

# 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 Review of Linear Algebra

For two vectors  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ , the **inner product**  $\langle \mathbf{x}, \mathbf{y} \rangle$  of  $\mathbf{x}$  and  $\mathbf{y}$  is

$$\mathbf{x} = \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix}, \quad \mathbf{y} = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}, \quad \langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{x}^T \mathbf{y} = \sum_{i=1}^n x_i y_i.$$

where  $\mathbf{x}^T$  is the transpose of  $\mathbf{x}$ .

The **length** or the  $\ell_2$  norm of a vector  $\mathbf{x}$  is  $\|\mathbf{x}\| = \|\mathbf{x}\|_2 = \sqrt{\mathbf{x}^T \mathbf{x}}$  and we have  $\|\mathbf{x}\|_2^2 = \mathbf{x}^T \mathbf{x}$ . Let  $\alpha$  be the angle between  $\mathbf{x}$  and  $\mathbf{y}$ . Then  $\mathbf{x}^T \mathbf{y} = \|\mathbf{x}\| \|\mathbf{y}\| \cos \alpha$ . If  $\mathbf{x}^T \mathbf{y} = 0$ , then the two are called **orthogonal**.

For a collection of vectors  $\mathbf{v}_1, \dots, \mathbf{v}_m$ , a **linear combination** of these is any vector of the form  $a_1 \mathbf{v}_1 + \dots + a_m \mathbf{v}_m$ ,  $a_i \in \mathbb{R}$ . The set of all linear combinations of  $\mathbf{v}_1, \dots, \mathbf{v}_m$  is their **span** and denoted as  $\text{Span}\{\mathbf{v}_1, \dots, \mathbf{v}_m\}$ . This is a **subspace** (think line, plane, or the whole space). For a matrix  $A$ , the span of the columns of  $A$  is the **column space** of  $A$ .

The vectors  $\mathbf{v}_1, \dots, \mathbf{v}_m$  are **linearly independent** if there is no vector among them that can be written as a linear combination of the others, and linearly dependent otherwise. The vectors are linearly independent if and only if the only values for  $a_1, \dots, a_m$  satisfying  $a_1 \mathbf{v}_1 + \dots + a_m \mathbf{v}_m = \mathbf{0}$  are  $a_1, \dots, a_m = 0$ . In particular, the columns of a matrix  $A$  are linearly independent if and only if the only vector  $\mathbf{a}$  satisfying  $A\mathbf{a} = \mathbf{0}$  is  $\mathbf{a} = \mathbf{0}$ .

The **inverse** of a square matrix  $A$  is a matrix  $A^{-1}$  such that  $AA^{-1} = A^{-1}A = I$ , where  $I$  is the **identity matrix**, which has 1s on the diagonal and 0s elsewhere. A matrix that has an inverse is called **invertible**. A square matrix is invertible  $\iff$  for all distinct vectors  $\mathbf{a}$  and  $\mathbf{b}$ , we have  $A\mathbf{a} \neq A\mathbf{b}$   $\iff$  the only solution to  $A\mathbf{x} = \mathbf{0}$  is  $\mathbf{x} = \mathbf{0}$   $\iff$  its columns are linearly independent  $\iff$  its determinant  $|A|$  is nonzero. We also have  $|A^{-1}| = \frac{1}{|A|}$ .

Given a subspace  $S$  (e.g., a plane or the column space of a matrix) and a vector  $\mathbf{y}$ , let  $\hat{\mathbf{y}}$  be the vector in the subspace that is closest to  $\mathbf{y}$ . That is, we find  $\hat{\mathbf{y}} \in S$  such that  $\|\mathbf{y} - \hat{\mathbf{y}}\|$  is minimized. Then  $\hat{\mathbf{y}}$  is called the **projection** of  $\mathbf{y}$  onto the subspace  $S$ .

**Lemma 66.** *Let  $\hat{\mathbf{y}}$  be the projection of a vector  $\mathbf{y}$  onto a subspace  $S$ . Then  $\mathbf{y} - \hat{\mathbf{y}}$  is orthogonal to every vector in  $S$ .*

*Proof.* Suppose that this is not the case. Then there is a nonzero vector  $\mathbf{v} \in S$  such that  $(\mathbf{y} - \hat{\mathbf{y}})^T \mathbf{v} \neq 0$ . We will show that this contradicts the minimality of  $\|\mathbf{y} - \hat{\mathbf{y}}\|$ . For any  $a \in \mathbb{R}$ ,

$$\begin{aligned} \|\mathbf{y} - \hat{\mathbf{y}} - a\mathbf{v}\|_2^2 &= (\mathbf{y} - \hat{\mathbf{y}} - a\mathbf{v})^T (\mathbf{y} - \hat{\mathbf{y}} - a\mathbf{v}) \\ &= \|\mathbf{y} - \hat{\mathbf{y}}\|_2^2 - 2a\mathbf{v}^T (\mathbf{y} - \hat{\mathbf{y}}) + a^2 \|\mathbf{v}\|_2^2. \end{aligned}$$

This is a convex function in  $a$ . So setting the derivative to 0 gives the value of  $a$  that minimizes the error:

$$\frac{\partial}{\partial a} \|\mathbf{y} - \hat{\mathbf{y}} - a\mathbf{v}\|_2^2 = -2\mathbf{v}^T (\mathbf{y} - \hat{\mathbf{y}}) + 2a\|\mathbf{v}\|_2^2 = 0 \Rightarrow a = \frac{\mathbf{v}^T (\mathbf{y} - \hat{\mathbf{y}})}{\|\mathbf{v}\|_2^2} \neq 0.$$

Let

$$\hat{\mathbf{y}}' = \hat{\mathbf{y}} + \frac{\mathbf{v}^T (\mathbf{y} - \hat{\mathbf{y}})}{\mathbf{v}^T \mathbf{v}} \mathbf{v},$$

and note that  $\hat{\mathbf{y}}'$  is also in  $S$  but it is closer to  $\mathbf{y}$  contradicting the optimality of  $\hat{\mathbf{y}}$ . □

## 4.2 Random vectors

A **random vector** is a vector of random variables. Consider the random vectors  $\mathbf{x}$  and  $\mathbf{y}$

$$\mathbf{x} = \begin{pmatrix} x_1 \\ \vdots \\ x_m \end{pmatrix}, \quad \mathbf{y} = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}.$$

The **expected value** of  $\mathbf{x}$  is

$$\mathbb{E} \mathbf{x} = \begin{pmatrix} \mathbb{E} x_1 \\ \vdots \\ \mathbb{E} x_m \end{pmatrix}.$$

The **correlation matrix** of  $\mathbf{x}$  and  $\mathbf{y}$  is the  $m \times n$  matrix  $\mathbb{E}[\mathbf{x}\mathbf{y}^T]$ , whose  $i, j$ th element is  $\mathbb{E}[x_i y_j]$ . The **cross-covariance matrix** of  $\mathbf{x}$  and  $\mathbf{y}$  is  $\text{Cov}(\mathbf{x}, \mathbf{y})$  is the matrix  $\mathbb{E}[(\mathbf{x} - \mathbb{E} \mathbf{x})^T (\mathbf{y} - \mathbb{E} \mathbf{y})^T]$ , whose  $i, j$ th element is  $\text{Cov}(x_i, y_j)$ . The covariance of a vector  $\mathbf{x}$  is  $\text{Cov}(\mathbf{x}) = \text{Cov}(\mathbf{x}, \mathbf{x})$ . The **conditional expectation**  $\mathbb{E}[\mathbf{x}|\mathbf{y}]$  of  $\mathbf{x}$  given  $\mathbf{y}$  is a vector whose  $i$ th element is  $\mathbb{E}[x_i|\mathbf{y}]$ .

For matrices  $A, B$ , deterministic vectors  $\mathbf{a}, \mathbf{b}$ , and random vectors  $\mathbf{x}, \mathbf{y}, \mathbf{w}, \mathbf{z}$ , we have [1]

- $\mathbb{E}[A\mathbf{x} + \mathbf{a}] = A \mathbb{E} \mathbf{x} + \mathbf{a}$
- $\text{Cov}(\mathbf{x}, \mathbf{y}) = \mathbb{E}[\mathbf{x}\mathbf{y}^T] - \mathbb{E} \mathbf{x} \mathbb{E} \mathbf{y}^T$
- $\mathbb{E}[(A\mathbf{x})(B\mathbf{y})^T] = A \mathbb{E}[\mathbf{x}\mathbf{y}^T] B^T$
- $\text{Cov}(A\mathbf{x} + \mathbf{a}, B\mathbf{y} + \mathbf{b}) = A \text{Cov}(\mathbf{x}, \mathbf{y}) B^T$
- $\text{Cov}(A\mathbf{x} + \mathbf{a}) = A \text{Cov}(\mathbf{x}) A^T$
- $\text{Cov}(\mathbf{w} + \mathbf{x}, \mathbf{y} + \mathbf{z}) = \text{Cov}(\mathbf{w}, \mathbf{y}) + \text{Cov}(\mathbf{w}, \mathbf{z}) + \text{Cov}(\mathbf{x}, \mathbf{y}) + \text{Cov}(\mathbf{x}, \mathbf{z})$

### 4.3 Gaussian Random Vectors (Joint Gaussian Distribution)

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

$$p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right).$$

**Definition 67.** 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 are jointly Gaussian if the vector obtained by concatenating them is jointly Gaussian.

**Example 68.** For example if  $\begin{pmatrix} x \\ y \end{pmatrix}$  is a Gaussian vector, then  $z = 2x + 3y$  is Gaussian. Furthermore,

$$\begin{aligned} \mathbb{E}[z] &= 2\mathbb{E}[x] + 3\mathbb{E}[y], \\ \text{Cov}(z) &= \text{Cov}(2x + 3y, 2x + 3y) = 4\text{Cov}(x, x) + 12\text{Cov}(x, y) + 9\text{Cov}(y, y) \\ &= 4\text{Var}(x) + 12\text{Cov}(x, y) + 9\text{Var}(y), \end{aligned}$$

which completely characterizes the distribution of  $z$ .

For an  $m$  dimensional Gaussian vector  $\mathbf{x}$ , the elements of  $\mathbf{x}$  are **independent** if and only if the covariance matrix is diagonal.

For an  $m$ -dimensional Gaussian random vector  $\mathbf{x}$ , assuming that the covariance matrix  $K = \text{Cov}(\mathbf{x})$  is invertible, we have

$$p(\mathbf{x}) = \frac{1}{(2\pi)^{m/2} |K|^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T K^{-1}(\mathbf{x} - \boldsymbol{\mu})\right).$$

### 4.4 Maximum likelihood for Gaussian Random Vectors

Let  $\mathbf{z}$  be a Gaussian random vector of dimension  $d$  with mean  $\boldsymbol{\mu}$  and covariance matrix  $K$ . If  $K$  is invertible, the pdf of  $\mathbf{z}$  can be written as

$$p(\mathbf{z}|\boldsymbol{\mu}, K) = \frac{1}{\sqrt{(2\pi)^d |K|}} \exp\left(-\frac{1}{2}(\mathbf{z} - \boldsymbol{\mu})^T K^{-1}(\mathbf{z} - \boldsymbol{\mu})\right), \quad \boldsymbol{\mu} = \mathbb{E}[\mathbf{z}], \quad K = \mathbb{E}[(\mathbf{z} - \boldsymbol{\mu})(\mathbf{z} - \boldsymbol{\mu})^T],$$

where  $|K|$  is the determinant of  $K$ .

Given a set of  $n$  iid samples  $\mathcal{D} = \{\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_n\}$ , where each  $\mathbf{z}_i$  is a  $d$ -dimensional vector, how can we estimate  $\boldsymbol{\mu}$  and  $K$  using maximum likelihood? Estimating these quantities allows us to find the distribution. In particular, if we can view  $z_d$  as the output variable and  $z_1, \dots, z_{d-1}$  as input variables, then we can estimate  $z_d$  based on  $z_1, \dots, z_{d-1}$  as  $\mathbb{E}[z_d|z_1, \dots, z_{d-1}]$ .

To estimate  $\boldsymbol{\mu}$  and  $K$ , we write

$$\ell(\boldsymbol{\mu}, K) = \ln p(\mathcal{D}; \boldsymbol{\mu}, K) = \sum_{i=1}^n \ln p(\mathbf{z}_i; \boldsymbol{\mu}, K) \doteq \frac{n}{2} \ln |K^{-1}| - \frac{1}{2} \sum_{i=1}^n (\mathbf{z}_i - \boldsymbol{\mu})^T K^{-1}(\mathbf{z}_i - \boldsymbol{\mu}),$$

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

As seen in the appendix (99-Appendix.tex), for a symmetric matrix  $A$ , we have  $\frac{d}{dv}(\mathbf{y}^T A \mathbf{y}) = 2\mathbf{y}^T A \frac{d\mathbf{y}}{dv}$ . Hence,

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

Setting this equal to zero yields

$$\hat{\boldsymbol{\mu}}_{ML} = \bar{\mathbf{z}} = \frac{1}{n} \sum_{i=1}^n \mathbf{z}_i.$$

**Exercise 69.** Using the facts

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

prove that

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