

## CHAPTER 1

## INTRODUCTION

In this chapter we introduce and contrast the matricial and geometric formulations of the so-called general linear model and introduce some notational conventions.

## 1. Orientation

Recall the classical framework of the *general linear model (GLM)*. One is given an  $n$ -dimensional random vector  $\mathbf{Y}^{n \times 1} = (Y_1, \dots, Y_n)^T$ , perhaps multivariate normally distributed, with covariance matrix  $(\text{Cov}(Y_i, Y_j))^{n \times n} = \sigma^2 \mathbf{I}^{n \times n}$  and mean vector  $\boldsymbol{\mu}^{n \times 1} = E(\mathbf{Y}) = (EY_1, \dots, EY_n)^T$  of the form

$$\boldsymbol{\mu} = \mathbf{X}\boldsymbol{\beta},$$

where  $\mathbf{X}^{n \times p}$  is known and  $\sigma^2$  and  $\boldsymbol{\beta}^{p \times 1} = (\beta_1, \dots, \beta_p)^T$  are unknown; in addition, the  $\beta_i$ 's may be subject to linear constraints  $\mathbf{R}\boldsymbol{\beta} = \mathbf{0}$ , where  $\mathbf{R}^{c \times p}$  is known.  $\mathbf{X}$  is called the *design*, or *regression matrix*, and  $\boldsymbol{\beta}$  is called the *parameter vector*.

**1.1 Example.** In the classical *two-sample problem*, one has

$$\mathbf{X}^T = \left( \underbrace{\begin{matrix} 1 & 1 & \dots & 1 \\ 0 & 0 & \dots & 0 \end{matrix}}_{n_1 \text{ times}} \underbrace{\begin{matrix} 0 & 0 & \dots & 0 \\ 1 & 1 & \dots & 1 \end{matrix}}_{n_2 \text{ times}} \right) \quad \text{and} \quad \boldsymbol{\beta} = \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix},$$

that is,

$$E(Y_i) = \begin{cases} \mu_1, & \text{if } 1 \leq i \leq n_1, \\ \mu_2, & \text{if } n_1 < i \leq n_1 + n_2 = n. \end{cases} \quad \bullet$$

**1.2 Example.** In *simple linear regression*, one has

$$\mathbf{X}^T = \left( \begin{matrix} 1 & 1 & \dots & 1 \\ x_1 & x_2 & \dots & x_n \end{matrix} \right) \quad \text{and} \quad \boldsymbol{\beta} = \begin{pmatrix} a \\ b \end{pmatrix},$$

that is,

$$E(Y_i) = a + bx_i \quad \text{for } i = 1, \dots, n. \quad \bullet$$

Typical problems are the estimation of linear combinations of the  $\beta_i$ 's, testing that some such linear combinations are 0 (or some other prescribed value), and the estimation of  $\sigma^2$ .

**1.3 Example.** In the two-sample problem, one is often interested in estimating the difference  $\mu_2 - \mu_1$  or in testing the null hypothesis that  $\mu_1 = \mu_2$ . •

**1.4 Example.** In simple linear regression, one seeks estimates of the intercept  $a$  and slope  $b$  and may want to test, for example, the hypothesis that  $b = 0$  or the hypothesis that  $b = 1$ . •

If you have had some prior statistical training, you may well have already encountered the resolution of these problems. You may know, for example, that provided  $\mathbf{X}$  is of full rank and no linear constraints are imposed on  $\boldsymbol{\beta}$ , the *best (minimum variance) linear unbiased estimator (BLUE)* of  $\sum_{1 \leq i \leq p} c_i \beta_i$  is  $\sum_{1 \leq i \leq p} c_i \hat{\beta}_i$ , where

$$(\hat{\beta}_1, \dots, \hat{\beta}_p)^T = \mathbf{C}\mathbf{X}^T\mathbf{Y}, \quad \text{with } \mathbf{C} = \mathbf{A}^{-1}, \quad \mathbf{A} = \mathbf{X}^T\mathbf{X};$$

this is called the *Gauss-Markov theorem*.

In this book we will be studying the GLM from a geometric point of view, using linear algebra in place of matrix algebra. Although we will not reach any conclusions that could not be obtained solely by matrix techniques, the basic ideas will emerge more clearly. With the added intuitive feeling and mathematical insight this provides, one will be better able to understand old results and formulate and prove new ones.

From a geometric perspective, the GLM may be described as follows, using some terms that will be defined in subsequent chapters. One is given a *random vector*  $Y$  taking values in some given *inner product space*  $(V, \langle \cdot, \cdot \rangle)$ . It is assumed that  $Y$  has a *weakly spherical covariance operator* and the *mean*  $\mu$  of  $Y$  lies in a given *manifold*  $M$  of  $V$ ; for purposes of testing, it is further assumed that  $Y$  is *normally distributed*. One desires to estimate  $\mu$  (or *linear functionals* of  $\mu$ ) and to test hypotheses such as  $\mu \in M_0$ , where  $M_0$  is a given *submanifold* of  $M$ . The Gauss-Markov theorem says that the BLUE of the linear functional  $\psi(\mu)$  is  $\psi(\hat{\mu})$ , where  $\hat{\mu}$  is the *orthogonal projection* of  $Y$  onto  $M$ . As we will see, this geometric description of the problem encompasses the matricial formulation of the GLM not only as it is set out above (take, for example,  $V = \mathbb{R}^n$ ,  $\langle \cdot, \cdot \rangle = \text{dot-product}$ ,  $Y = \mathbf{Y}$ ,  $\mu = \boldsymbol{\mu}$ , and  $M =$  the subspace of  $\mathbb{R}^n$  spanned by the columns of the design matrix  $\mathbf{X}$ ), but also in cases where  $\mathbf{X}$  is of less than full rank and/or linear constraints are imposed on the  $\beta_i$ 's.

## 2. An illustrative example

To illustrate the differences between the matricial and geometric approaches, we compare the ways in which one establishes the independence of

$$\bar{Y} = \hat{\mu} = \frac{\sum_{1 \leq i \leq n} Y_i}{n} \quad \text{and} \quad s^2 = \hat{\sigma}^2 = \frac{1}{n-1} \sum_{1 \leq i \leq n} (Y_i - \bar{Y})^2$$

in the one-sample problem

$$\mathbf{Y}^{n \times 1} \sim N(\mu \mathbf{e}, \sigma^2 \mathbf{I}^{n \times n}) \quad \text{with} \quad \mathbf{e} = (1, 1, \dots, 1)^T. \quad (2.1)$$

(The vector  $\mathbf{e}$  is called the *equiangular vector*.)

The classical matrix proof, which uses some facts about multivariate normal distributions, runs like this. Let  $\mathbf{B}^{n \times n} = (b_{ij})$  be the matrix

$$\begin{pmatrix} \frac{1}{\sqrt{n}} & \frac{1}{\sqrt{n}} & \frac{1}{\sqrt{n}} & \frac{1}{\sqrt{n}} & \cdots & \frac{1}{\sqrt{n}} & \frac{1}{\sqrt{n}} \\ \frac{1}{\sqrt{2}} & \frac{-1}{\sqrt{2}} & 0 & 0 & \cdots & 0 & 0 \\ \frac{1}{\sqrt{6}} & \frac{1}{\sqrt{6}} & \frac{-2}{\sqrt{6}} & 0 & \cdots & 0 & 0 \\ \frac{1}{\sqrt{12}} & \frac{1}{\sqrt{12}} & \frac{1}{\sqrt{12}} & \frac{-3}{\sqrt{12}} & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \cdots & \vdots & \vdots \\ \frac{1}{\sqrt{n(n-1)}} & \frac{1}{\sqrt{n(n-1)}} & \frac{1}{\sqrt{n(n-1)}} & \frac{1}{\sqrt{n(n-1)}} & \cdots & \frac{1}{\sqrt{n(n-1)}} & \frac{-(n-1)}{\sqrt{n(n-1)}} \end{pmatrix}.$$

Note that the rows (and columns) of  $\mathbf{B}$  are orthonormal ( $\sum_{1 \leq k \leq n} b_{ik} b_{jk} = \delta_{ij} \equiv \begin{cases} 1, & \text{if } i = j \\ 0, & \text{if } i \neq j \end{cases}$ ) and that the first row is  $\frac{1}{\sqrt{n}} \mathbf{e}^T$ . Set

$$\mathbf{Z} = \mathbf{B}\mathbf{Y}.$$

Then

$$\mathbf{Z} \sim N(\boldsymbol{\nu}, \boldsymbol{\Sigma})$$

with

$$\boldsymbol{\nu} = \mathbf{B}(\mu \mathbf{e}) = \mu \mathbf{B}\mathbf{e} = (\sqrt{n} \mu, 0, \dots, 0)^T$$

and

$$\boldsymbol{\Sigma} = \mathbf{B}(\sigma^2 \mathbf{I})\mathbf{B}^T = \sigma^2 \mathbf{B}\mathbf{B}^T = \sigma^2 \mathbf{I};$$

that is,  $Z_1, Z_2, \dots, Z_n$  are independent normal random variables, each with variance  $\sigma^2$ ,  $E(Z_1) = \sqrt{n} \mu$ , and  $E(Z_j) = 0$  for  $2 \leq j \leq n$ . Moreover,

$$Z_1 = \frac{1}{\sqrt{n}} \sum_{1 \leq i \leq n} Y_i = \sqrt{n} \bar{Y}, \quad \text{or} \quad \bar{Y} = \frac{Z_1}{\sqrt{n}},$$

while

$$\begin{aligned} (n-1)s^2 &= \sum_{1 \leq i \leq n} (Y_i - \bar{Y})^2 = \sum_{1 \leq i \leq n} Y_i^2 - n\bar{Y}^2 \\ &= \sum_{1 \leq i \leq n} Y_i^2 - Z_1^2 = \sum_{1 \leq i \leq n} Z_i^2 - Z_1^2 = \sum_{2 \leq i \leq n} Z_i^2 \end{aligned} \quad (2.2)$$

because

$$\sum_{1 \leq i \leq n} Z_i^2 = \mathbf{Z}^T \mathbf{Z} = \mathbf{Y}^T \mathbf{B}^T \mathbf{B} \mathbf{Y} = \mathbf{Y}^T \mathbf{Y} = \sum_{1 \leq i \leq n} Y_i^2.$$

This gives the independence of  $\bar{Y}$  and  $s^2$ , and it is an easy step to get the marginal distributions:  $\bar{Y} \sim N(\mu, \sigma^2/n)$  and  $(n-1)s^2/\sigma^2 \sim \chi_{n-1}^2$ .

What is the nature of the transformation  $\mathbf{Z} = \mathbf{B}\mathbf{Y}$ ? Let  $\mathbf{b}_1 = \mathbf{e}/\sqrt{n}$ ,  $\mathbf{b}_2, \dots, \mathbf{b}_n$  denote the transposes of the rows of  $\mathbf{B}$ . The coordinates of  $\mathbf{Y} = \sum_{1 \leq j \leq n} C_j \mathbf{b}_j$  with respect to this new orthonormal basis for  $\mathbb{R}^n$  are given by

$$C_i = \mathbf{b}_i^T \mathbf{Y} = Z_i, \quad i = 1, \dots, n.$$

The effect of the change of coordinates  $\mathbf{Y} \rightarrow \mathbf{Z}$  is to split  $\mathbf{Y}$  into its components along, and orthogonal to, the equiangular vector  $\mathbf{e}$ .

Now I will show you the geometric proof, which uses some properties of (weakly) spherical normal random vectors taking values in an inner product space  $(V, \langle \cdot, \cdot \rangle)$ , here  $(\mathbb{R}^n, \text{dot-product})$ . The assumptions imply that  $\mathbf{Y}$  is spherical normally distributed about its mean  $E(\mathbf{Y})$  and  $E(\mathbf{Y})$  lies in the manifold  $M$  spanned by  $\mathbf{e}$ . Let  $P_M$  denote orthogonal projection onto  $M$  and  $Q_M$  orthogonal projection onto the orthogonal complement  $M^\perp$  of  $M$ . Basic distribution theory says that  $P_M \mathbf{Y}$  and  $Q_M \mathbf{Y}$  are independent. But

$$P_M \mathbf{Y} = \frac{\langle \mathbf{e}, \mathbf{Y} \rangle}{\langle \mathbf{e}, \mathbf{e} \rangle} \mathbf{e} = \bar{Y} \mathbf{e} \quad (2.3)$$

and

$$Q_M \mathbf{Y} = \mathbf{Y} - P_M \mathbf{Y} = \mathbf{Y} - \bar{Y} \mathbf{e} = (Y_1 - \bar{Y}, \dots, Y_n - \bar{Y})^T;$$

it follows that  $\bar{Y}$  and  $(n-1)s^2 = \sum_{1 \leq i \leq n} (Y_i - \bar{Y})^2 = \|Q_M \mathbf{Y}\|^2$  are independent. Again it is an easy matter to get the marginal distributions.

To my way of thinking, granted the technical apparatus, the second proof is clearer, being more to the point. The first proof does the same things, but (to the uninitiated) in an obscure manner.

### 3. Notational conventions

The chapters are organized into sections. Within each section of the current chapter, enumerated items are numbered consecutively in the form

$$(\text{section\_number.item\_number}).$$

References to items in a different chapter take the expanded form

$$(\text{chapter\_number.section\_number.item\_number}).$$

For example, (2.4) refers to the 4<sup>th</sup> numbered item (which may be an example, exercise, theorem, formula, or whatever) in the 2<sup>nd</sup> section of the current chapter, while (6.1.3) refers to the 3<sup>rd</sup> numbered item in the 1<sup>st</sup> section of the 6<sup>th</sup> chapter.

## SECTION 3. NOTATIONAL CONVENTIONS

5

Each exercise is assigned a difficulty level using the syntax

**Exercise** [ $d$ ],

where  $d$  is an integer in the range 1 to 5 — the larger is  $d$ , the harder the exercise. The value of  $d$  depends both on the intrinsic difficulty of the exercise and the length of time needed to write up the solution.

To help distinguish between the matricial and geometric points of view, matrices, including row and column vectors, are written in *italic boldface* type while linear transformations and elements of abstract vector spaces are written simply in *italic* type. We speak, for example, of the design matrix  $\mathbf{X}$  but of vectors  $v$  and  $w$  in an inner product space  $V$ .

The end of a proof is marked by a ■, of an example by a ●, of an exercise by a ◇, and of a part of problem set by a ○.

## CHAPTER 2

## TOPICS IN LINEAR ALGEBRA

In this chapter we discuss some topics from linear algebra that play a central role in the geometrical analysis of the GLM. The notion of orthogonal projection in an inner product space is introduced in Section 2.1 and studied in Section 2.2. A class of orthogonal decompositions that are useful in the design of experiments is studied in Section 2.3. The spectral representation of self-adjoint transformations is developed in Section 2.4. Linear and bilinear functionals are discussed in Section 2.5. The chapter closes with a problem set in Section 2.6, followed by an appendix containing a brief review of the basic definitions and facts from linear algebra with which we presume the reader is already familiar.

## 1. Orthogonal projections

Throughout this book we operate in the context of a finite-dimensional *inner product space*  $(V, \langle \cdot, \cdot \rangle)$  —  $V$  is a finite-dimensional vector space and  $\langle \cdot, \cdot \rangle: V \times V \rightarrow \mathbb{R}$  is an *inner product*:

- (i) (*positive-definiteness*)  $\langle x, x \rangle \geq 0$  for all  $x \in V$  and  $\langle x, x \rangle = 0$  if and only if  $x = 0$ .
- (ii) (*symmetry*)  $\langle x, y \rangle = \langle y, x \rangle$  for all  $x, y \in V$ .
- (iii) (*bilinearity*) For all  $c_1, c_2 \in \mathbb{R}$  and  $x_1, x_2, x, y_1, y_2, y \in V$ , one has

$$\begin{aligned}\langle c_1x_1 + c_2x_2, y \rangle &= c_1\langle x_1, y \rangle + c_2\langle x_2, y \rangle \\ \langle x, c_1y_1 + c_2y_2 \rangle &= c_1\langle x, y_1 \rangle + c_2\langle x, y_2 \rangle.\end{aligned}$$

The canonical example is  $V = \mathbb{R}^n$  endowed with the dot-product

$$\langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{x} \cdot \mathbf{y} = \sum_{1 \leq i \leq n} x_i y_i = \mathbf{x}^T \mathbf{y}$$

for  $\mathbf{x} = (x_1, \dots, x_n)^T$  and  $\mathbf{y} = (y_1, \dots, y_n)^T$ . Unless specifically stated to the contrary, we always view  $\mathbb{R}^n$  as endowed with the dot-product.

Two vectors  $x$  and  $y$  in  $V$  are said to be *perpendicular*, or *orthogonal* (with respect to  $\langle \cdot, \cdot \rangle$ ), if

$$\langle x, y \rangle = 0;$$

one writes

$$x \perp y.$$

The quantity

$$\|x\| = \sqrt{\langle x, x \rangle}$$

is called the *length*, or *norm*, of  $x$ . The squared length of the sum of two vectors is given by

$$\begin{aligned} \|x + y\|^2 &= \langle x + y, x + y \rangle = \langle x, x + y \rangle + \langle y, x + y \rangle \\ &= \|x\|^2 + 2\langle x, y \rangle + \|y\|^2, \end{aligned}$$

which reduces to the *Pythagorean theorem*,

$$\|x + y\|^2 = \|x\|^2 + \|y\|^2, \quad (1.1)$$

when  $x \perp y$ .

**1.2 Exercise** [1]. Let  $v_1$  and  $v_2$  be two nonzero vectors in  $\mathbb{R}^2$  and let  $\theta$  be the angle between them, measured counter-clockwise from  $v_1$  to  $v_2$ . Show that

$$\cos(\theta) = \frac{\langle v_1, v_2 \rangle}{\|v_1\| \|v_2\|}$$

and deduce that  $v_1 \perp v_2$  if and only if  $\theta = 90^\circ$  or  $270^\circ$ .

[Hint: Use the identity  $\cos(\theta_2 - \theta_1) = \cos(\theta_1)\cos(\theta_2) + \sin(\theta_1)\sin(\theta_2)$ .]  $\diamond$

**1.3 Exercise** [3]. Let  $x$ ,  $y$ , and  $z$  be the vectors in  $\mathbb{R}^n$  given by

$$x_i = 1, \quad y_i = i - \frac{n+1}{2}, \quad \text{and} \quad z_i = \left(i - \frac{n+1}{2}\right)^2 - \frac{n^2-1}{12} \quad \text{for } 1 \leq i \leq n. \quad (1.4)$$

Show that  $x$ ,  $y$ , and  $z$  are mutually orthogonal and span the same subspace as do  $(1, 1, \dots, 1)^T$ ,  $(1, 2, \dots, n)^T$ , and  $(1, 4, \dots, n^2)^T$ . Exhibit  $x$ ,  $y$ , and  $z$  explicitly for  $n = 5$ .  $\diamond$

**1.5 Exercise** [1]. Let  $\mathcal{O}: V \rightarrow V$  be a linear transformation. Show that  $\mathcal{O}$  preserves lengths, that is,

$$\|\mathcal{O}x\| = \|x\| \quad \text{for all } x \in V$$

if and only if it preserves inner products, that is,

$$\langle \mathcal{O}x, \mathcal{O}y \rangle = \langle x, y \rangle \quad \text{for all } x, y \in V.$$

Such a transformation is said to be *orthogonal*.

[Hint: Observe that

$$\langle x, y \rangle = \frac{\|x + y\|^2 - \|x\|^2 - \|y\|^2}{2} \quad (1.6)$$

for  $x, y \in V$ .]  $\diamond$

**1.7 Exercise** [2]. (1) Let  $v_1$  and  $v_2$  be elements of  $V$ . Show that  $v_1 = v_2$  if and only if  $\langle v_1, w \rangle = \langle v_2, w \rangle$  for each  $w \in V$ , or just for each  $w$  in some basis for  $V$ . (2) Let  $T_1$  and  $T_2$  be two linear transformations of  $V$ . Show that  $T_1 = T_2$  if and only if  $\langle v, T_1 w \rangle = \langle v, T_2 w \rangle$  for each  $v$  and  $w$  in  $V$ , or just for each  $v$  and  $w$  in some basis for  $V$ .  $\diamond$

As intimated in the Introduction, the notion of orthogonal projection onto subspaces of  $V$  plays a key role in the study of the GLM. We begin our study of projections with the following seminal result.

**1.8 Theorem.** *Suppose that  $M$  is a subspace of  $V$  and that  $x \in V$ . There is exactly one vector  $m \in M$  such that the residual  $x - m$  is orthogonal to  $M$ :*

$$(x - m) \perp y \quad \text{for all } y \in M, \quad (1.9)$$

or, equivalently, such that  $m$  is closest to  $x$ :

$$\|x - m\| = \inf\{\|x - y\| : y \in M\}. \quad (1.10)$$

The proof will be given shortly. The unique  $m \in M$  such that (1.9) and (1.10) hold is called the *orthogonal projection of  $x$  onto  $M$* , written  $P_M x$ , and the mapping  $P_M$  that sends  $x \in V$  to  $P_M x$  is called *orthogonal projection onto  $M$* . In the context of  $\mathbb{R}^n$  with  $\mathbf{x} = (x_i)$  and  $\mathbf{m} = (m_i)$ , (1.9) reads

$$\sum_{1 \leq i \leq n} (x_i - m_i) y_i = 0 \quad \text{for all } \mathbf{y} = (y_i) \in M,$$

while (1.10) is the *least squares characterization* of  $\mathbf{m}$ :

$$\sum_{1 \leq i \leq n} (x_i - m_i)^2 = \inf\left\{ \sum_{1 \leq i \leq n} (x_i - y_i)^2 : \mathbf{y} \in M \right\}.$$

**1.11 Exercise** [1]. Let  $M$  be a subspace of  $V$ . (a) Show that

$$\|x - P_M x\|^2 = \|x\|^2 - \|P_M x\|^2 \quad (1.12)$$

for all  $x \in V$ . (b) Deduce that

$$\|P_M x\| \leq \|x\| \quad (1.13)$$

for all  $x \in V$ , with equality holding if and only if  $x \in M$ .  $\diamond$

**Proof of Theorem 1.8.** (1.9) *implies* (1.10): Suppose  $m \in M$  satisfies (1.9). Then for all  $y \in M$  the Pythagorean theorem gives

$$\|x - y\|^2 = \|(x - m) + (m - y)\|^2 = \|x - m\|^2 + \|m - y\|^2,$$

so  $m$  satisfies (1.10).

(1.10) *implies* (1.9): Suppose  $m \in M$  satisfies (1.10). Then, for any  $0 \neq y \in M$  and any  $\delta \in \mathbb{R}$ ,

$$\|x - m\|^2 \leq \|x - m + \delta y\|^2 = \|x - m\|^2 + 2\delta \langle x - m, y \rangle + \delta^2 \|y\|^2;$$

this relation forces  $\langle x - m, y \rangle = 0$ .



*Uniqueness:* If

$$m_1 + y_1 = x = m_2 + y_2$$

with  $m_i \in M$  and  $y_i \perp M$  for  $i = 1$  and  $2$ , then the vector

$$m_1 - m_2 = y_2 - y_1$$

lies in  $M$  and is perpendicular to  $M$ , and so it is perpendicular to itself:

$$0 = \langle m_1 - m_2, m_1 - m_2 \rangle = \|m_1 - m_2\|^2,$$

whence  $m_1 = m_2$  by the positive-definiteness of  $\langle \cdot, \cdot \rangle$ .

*Existence:* We will show in a moment that  $M$  has a basis  $m_1, \dots, m_k$  consisting of mutually orthogonal vectors. For any such basis the generic  $m = \sum_{1 \leq j \leq k} c_j m_j$  in  $M$  satisfies

$$(x - m) \perp M$$

if and only if  $x - m$  is orthogonal to each  $m_i$ , that is, if and only if

$$\langle m_i, x \rangle = \langle m_i, \sum_{1 \leq j \leq k} c_j m_j \rangle = \sum_{1 \leq j \leq k} \langle m_i, m_j \rangle c_j = \langle m_i, m_i \rangle c_i$$

for  $1 \leq i \leq k$ . It follows that we can take

$$P_M x = \sum_{1 \leq i \leq k} \frac{\langle m_i, x \rangle}{\langle m_i, m_i \rangle} m_i. \tag{1.14}$$

To produce an orthogonal basis for  $M$ , let  $m_1^*, \dots, m_k^*$  be any basis for  $M$  and inductively define new basis vectors  $m_1, \dots, m_k$  by the recipe  $m_1 = m_1^*$  and

$$\begin{aligned} m_j &= m_j^* - P_{[m_1^*, \dots, m_{j-1}^*]} m_j^* = m_j^* - P_{[m_1, \dots, m_{j-1}]} m_j^* \\ &= m_j^* - \sum_{1 \leq i \leq j-1} \frac{\langle m_j^*, m_i \rangle}{\langle m_i, m_i \rangle} m_i \end{aligned} \tag{1.15}$$

for  $j = 2, \dots, k$ ; here  $[m_1^*, \dots, m_{j-1}^*]$  denotes the span of  $m_1^*, \dots, m_{j-1}^*$  and  $[m_1, \dots, m_{j-1}]$  denotes the (identical) span of  $m_1, \dots, m_{j-1}$  (see Subsection 2.6.1). ■

The recursive scheme for cranking out the  $m_j$ 's above is called *Gram-Schmidt orthogonalization*. As a special case of (1.14) we have the following simple, yet key, formula for projecting onto a one-dimensional space:

$$P_{[m]} x = \frac{\langle x, m \rangle}{\langle m, m \rangle} m \quad \text{for } m \neq 0. \tag{1.16}$$

**1.17 Example.** In the context of  $\mathbb{R}^n$  take  $\mathbf{m} = \mathbf{e} \equiv (1, \dots, 1)^T$ . Formula (1.16) then reads

$$P_{[\mathbf{e}]} \mathbf{x} = \frac{\langle \mathbf{x}, \mathbf{e} \rangle}{\langle \mathbf{e}, \mathbf{e} \rangle} \mathbf{e} = \frac{\sum_i x_i}{\sum_i 1} \mathbf{e} = \bar{x} \mathbf{e} = (\bar{x}, \dots, \bar{x})^T;$$

we used this result in the Introduction (see (1.2.3)). Formula (1.12) reads

$$\sum_i (x_i - \bar{x})^2 = \|\mathbf{x} - \bar{x}\mathbf{e}\|^2 = \|\mathbf{x}\|^2 - \|\bar{x}\mathbf{e}\|^2 = \sum_i x_i^2 - n\bar{x}^2;$$

this is just the computing formula for  $(n-1)s^2$  used in (1.2.2). According to (1.10),  $c = \bar{x}$  minimizes the sum of squares  $\sum_i (x_i - c)^2$ . •

**1.18 Exercise** [2]. Use the preceding techniques to compute  $P_M \mathbf{y}$  for  $\mathbf{y} \in \mathbb{R}^n$  and  $M = [\mathbf{e}, (x_1, \dots, x_n)^T]$ , the manifold spanned by the columns of the design matrix for simple linear regression. ◇

**1.19 Exercise** [3]. Show that the transposes of the rows of the matrix  $\mathbf{B}$  of Section 1.2 result from first applying the Gram-Schmidt orthogonalization scheme to the vectors  $\mathbf{e}$ ,  $(1, -1, 0, \dots, 0)^T$ ,  $(0, 1, -1, 0, \dots, 0)^T$ ,  $\dots$ ,  $(0, \dots, 0, 1, -1)^T$  in  $\mathbb{R}^n$  and then normalizing to unit length. ◇

**1.20 Exercise** [2]. Suppose  $x, y \in V$ . Prove the *Cauchy-Schwarz inequality*:

$$|\langle x, y \rangle| \leq \|x\| \|y\|, \quad (1.21)$$

with equality if and only if  $x$  and  $y$  are linearly dependent. Deduce *Minkowski's inequality*:

$$\|x + y\| \leq \|x\| + \|y\|. \quad (1.22)$$

[Hint: For (1.21), take  $m = y$  in (1.16) and use part (b) of Exercise 1.11.] ◇

**1.23 Exercise** [4]. Define  $d: V \times V \rightarrow \mathbb{R}$  by  $d(x, y) = \|y - x\|$ . Show that  $d$  is a metric on  $V$  such that the set  $\{x \in V : \|x\| = 1\}$  is compact. ◇

**1.24 Exercise** [3]. Show that for any two subspaces  $M$  and  $N$  of  $V$ ,

$$\sup\{\|P_M x\| : x \in N \text{ and } \|x\| = 1\} \leq 1, \quad (1.25)$$

with equality holding if and only if  $M$  and  $N$  have a nonzero vector in common.

[Hint: A continuous real-valued function on a compact set attains its maximum.] ◇

To close this section we generalize (1.14) to cover the case of an arbitrary basis for  $M$ . Suppose then that the basis vectors  $m_1, \dots, m_k$  are not necessarily orthogonal and let  $x \in V$ . As in the derivation of (1.14),

$$P_M x = \sum_j c_j m_j$$

is determined by the condition

$$m_i \perp (x - P_M x) \quad \text{for } i = 1, \dots, k,$$

that is, by the so-called *normal equations*

$$\sum_{1 \leq j \leq k} \langle m_i, m_j \rangle c_j = \langle m_i, x \rangle, \quad i = 1, \dots, k. \quad (1.26)$$