \boldsymbol{\theta} - 2\boldsymbol{\theta}^T \mathbf{S}_0^{-1} \boldsymbol{\mu}_0 + n\boldsymbol{\theta}^T \mathbf{K}^{-1} \boldsymbol{\theta} - 2\boldsymbol{\theta}^T \mathbf{K}^{-1} \sum_{i=1}^n \mathbf{x}_i \doteq \boldsymbol{\theta}^T \mathbf{S}_n^{-1} \boldsymbol{\theta} - 2\boldsymbol{\theta}^T \mathbf{S}_n^{-1} \hat{\boldsymbol{\theta}}_n. \quad (4.19)$$

Here, we have used the fact that

$$(\mathbf{a} - \mathbf{b})^T \mathbf{A}(\mathbf{a} - \mathbf{b}) = \mathbf{a}^T \mathbf{A} \mathbf{a} - \mathbf{a}^T \mathbf{A} \mathbf{b} - \mathbf{b}^T \mathbf{A} \mathbf{a} + \mathbf{b}^T \mathbf{A} \mathbf{b} = \mathbf{a}^T \mathbf{A} \mathbf{a} - 2\mathbf{a}^T \mathbf{A} \mathbf{b} + \mathbf{b}^T \mathbf{A} \mathbf{b},$$

for vectors \mathbf{a}, \mathbf{b} and a symmetric matrix \mathbf{A} . Note that $\mathbf{a}^T \mathbf{A} \mathbf{b} = \mathbf{b}^T \mathbf{A} \mathbf{a}$, as both sides are scalars and $\mathbf{a}^T \mathbf{A} \mathbf{b} = (\mathbf{a}^T \mathbf{A} \mathbf{b})^T = \mathbf{b}^T \mathbf{A} \mathbf{a}$.

We now collect the terms of the form $\boldsymbol{\theta}^T \mathbf{A} \boldsymbol{\theta}$,

$$\boldsymbol{\theta}^T (\mathbf{S}_0^{-1} + n\mathbf{K}^{-1}) \boldsymbol{\theta} - 2\boldsymbol{\theta}^T (\mathbf{S}_0^{-1} \boldsymbol{\mu}_0 + \mathbf{K}^{-1} \sum_{i=1}^n \mathbf{x}_i) \doteq \boldsymbol{\theta}^T \mathbf{S}_n^{-1} \boldsymbol{\theta} - 2\boldsymbol{\theta}^T \mathbf{S}_n^{-1} \hat{\boldsymbol{\theta}}_n, \quad (4.20)$$

leading to the following values for the parameters of the posterior distribution $\boldsymbol{\Theta}|\mathbf{x}_1^n \sim \mathcal{N}(\hat{\boldsymbol{\theta}}_n, \mathbf{S}_n)$,

$$\mathbf{S}_n^{-1} = \mathbf{S}_0^{-1} + n\mathbf{K}^{-1}, \quad (4.21)$$

$$\hat{\boldsymbol{\theta}}_n = \mathbf{S}_n (\mathbf{S}_0^{-1} \boldsymbol{\mu}_0 + n\mathbf{K}^{-1} \bar{\mathbf{x}}) \quad (4.22)$$

$$= (\mathbf{S}_0^{-1} + n\mathbf{K}^{-1})^{-1} (\mathbf{S}_0^{-1} \boldsymbol{\mu}_0 + n\mathbf{K}^{-1} \bar{\mathbf{x}}), \quad (4.23)$$

where $\bar{\mathbf{x}}$ is $\sum_{i=1}^n \mathbf{x}_i / n$. The posterior mean, $\hat{\boldsymbol{\theta}}_n$, which we can also view as a point estimate, is the weighted average of the prior mean $\boldsymbol{\mu}_0$ and what is suggested by the data $\bar{\mathbf{x}}$.

Exercise 4.4. Find $\hat{\boldsymbol{\theta}}_n$ and \mathbf{S}_n^{-1} when $\mathbf{S}_0 = s^2 \mathbf{I}$ and $\mathbf{K} = \sigma^2 \mathbf{I}$ and interpret the results. \triangle