

## *Research Article*

# **From Equivalent Linear Equations to Gauss-Markov Theorem**

**Czesław Stępnia**

*Institute of Mathematics, University of Rzeszów, Rejtana 16 A, 35-959 Rzeszów, Poland*

Correspondence should be addressed to Czesław Stępnia, [stepniak@umcs.lublin.pl](mailto:stepniak@umcs.lublin.pl)

Received 17 December 2009; Revised 20 May 2010; Accepted 27 June 2010

Academic Editor: Andrei Volodin

Copyright © 2010 Czesław Stępnia. This is an open access article distributed under the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

Gauss-Markov theorem reduces linear unbiased estimation to the Least Squares Solution of inconsistent linear equations while the normal equations reduce the second one to the usual solution of consistent linear equations. It is rather surprising that the second algebraic result is usually derived in a differential way. To avoid this dissonance, we state and use an auxiliary result on equivalence of two systems of linear equations. This places us in a convenient position to attack on the main problems in the Gauss-Markov model in an easy way.

## **1. Introduction**

The Gauss-Markov theorem is the most classical achievement in statistics. Its role in statistics is comparable with that of the Pythagorean theorem in geometry. In fact, there are close relations between both of them.

The Gauss-Markov theorem is presented in many books and derived in many ways. The most popular approaches involve

- (i) geometry (cf., Kruskal [1, 2]),
- (ii) differential calculus (cf., Scheffé [3], Rao [4]),
- (iii) generalized inverse matrices (cf., Rao and Mitra [5], Bapat [6]),
- (iv) projection operators (see Seber [7]).

We presume that such a big market has many clients. This paper is intended for some of them. Our consideration is straightforward and self-contained. Moreover, it needs only moderate prerequisites.

The main tool used in this paper is equivalence of two systems of linear equations.

## 2. Preliminaries

For any matrix  $\mathbf{A}$  of  $n \times p$  define the sets

$$\begin{aligned}\mathcal{R}(\mathbf{A}) &= \{\mathbf{a} \in R^n : \mathbf{a} = \mathbf{A}\mathbf{x} \text{ for some } \mathbf{x} \in R^p\} \text{ (i.e., the range of } \mathbf{A}), \\ \mathcal{N}(\mathbf{A}) &= \{\mathbf{x} \in R^p : \mathbf{A}\mathbf{x} = \mathbf{0}\} \text{ (i.e., the kernel of } \mathbf{A}).\end{aligned}\quad (2.1)$$

We note that

$$\mathbf{a}^T \mathbf{b} = 0, \quad \forall \mathbf{a} \in \mathcal{R}(\mathbf{A}), \mathbf{b} \in \mathcal{N}(\mathbf{A}^T). \quad (2.2)$$

It is clear that the range  $\mathcal{R}(\mathbf{A})$  constitutes  $r$ -dimensional linear space in  $R^n$  spanned by the columns of  $\mathbf{A}$ , where  $r = \text{rank}(\mathbf{A})$ , while  $\mathcal{N}(\mathbf{A}^T)$  constitutes  $(n - r)$ -dimensional space of all vectors being orthogonal to any vector in  $\mathcal{R}(\mathbf{A})$  relative to the usual inner product  $(\mathbf{a}, \mathbf{b}) = \mathbf{a}^T \mathbf{b}$ . Thus, any vector  $\mathbf{y} \in R^n$  may be presented in the form

$$\mathbf{y} = \mathbf{y}_1 + \mathbf{y}_2, \quad \text{where } \mathbf{y}_1 \in \mathcal{R}(\mathbf{A}), \mathbf{y}_2 \in \mathcal{N}(\mathbf{A}^T) \text{ are orthogonal.} \quad (2.3)$$

Since  $\mathbf{A}^T \mathbf{A}\mathbf{x} = \mathbf{0}$  if and only if  $\mathbf{x}^T \mathbf{A}^T \mathbf{A}\mathbf{x} = 0$  and, hence,  $\mathbf{A}\mathbf{x} = \mathbf{0}$ , we get  $\mathcal{N}(\mathbf{A}\mathbf{A}^T) = \mathcal{N}(\mathbf{A}^T)$ .

Denote by  $\mathbf{P} = \mathbf{P}_A$  the linear operator from  $R^n$  onto  $\mathcal{R}(\mathbf{A})$  defined by

$$\mathbf{P}\mathbf{y} = \begin{cases} \mathbf{y}, & \text{if } \mathbf{y} \in \mathcal{R}(\mathbf{A}), \\ \mathbf{0}, & \text{if } \mathbf{y} \in \mathcal{N}(\mathbf{A}^T) \end{cases} \quad (2.4)$$

(i.e., the *orthogonal projector* onto  $\mathcal{R}(\mathbf{A})$ ). It follows from definition (2.4) that  $\mathbf{P}\mathbf{P} = \mathbf{P}$ . The following lemma (see [8]) will be a key tool in the further consideration.

**Lemma 2.1.** *For any matrix  $\mathbf{A}$  and for any vector  $\mathbf{b} \in \mathcal{R}(\mathbf{A})$ , the following are equivalent:*

- (i)  $\mathbf{A}\mathbf{x} = \mathbf{b}$ ,
- (ii)  $\mathbf{A}^T \mathbf{A}\mathbf{x} = \mathbf{A}^T \mathbf{b}$ .

*Proof.* (i) $\Rightarrow$ (ii) is evident (without any condition on  $\mathbf{b}$ ).

(ii) $\Rightarrow$ (i). By the assumption that  $\mathbf{b} \in \mathcal{R}(\mathbf{A})$ , we get  $\mathbf{b} = \mathbf{A}\mathbf{c}$  for some  $\mathbf{c}$ . Thus, (ii) reduces to  $\mathbf{A}^T \mathbf{A}\mathbf{x} = \mathbf{A}^T \mathbf{A}\mathbf{c}$  and its general solution is  $\mathbf{x} = \mathbf{c} + \mathbf{x}_0$ , where  $\mathbf{x}_0 \in \mathcal{N}(\mathbf{A}^T \mathbf{A}) = \mathcal{N}(\mathbf{A})$ . Therefore,  $\mathbf{x}$  is a solution of (i).  $\square$

*Remark 2.2.* The assumption that  $\mathbf{b} \in \mathcal{R}(\mathbf{A})$  in Lemma 2.1 is essential. To see this, let us set

$$\mathbf{A} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 1 & 0 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} 1 \\ 0 \\ -1 \end{bmatrix}. \quad (2.5)$$

Then,  $\mathbf{A}^T \mathbf{A} = \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix}$  and  $\mathbf{A}^T \mathbf{b} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$ . Thus,  $\mathbf{A}^T \mathbf{A}\mathbf{x} = \mathbf{A}^T \mathbf{b}$  has a solution  $\mathbf{x} = [0, 0]^T$ , while  $\mathbf{A}\mathbf{x} = \mathbf{b}$  is inconsistent.

### 3. Least Squares Solution

For any matrix  $\mathbf{A}$  of  $n \times p$  and for any vector  $\mathbf{b} \in R^n$ , consider the linear equation

$$\mathbf{Ax} = \mathbf{b}. \quad (3.1)$$

Equation (3.1) may be consistent (if  $\mathbf{b} \in \mathcal{R}(\mathbf{A})$ ) or inconsistent (if not). In the second case we are seeking for such  $\mathbf{x}$  that the residual vector  $\mathbf{b} - \mathbf{Ax}$  be as small as possible.

*Definition 3.1.* Any vector  $\hat{\mathbf{x}} \in R^p$  is said to be the Least Squares Solution (LSS) of (3.1) if

$$(\mathbf{b} - \mathbf{A}\hat{\mathbf{x}})^T(\mathbf{b} - \mathbf{A}\hat{\mathbf{x}}) \leq (\mathbf{b} - \mathbf{Ax})^T(\mathbf{b} - \mathbf{Ax}), \quad \text{for any } \mathbf{x} \in R^p. \quad (3.2)$$

The following theorem shows that this definition is not empty and reduces the LSS of an inconsistent equation (3.1) to the ordinary solution of a consistent one.

**Theorem 3.2.** (a) Equation (3.1) has at least one LSS.

(b) Vector  $\mathbf{x} \in R^p$  is an LSS of (3.1) if and only if

$$\mathbf{A}^T \mathbf{Ax} = \mathbf{A}^T \mathbf{b}. \quad (3.3)$$

(c) Condition (3.3) is equivalent to

$$\mathbf{Ax} = \mathbf{Pb}, \quad (3.4)$$

where  $\mathbf{P} = \mathbf{P}_A$  is the orthogonal projector onto  $\mathcal{R}(\mathbf{A})$  defined by (2.4).

(d) General solution of (3.3) may be presented in the form  $\mathbf{x} = \mathbf{x}_0 + \mathbf{x}_1$ , where  $\mathbf{x}_0$  is a particular solution, while  $\mathbf{x}_1 \in \mathcal{N}(\mathbf{A})$ .

*Remark 3.3.* In the statistical literature, (3.3) is said to be normal.

*Proof.* By properties of the projector  $\mathbf{P}$ , we get

$$\begin{aligned} (\mathbf{b} - \mathbf{Ax})^T(\mathbf{b} - \mathbf{Ax}) &= [\mathbf{Pb} + (\mathbf{I} - \mathbf{P})\mathbf{b} - \mathbf{Ax}]^T [\mathbf{Pb} + (\mathbf{I} - \mathbf{P})\mathbf{b} - \mathbf{Ax}] \\ &= (\mathbf{Pb} - \mathbf{Ax})^T(\mathbf{Pb} - \mathbf{Ax}) + [(\mathbf{I} - \mathbf{P})\mathbf{b}]^T [(\mathbf{I} - \mathbf{P})\mathbf{b}] \\ &= (\mathbf{Pb} - \mathbf{Ax})^T(\mathbf{Pb} - \mathbf{Ax}) + \mathbf{b}^T(\mathbf{I} - \mathbf{P})\mathbf{b} \\ &\geq \mathbf{b}^T(\mathbf{I} - \mathbf{P})\mathbf{b} \end{aligned} \quad (3.5)$$

with the equality if and only if (3.4) holds. Moreover, by definition of  $\mathbf{P}$ , (3.4) is consistent and, by Lemma 2.1, it is equivalent to (3.3).

Statement (d) follows directly from definition of kernel.  $\square$

#### 4. Gauss-Markov Model and Gauss-Markov Theorem

Let  $\mathbf{y}$  be an arbitrary random vector in  $R^n$  with finite second moment  $E(\mathbf{y}^T \mathbf{y})$ . Then there exist a unique vector  $\boldsymbol{\mu} \in R^n$  and a unique symmetric nonnegative definite matrix  $\mathbf{V}$  of  $n \times n$  such that

$$\begin{aligned} E(\mathbf{a}^T \mathbf{y}) &= \mathbf{a}^T \boldsymbol{\mu}, \\ \text{Cov}(\mathbf{A}^T \mathbf{y}, \mathbf{B}^T \mathbf{y}) &= \mathbf{A}^T \mathbf{V} \mathbf{B} \end{aligned} \quad (4.1)$$

for all vectors  $\mathbf{a} \in R^n$  and all matrices  $\mathbf{A}$  and  $\mathbf{B}$  of  $n$  rows. Traditionally, such  $\boldsymbol{\mu}$  and  $\mathbf{V}$  are called the expectation and the dispersion of the random vector  $\mathbf{y}$ .

As usual, we will assume that  $\boldsymbol{\mu}$  and  $\mathbf{V}$  have the representations

$$\begin{aligned} \boldsymbol{\mu} &= \mathbf{X}\boldsymbol{\beta}, \\ \mathbf{V} &= \sigma^2 \mathbf{I}_n, \end{aligned} \quad (4.2)$$

where  $\mathbf{X}$  is a given matrix of  $n \times p$  while  $\boldsymbol{\beta} \in R^p$  and  $\sigma^2 > 0$  are unknown parameters. We will refer to the structure  $(\mathbf{y}, \mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I}_n)$  as to the standard Gauss-Markov model. In the context of the model we will consider unbiased estimation of the parametric vector  $\boldsymbol{\Psi} = \mathbf{C}^T \boldsymbol{\beta}$ , where  $\mathbf{C}$  is of  $p \times q$ , by estimators of the form  $\hat{\boldsymbol{\Psi}} = \mathbf{D}^T \mathbf{y}$ , where  $\mathbf{D}$  is of  $n \times p$  matrix. Since  $\mathbf{D}^T \mathbf{y}$  is unbiased if and only if  $\mathbf{D}^T \mathbf{X}\boldsymbol{\beta} = \mathbf{C}^T \boldsymbol{\beta}$  for all  $\boldsymbol{\beta}$ ,  $\mathbf{C}^T \boldsymbol{\beta}$  is estimable if and only if

$$\mathbf{C}^T = \mathbf{D}^T \mathbf{X}, \quad \text{for some } \mathbf{D}. \quad (4.3)$$

Without loss of generality, we may and will assume that  $\mathcal{R}(\mathbf{D}) \subseteq \mathcal{R}(\mathbf{X})$ . We note that such a matrix  $\mathbf{D}$  is uniquely determined by  $\mathbf{C}$ .

The well-known Gauss-Markov theorem provides a constructive way for estimation of the function  $\boldsymbol{\Psi}$ . It is based on a solution of the normal equation  $\mathbf{X}^T \mathbf{X}\boldsymbol{\beta} = \mathbf{X}^T \mathbf{y}$  which plays the role of the estimator for  $\boldsymbol{\beta}$ .

**Theorem 4.1.** *For any estimable  $\boldsymbol{\Psi} = \mathbf{C}^T \boldsymbol{\beta}$  in the standard Gauss-Markov model  $(\mathbf{y}, \mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I}_n)$ , there exists a unique linear unbiased estimator with minimal dispersion. This estimator, called the Least Squares Estimator (LSE) of  $\boldsymbol{\Psi}$ , may be presented in the form  $\mathbf{C}^T \hat{\boldsymbol{\beta}}$ , where  $\hat{\boldsymbol{\beta}}$  is an arbitrary LSS of  $\mathbf{X}\boldsymbol{\beta} = \mathbf{y}$  or, equivalently, it is a solution of the normal equation*

$$\mathbf{X}^T \mathbf{X} \hat{\boldsymbol{\beta}} = \mathbf{X}^T \mathbf{y}. \quad (4.4)$$

*Proof.* By Theorem 3.2 the condition  $\mathbf{X}^T \mathbf{X}\boldsymbol{\beta} = \mathbf{X}^T \mathbf{y}$  is equivalent to  $\mathbf{P}_{\mathbf{X}} \mathbf{y} = \mathbf{X}\boldsymbol{\beta} = \mathbf{E}\mathbf{y}$ . Therefore, by (4.3), for any estimable  $\boldsymbol{\Psi} = \mathbf{C}^T \boldsymbol{\beta}$  and for any solution  $\hat{\boldsymbol{\beta}}$  of (4.4), the statistic  $\mathbf{C}^T \hat{\boldsymbol{\beta}} = \mathbf{D}^T \mathbf{X} \hat{\boldsymbol{\beta}}$  is unbiased. On the other hand,

$$E(\mathbf{D}_1^T \mathbf{y}) = E(\mathbf{D}_2^T \mathbf{y}), \quad \text{iff } \mathcal{R}(\mathbf{D}_1 - \mathbf{D}_2) \subseteq \mathcal{N}(\mathbf{X}^T). \quad (4.5)$$

Hence, any unbiased estimator of  $\Psi$  may be presented in the form  $\mathbf{D}^T \mathbf{X} \hat{\beta} + \mathbf{B}^T \mathbf{y}$ , where the first component is the LSE of  $\Psi$  while  $\mathbf{B}^T \mathbf{D} = \mathbf{B}^T \mathbf{X} = \mathbf{0}$ . In particular the components are not correlated. Therefore, the variance of the sum is greater than the variance of the LSE  $\mathbf{D}^T \mathbf{X} \hat{\beta}$ , unless  $\mathbf{B}^T \mathbf{y} \neq \mathbf{0}$ . Moreover, by Theorem 3.2(d) this estimator is invariant with respect to the choice of the LSS  $\hat{\beta}$ . In consequence, the LSE of the function  $\Psi$  is unique.  $\square$

## Acknowledgment

Thanks are due to a reviewer for his comments leading to the improvement in the presentation of this paper.

## References

- [1] W. Kruskal, "The coordinate-free approach to Gauss-Markov estimation, and its application to missing and extra observations," in *Proceedings of 4th Berkeley Symposium on Mathematical Statistics and Probability*, vol. 1, pp. 435–451, University of California Press, Berkeley, Calif, USA, 1961.
- [2] W. Kruskal, "When are Gauss-Markov and least squares estimators identical? A coordinate-free approach," *Annals of Mathematical Statistics*, vol. 39, pp. 70–75, 1968.
- [3] H. Scheffé, *The Analysis of Variance*, John Wiley & Sons, New York, NY, USA, 1959.
- [4] C. R. Rao, *Linear Statistical Inference and Its Applications*, John Wiley & Sons, New York, NY, USA, 2nd edition, 1973.
- [5] C. R. Rao and S. K. Mitra, *Generalized Inverse of Matrices and Its Applications*, John Wiley & Sons, New York, NY, USA, 1971.
- [6] R. B. Bapat, *Linear Algebra and Linear Models*, Universitext, Springer, New York, NY, USA, 2nd edition, 2000.
- [7] G. A. F. Seber, *Linear Regression Analysis*, John Wiley & Sons, New York, NY, USA, 1977.
- [8] C. Stepniak, "Through a generalized inverse," *Demonstratio Mathematica*, vol. 41, no. 2, pp. 291–296, 2008.