# The Gauss-Markov Theorem: Beyond the BLUE

Leigh J. Halliwell, FCAS, MAAA

**Abstract:** Until now the Gauss-Markov theorem has been the handmaid of least squares; it has served as a proof that the least-squares method produces the Best Linear Unbiased Estimator (BLUE). This theoretical paper shows that it can be, and should be, reformulated as the solution to the problem of the minimization of a quadratic form subject to a linear constraint. The whole theory of linear statistical modeling, from basic to complicated, receives a clean and efficient development on the basis of this reformulation; estimates and predictions based thereon are BLUE from the start, rather than BLUE by subsequent proof. With an intermediate-level background in matrix algebra the reader will understand the frequent interpretations of this development in terms of an n-dimensional projective geometry. Because this paper elevates BLUE to its true role, "Beyond the BLUE" really means "To the True BLUE."

**Keywords**: Gauss-Markov, BLUE, linear model, projection, distance metric

## 1. INTRODUCTION

The many treatments of the Gauss-Markov theorem (e.g., Judge [1988, 202-206], Halliwell [2007, Appendix B], and Wikipedia) lead one to believe that the theorem is no more than a proof that $\hat{\boldsymbol{\beta}} = (X'X)^{-1}X'\mathbf{y}$ is best linear unbiased estimator (BLUE) of β in the model $\mathbf{y} = X\beta + \mathbf{e}$, where $Var[\mathbf{e}] = \sigma^2 I$. In this capacity the theorem is impressive enough; however, with a little abstraction it becomes much more, as we shall see in the following eleven sections.

## 2. STATEMENT OF THE THEOREM AND ITS PROOF

The Gauss-Markov theorem is essentially the solution to a constrained-optimization problem, more exactly, to the problem of minimizing a quadratic form subject to a linear constraint. Here is our formulation of the theorem:

<u>The Gauss-Markov Theorem</u>: If symmetric $\Sigma_{n \times n}$ is positive-definite and $A_{m \times n}$ is of full-row rank, then $\Phi(W) = W'\Sigma^{-1}W$ can be minimized subject to the linear constraint $A_{m \times n} W_{n \times p} = B_{m \times p}$. The value $W^* = \Sigma A'(A\Sigma A')^{-1}B$ uniquely minimizes $\Phi$ at $\Phi(W^*) = B'(A\Sigma A')^{-1}B$.

To prove the theorem, we take for granted two theorems about positive-definite matrices.[1]
First, positive-definite matrices have inverses; the inverses also are positive-definite.
Therefore, symmetric $\Sigma^{-1}$ exists, and is positive-definite. Second, if $A_{m \times n}$ is of full-row
rank and $T_{n \times n}$ is positive-definite, then $ATA'$ is positive-definite. From these it follows that
$A \Sigma A'$ is positive-definite and invertible; hence, $W^* = \Sigma A' (A \Sigma A')^{-1} B$ exists. Moreover,
$W^*$ satisfies the constraint, since $A W^* = A \Sigma A' (A \Sigma A')^{-1} B = I_m B = B$.

Now if $W_1$ satisfies the constraint, then:

$$
\begin{aligned}
W_1' \Sigma^{-1} W^* &= W_1' \Sigma^{-1} \Sigma A' (A \Sigma A')^{-1} B \\
&= W_1' A' (A \Sigma A')^{-1} B \\
&= (A W_1)' (A \Sigma A')^{-1} B \\
&= B' (A \Sigma A')^{-1} B
\end{aligned}
$$

And since $W^*$ is an allowable instance of $W_1$, we have the following chain of equalities:

$$
W_1' \Sigma^{-1} W^* = B' (A \Sigma A')^{-1} B = W^{*'} \Sigma^{-1} W^* = \left( W^{*'} \Sigma^{-1} W^* \right)' = \left( W_1' \Sigma^{-1} W^* \right)' = W^{*'} \Sigma^{-1} W_1
$$

As the heart of the Gauss-Markov proof:

$$
\begin{aligned}
\Phi(W_1) - \Phi(W^*) &= W_1' \Sigma^{-1} W_1 - W^{*'} \Sigma^{-1} W^* \\
&= W_1' \Sigma^{-1} W_1 - W^{*'} \Sigma^{-1} W^* - W^{*'} \Sigma^{-1} W^* + W^{*'} \Sigma^{-1} W^* \\
&= W_1' \Sigma^{-1} W_1 - W_1' \Sigma^{-1} W^* - W^{*'} \Sigma^{-1} W_1 + W^{*'} \Sigma^{-1} W^* \\
&= (W_1 - W^*)' \Sigma^{-1} (W_1 - W^*) \\
&\geq 0_{p \times p}
\end{aligned}
$$

---

[1] For a review of positive-definite and non-negative-definite (or positive-semi-definite) matrices see Judge [1988, Appendix A.14] and Halliwell [1997, Appendix A].

The last line is to be taken in a matrix-definite sense, viz., that the difference $\Phi(W_1) - \Phi(W^*)$ is the non-negative-definite matrix $(W_1 - W^*)'\Sigma^{-1}(W_1 - W^*)$. And because $\Sigma^{-1}$ is positive-definite, the difference equals the zero matrix $(0_{p \times p})$ if and only if $W_1 = W^*$. Therefore, $W^* = \Sigma A'(A\Sigma A')^{-1}B$ uniquely minimizes $\Phi(W) = W'\Sigma^{-1}W$ subject to $AW = B$. Furthermore, the minimum is $\Phi(W^*) = B'(A\Sigma A')^{-1}B$.

## 3. GEOMETRICAL INTERPRETATION WITH A DISTANCE METRIC

A geometrical interpretation of the theorem will prove helpful. Again, let $W_1$ satisfy the constraint, and let $W^* = \Sigma A'(A\Sigma A')^{-1}B$. From the chain of equalities, we derive:

$$(W_1 - W^*)'\Sigma^{-1}W^* = W_1'\Sigma^{-1}W^* - W^{*'}\Sigma^{-1}W^* = 0_{p \times p} = (0_{p \times p})' = W^{*'}\Sigma^{-1}(W_1 - W^*)$$

These are unusual quadratic forms. The usual quadratic form is $y'\Sigma x$, where the factors before and after $\Sigma^{-1}$ are $1 \times n$ and $n \times 1$ vectors. Here the form is $Y'\Sigma^{-1}X$, where the factors before and after $\Sigma^{-1}$ are $p \times n$ and $n \times p$ matrices, and the integer $p$ may exceed one.

But for now, consider the usual quadratic form in the special case that $\Sigma = I_n$. Actuaries know that $x'I_n x = x'x = \sum_{i=1}^{n} x_i^2$ is the square of the distance from the origin of $\Re^n$ to x (or the area of a square the length of whose sides is that distance). Less well known is that $y'x = x'y = \sum_{i=1}^{n} x_i y_i$ represents the area of a rectangle, the length of one of whose sides is

the length of the projection of one vector onto the other. Most will recognize, however, an equivalent interpretation, viz., that $y'x = x'y = 0$ if and only if $x \perp y$. The standard (Euclidean) definition of the distance from $x$ to $y$ is $d(x,y) = \sqrt{(y-x)'(y-x)}$. It has the three properties of a metric on $\mathfrak{R}^n$:

1. $d(x,y) \geq 0; \; d(x,y) = 0 \Leftrightarrow x = y \quad$ non-negativity; trivially zero
2. $d(y,x) = d(x,y) \qquad\qquad\qquad$ symmetry
3. $d(x,y) + d(y,z) \geq d(x,z) \qquad\quad$ triangle inequality

But for any positive-definite matrix $\Sigma_{n \times n}$, one can define a valid "$\Sigma$ metric" on $\mathfrak{R}^n$ as

$d_\Sigma(x,y) = \sqrt{(y-x)'\Sigma^{-1}(y-x)}$, which is valid in that it possesses these three properties.[2]

The matrix $\Sigma$ represents a combination of scaling and rotating the axes of $\mathfrak{R}^n$.

So what is special in the Gauss-Markov theorem about $W^* = \Sigma A'(A\Sigma A')^{-1}B$? Adapting the concept of perpendicularity to a metric, we have:

$$\left(W_1 - W^*\right)'\Sigma^{-1}W^* = 0_{p \times p}$$

---

[2] Some confusion results from using the inverse of $\Sigma$ in the quadratic form; one must think twice to determine whether something is a $\Sigma$ metric or a $\Sigma^{-1}$ metric. However, consider the usual formula for the ellipse whose major semi-axis is two units and minor semi-axis is one: $(x_1/2)^2 + (x_2/1)^2 = 1^2$. As a quadratic form this would be:

$$[x_1 \quad x_2]\begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix}^{-1}\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = 1^2$$

It seemed more natural to call this a [2 0; 0 1] metric (as if to say, "Two units on the first axis count as one unit on the second."), rather than to call it a [½ 0; 0 1] metric. This ellipse is the set of points in $\mathfrak{R}^2$ whose distance from the origin is one unit according to the 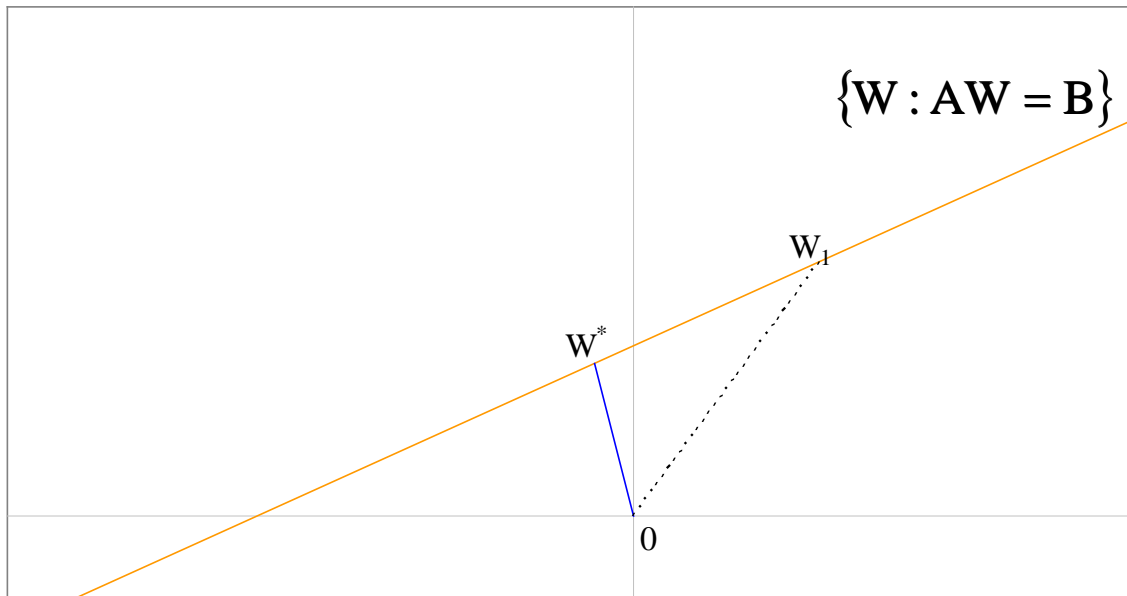[2 0; 0 1] metric. It may help some readers to know that $d_\Sigma(x,y) = \sqrt{(y-x)'\Sigma^{-1}(y-x)}$ is called the "Mahalanobis distance" (cf. Wikipedia), in whose definition the $\Sigma$ matrix is inverted. Appendix A provides a proof of the triangle inequality, as well as a justification of the geometric interpretation of $x'y$ as the product of the length of x and the length of the projection of y onto x.

This means that according to the $\Sigma$ metric $\mathbf{W}^*$ is perpendicular to $\mathbf{W}_1 - \mathbf{W}^*$ (and vice versa). In mathematical notation, $\mathbf{W}^* \underset{\Sigma}{\perp} \left( \mathbf{W}_1 - \mathbf{W}^* \right)$. The heart of the Gauss-Markov theorem, expressed above as $\Phi(\mathbf{W}_1) - \Phi(\mathbf{W}^*) = \left( \mathbf{W}_1 - \mathbf{W}^* \right)' \Sigma^{-1} \left( \mathbf{W}_1 - \mathbf{W}^* \right)$, is really just the Pythagorean theorem adapted to the $\Sigma$ metric:

$$\mathbf{W}^{*'} \Sigma^{-1} \mathbf{W}^* + \left( \mathbf{W}_1 - \mathbf{W}^* \right)' \Sigma^{-1} \left( \mathbf{W}_1 - \mathbf{W}^* \right) = \mathbf{W}_1' \Sigma^{-1} \mathbf{W}_1$$

$\mathbf{W}^*$ is the element of the constraint set closest to the origin according to the $\Sigma$ metric. The following diagram clarifies this:



The orange line represents the constraint set.[3] The origin, $\mathbf{W}^*$, and $\mathbf{W}_1$ form the $\Sigma$-right triangle, of which $\mathbf{W}_1$ is the hypotenuse, and $\mathbf{W}^*$ and $\mathbf{W}_1 - \mathbf{W}^*$ are the legs. The salient point is that the $\Sigma$ area of the square with side $\mathbf{W}^*$ is less than or equal to that of the square with side $\mathbf{W}_1$, or $\mathbf{W}^{*'} \Sigma \mathbf{W}^* \leq \mathbf{W}_1' \Sigma \mathbf{W}_1$, and equal if and only if $\mathbf{W}_1 = \mathbf{W}^*$. This is valid even

---

[3] Since the constraint on $\mathbf{W}$ is linear, the constraint set is a hyperplane (technically, an affine apace). The Gauss-Markov theorem requires a *linear* constraint; constraints involving curvature are inadmissible.

when the area concept is abstracted from a non-negative scalar to a non-negative-definite matrix.

This ends the geometric interpretation. Gauss-Markov reasoning happens whenever a quadratic form is to be minimized subject to a linear constraint. Gauss-Markov/BLUE proofs are abstractions of what we all learned in plane Geometry, viz., that the shortest distance from a point to a straight line is along a line segment perpendicular to the line. Lines are abstracted into linear constraints and distance is abstracted into a $\Sigma$ metric.

It is hardly necessary to memorize the formula for $\mathbf{W}^*$. With the following heuristic reasoning one can derive it on the fly. Since $\mathbf{A}_{m \times n}$ is of full-row rank (or of rank $m$), the $m \times m$ matrix $\mathbf{A}\mathbf{A}'$ is invertible. In fact, as stated above, for any positive-definite $\mathbf{T}_{n \times n}$, $\mathbf{A}\mathbf{T}\mathbf{A}'$ is invertible. Thus, there is a family of "right inverses" of $\mathbf{A}$ that have the form $\mathbf{T}\mathbf{A}'(\mathbf{A}\mathbf{T}\mathbf{A}')^{-1}$. $\mathbf{W}^*$ will be the matrix product of one of these right inverses and $\mathbf{B}$, i.e., $\mathbf{W}^* = \mathbf{T}\mathbf{A}'(\mathbf{A}\mathbf{T}\mathbf{A}')^{-1}\mathbf{B}$. Since we seek to minimize $\mathbf{W}'\Sigma^{-1}\mathbf{W}$, distance is measured according to a $\mathbf{T} = \Sigma$ metric. According to this metric $\mathbf{W}^* = \Sigma\mathbf{A}'(\mathbf{A}\Sigma\mathbf{A}')^{-1}\mathbf{B}$ is the element of the constraint set closest to the origin.

## 4. PROJECTION INTO THE CONSTRAINT SPACE

In the interest of economy and precision, let us introduce some more formalism. Our 'W' variables denote elements of $\mathfrak{R}^{n \times p}$, the real space of $n \times p$ dimensions. Let us use '$\Omega$' to denote the constraint set: $\Omega = \{\mathbf{W} \in \mathfrak{R}^{n \times p} : \mathbf{A}\mathbf{W} = \mathbf{B}\}$. Obviously, $\Omega \subseteq \mathfrak{R}^{n \times p}$; but it is not empty under the assumption that $\mathbf{A}_{m \times n}$ is of full-row rank. In fact, we have just seen that

$W^* = \Sigma A'(A\Sigma A')^{-1}B$ is the element of $\Omega$ closest to the origin of $\Re^{n \times p}$ according to the $\Sigma$ metric. We may say that $W^* = \Sigma A'(A\Sigma A')^{-1}B$ is the $\Sigma$ projection of the origin into $\Omega$. In general, what is the $\Sigma$ projection of *any* element of $\Re^{n \times p}$ into $\Omega$?

Using '$P$' for projection, we define $P(W_0; \Omega, \Sigma)$ as the function which projects $W_0 \in \Re^{n \times p}$ into $\Omega$ according to the $\Sigma$ metric. As before, $\Omega$ is the non-empty solution set of the linear constraint $AW = B$, and $\Sigma$ is positive-definite. When these parameters are understood, we will use the abbreviation $P(W_0)$. So $P(W_0; \Omega, \Sigma)$ is *an* element of $\Omega$ that minimizes the $\Sigma$-metric distance from $W_0$ to $\Omega$. Equivalently, it minimizes the quadratic form

$$\Phi(W) = (W - W_0)'\Sigma^{-1}(W - W_0) \text{ subject to } AW = B.$$

We could argue from scratch as in Section 2, but the following analysis is more insightful. The constraint $AW = B$ is equivalent to $A(W - W_0) = B - AW_0$. So the projection problem is to minimize $(W - W_0)'\Sigma^{-1}(W - W_0)$ subject to $A(W - W_0) = B - AW_0$. This is the Gauss-Markov problem with two changes in variables:

$$W \quad \rightarrow \quad W - W_0$$
$$B \quad \rightarrow \quad B - AW_0$$

Hence, the Gauss-Markov theorem states that $(W - W_0)^* = \Sigma A'(A\Sigma A')^{-1}(B - AW_0)$ uniquely minimizes $\Phi(W - W_0) = (W - W_0)'\Sigma^{-1}(W - W_0)$. But since $W_0$ is a constant, $(W - W_0)^* = W^* - W_0$, or $W^* = W_0 + (W - W_0)^* = W_0 + \Sigma A'(A\Sigma A')^{-1}(B - AW_0)$. So there is not just *an* element of projection, but a unique element:

$$P(W; \Omega, \Sigma) = W + \Sigma A'(A\Sigma A')^{-1}(B - AW)$$

$$= \left\{ I_n - \Sigma A'(A\Sigma A')^{-1}A \right\}W + \Sigma A'(A\Sigma A')^{-1}B$$

As a check:

$$AP(W; \Omega, \Sigma) = A\left\{ I_n - \Sigma A'(A\Sigma A')^{-1}A \right\}W + A\Sigma A'(A\Sigma A')^{-1}B$$

$$= \left\{ A - A\Sigma A'(A\Sigma A')^{-1}A \right\}W + A\Sigma A'(A\Sigma A')^{-1}B$$

$$= \left\{ A - I_m A \right\}W + I_m B$$

$$= B$$

Hence, for all $W \in \mathfrak{R}^{n \times p}$, $AP(W) \in \Omega$. So $P$ is a mapping from $\mathfrak{R}^{n \times p}$ into $\Omega$, i.e.,

$P : \mathfrak{R}^{n \times p} \to \Omega$. In particular, the mapping of the origin is:

$$P(0_{n \times p}) = \left\{ I_n - \Sigma A'(A\Sigma A')^{-1}A \right\}0 + \Sigma A'(A\Sigma A')^{-1}B = \Sigma A'(A\Sigma A')^{-1}B,$$

which is the $'W^*'$ of the theorem itself. Accordingly, we may employ the formulation

$P(W) = \left\{ I_n - \Sigma A'(A\Sigma A')^{-1}A \right\}W + P(0).$

$P$ maps element $W \in \mathfrak{R}^{n \times p}$ to the closest element of constraint set $\Omega$ according to the $\Sigma$ metric. Geometrically, $P$ sends a $\Sigma$ perpendicular from $W$ into $\Omega$; in symbols, $P(W) - W \underset{\Sigma}{\perp} W_1 - P(W)$, for every $W_1 \in \Omega$, as the following algebra shows:

$$\left(P(\mathbf{W}) - \mathbf{W}\right)' \Sigma^{-1}\left(\mathbf{W}_1 - P(\mathbf{W})\right)$$

$$= \left(\left\{\mathbf{I}_n - \Sigma\mathbf{A}'(\mathbf{A}\Sigma\mathbf{A}')^{-1}\mathbf{A}\right\}\mathbf{W} + P(0) - \mathbf{W}\right)' \Sigma^{-1}\left(\mathbf{W}_1 - P(\mathbf{W})\right)$$

$$= \left(P(0) - \Sigma\mathbf{A}'(\mathbf{A}\Sigma\mathbf{A}')^{-1}\mathbf{A}\mathbf{W}\right)' \Sigma^{-1}\left(\mathbf{W}_1 - P(\mathbf{W})\right)$$

$$= \left(\Sigma\mathbf{A}'(\mathbf{A}\Sigma\mathbf{A}')^{-1}\mathbf{B} - \Sigma\mathbf{A}'(\mathbf{A}\Sigma\mathbf{A}')^{-1}\mathbf{A}\mathbf{W}\right)' \Sigma^{-1}\left(\mathbf{W}_1 - P(\mathbf{W})\right)$$

$$= \left(\Sigma\mathbf{A}'(\mathbf{A}\Sigma\mathbf{A}')^{-1}(\mathbf{B} - \mathbf{A}\mathbf{W})\right)' \Sigma^{-1}\left(\mathbf{W}_1 - P(\mathbf{W})\right)$$

$$= (\mathbf{B} - \mathbf{A}\mathbf{W})'(\mathbf{A}\Sigma\mathbf{A}')^{-1}\mathbf{A}\Sigma\Sigma^{-1}\left(\mathbf{W}_1 - P(\mathbf{W})\right)$$

$$= (\mathbf{B} - \mathbf{A}\mathbf{W})'(\mathbf{A}\Sigma\mathbf{A}')^{-1}\mathbf{A}\left(\mathbf{W}_1 - P(\mathbf{W})\right)$$

$$= (\mathbf{B} - \mathbf{A}\mathbf{W})'(\mathbf{A}\Sigma\mathbf{A}')^{-1}\left(\mathbf{A}\mathbf{W}_1 - \mathbf{A}P(\mathbf{W})\right)$$

$$= (\mathbf{B} - \mathbf{A}\mathbf{W})'(\mathbf{A}\Sigma\mathbf{A}')^{-1}(\mathbf{B} - \mathbf{B})$$

$$= 0$$

Because of the first property of a metric (zero-triviality), within the restricted domain $\Omega$, $P$ is the identity mapping. Hence, not only is $P$ a mapping *into* the constraint set $\Omega$; it is also a mapping *onto* $\Omega$. Nonetheless, we will prove it algebraically. If $\mathbf{W} \in \Omega$:

$$P(\mathbf{W}) = \left\{\mathbf{I}_n - \Sigma\mathbf{A}'(\mathbf{A}\Sigma\mathbf{A}')^{-1}\mathbf{A}\right\}\mathbf{W} + \Sigma\mathbf{A}'(\mathbf{A}\Sigma\mathbf{A}')^{-1}\mathbf{B}$$

$$= \mathbf{W} - \Sigma\mathbf{A}'(\mathbf{A}\Sigma\mathbf{A}')^{-1}(\mathbf{A}\mathbf{W}) + \Sigma\mathbf{A}'(\mathbf{A}\Sigma\mathbf{A}')^{-1}\mathbf{B}$$

$$= \mathbf{W} - \Sigma\mathbf{A}'(\mathbf{A}\Sigma\mathbf{A}')^{-1}(\mathbf{B}) + \Sigma\mathbf{A}'(\mathbf{A}\Sigma\mathbf{A}')^{-1}\mathbf{B}$$

$$= \mathbf{W}$$

Conversely, if $P(\mathbf{W}) = \mathbf{W}$, then $\mathbf{A}\mathbf{W} = \mathbf{A}P(\mathbf{W}) = \mathbf{B}$ and $\mathbf{W} \in \Omega$. Therefore, $P$ is a many-to-one mapping from $\Re^{n \times p}$ onto constraint set $\Omega$, and an element of $\Re^{n \times p}$ belongs to $\Omega$ if and only if $P$ acts upon it as an identity mapping.

Just as the $\Sigma$-metric right inverse $\Sigma A'(A\Sigma A')^{-1}$ is conspicuous in the formula $P(0_{n \times p}) = \Sigma A'(A\Sigma A')^{-1} B$, so too is it conspicuous in the formula for what we will call the "$\Sigma$-projection matrix" $I_n - \Sigma A'(A\Sigma A')^{-1} A$.[4] Since $\Sigma$ is positive-definite, it can be Cholesky-decomposed as $\Sigma = QQ'$ for some non-singular $Q_{n \times n}$. Then the matrix can be factored as $I_n - \Sigma A'(A\Sigma A')^{-1} A = Q\{I_n - Q'A'(AQQ'A')^{-1} AQ\}Q^{-1} = QMQ^{-1}$. So the rank of the matrix is the rank of $M = I_n - Q'A'(AQQ'A')^{-1} AQ$. But $M$ is a (symmetric) idempotent matrix (i.e., $M = MM'$), and the rank of an idempotent matrix equals its trace (Judge [1988, Appendix A.4 and A.12] and Halliwell [1997, Appendix B, 317; also Note 3]). Employing basic theorems about the trace operator, we derive:

$$
\begin{aligned}
rank\left(I_n - \Sigma A'(A\Sigma A')^{-1} A\right) &= rank\left(Q\{I_n - Q'A'(AQQ'A')^{-1} AQ\}Q^{-1}\right) \\
&= rank\left(I_n - Q'A'(AQQ'A')^{-1} AQ\right) \\
&= Tr\left(I_n - Q'A'(AQQ'A')^{-1} AQ\right) \\
&= Tr(I_n) - Tr\left(Q'A'(AQQ'A')^{-1} AQ\right) \\
&= Tr(I_n) - Tr\left((AQQ'A')^{-1} AQQ'A'\right) \\
&= Tr(I_n) - Tr(I_m) \\
&= n - m
\end{aligned}
$$

So every column of the elements of constraint set $\Omega$ has $m$ fewer degrees of freedom than the columns of $\Re^{n \times p}$; in a sense, the dimensionality of $\Omega$ is $\Re^{(n-m) \times p}$. This can be surmised from the full-row rank of $A_{m \times n}$, which imposes $m$ independent restrictions on the elements of $\Re^{n \times p}$ that belong to $\Xi$. This will prove useful in Section 9, in which we will treat linear statistical models with parameter constraints.

---

[4] The projection matrix shows to its greatest effect in the homogeneous form, i.e., in the differential form, $P(W_2) - P(W_1) = \{I_n - \Sigma A'(A\Sigma A')^{-1} A\}(W_2 - W_1)$.

## 5. VARIANCE AS A METRIC

In this section we will show how the variance of a random vector serves as its natural metric.

Let $\mathbf{x}$ be an $n \times 1$ random vector with mean $\mu$ and non-degenerate variance $\Sigma$. Since the variance is non-degenerate, the variance of every non-zero linear combination of its elements is positive, i.e., $\Sigma$ is positive-definite and $\Sigma^{-1}$ exists. The variance of a random vector is a measure of its ability to differ from its mean. So the distances of random vectors from their means should somehow be invariant, when their variances serve as their distance metrics.

The square of the $\Sigma$-metric distance of $\mathbf{x}$ from its mean is $d_\Sigma^2(\mu, \mathbf{x}) = (\mathbf{x} - \mu)' \Sigma^{-1} (\mathbf{x} - \mu)$. And by definition, the variance of $\mathbf{x}$ is $Var[\mathbf{x}] = \Sigma = E\left[ (\mathbf{x} - \mu)(\mathbf{x} - \mu)' \right]$. Using again the trace-operator theorems of the previous section, we find:

$$
\begin{aligned}
E[d_\Sigma^2(\mu, \mathbf{x})] &= E\left[ (\mathbf{x} - \mu)' \Sigma^{-1} (\mathbf{x} - \mu) \right] \\
&= Tr\left( E\left[ (\mathbf{x} - \mu)' \Sigma^{-1} (\mathbf{x} - \mu) \right] \right) \\
&= E\left[ Tr\left( (\mathbf{x} - \mu)' \Sigma^{-1} (\mathbf{x} - \mu) \right) \right] \\
&= E\left[ Tr\left( \Sigma^{-1} (\mathbf{x} - \mu)(\mathbf{x} - \mu)' \right) \right] \\
&= Tr\left( E\left[ \Sigma^{-1} (\mathbf{x} - \mu)(\mathbf{x} - \mu)' \right] \right) \\
&= Tr\left( \Sigma^{-1} E\left[ (\mathbf{x} - \mu)(\mathbf{x} - \mu)' \right] \right) \\
&= Tr\left( \Sigma^{-1} \Sigma \right) \\
&= Tr(I_n) = n
\end{aligned}
$$

Hence, a random vector's variance is its natural metric, according to which its expected squared distance from its mean equals its dimensionality *n*, the degrees of its freedom.[5]

## 6. PROJECTIONS OF RANDOM VECTORS

As in the previous section, let **x** be an $n \times 1$ random vector with mean $\mu$ and variance $\Sigma$, i.e., $\mathbf{x} \sim (\mu, \Sigma)$. Since the Gauss-Markov theorem has to do with abstract projective geometry, we may inquire about the moments of the $\Sigma$ projection of **x** into the constraint space $\Omega = \{ \mathbf{x} \in \Re^n : A_{m \times n} \mathbf{x}_{n \times 1} = \mathbf{b}_{m \times 1} \}$.

The $\Sigma$ projection is $P(\mathbf{x}; \Omega, \Sigma) = \{ I_n - \Sigma A'(A\Sigma A')^{-1} A \} \mathbf{x} + \Sigma A'(A\Sigma A')^{-1} \mathbf{b}$. Therefore:

$$
\begin{aligned}
E[P(\mathbf{x}; \Omega, \Sigma)] &= E\left[ \{ I_n - \Sigma A'(A\Sigma A')^{-1} A \} \mathbf{x} + \Sigma A'(A\Sigma A')^{-1} \mathbf{b} \right] \\
&= \{ I_n - \Sigma A'(A\Sigma A')^{-1} A \} E[\mathbf{x}] + \Sigma A'(A\Sigma A')^{-1} \mathbf{b} \\
&= \{ I_n - \Sigma A'(A\Sigma A')^{-1} A \} \mu + \Sigma A'(A\Sigma A')^{-1} \mathbf{b} \\
&= \mu + \Sigma A'(A\Sigma A')^{-1} (\mathbf{b} - \mu)
\end{aligned}
$$

The variance follows from the standard formula $Var[\mathbf{Qx}] = Q\,Var[\mathbf{x}]Q'$:

---

[5] Moreover, if **x** is multivariate normal, or if $\mathbf{x} \sim N(\mu, \Sigma)$, then $(\mathbf{x} - \mu)' \Sigma^{-1} (\mathbf{x} - \mu) \sim \chi_n^2$ (Judge [1988], §2.5.9). The multivariate normal distribution is unique in that its probability distribution is a function of its variance metric: $f_{\mathbf{x}}(\mathbf{x}) \propto e^{-\frac{1}{2}(\mathbf{x}-\mu)'\Sigma^{-1}(\mathbf{x}-\mu)}$. Most likely, this is the ultimate reason why normality is preserved under any linear transformation.

$$Var[P(\mathbf{x}; \Omega, \Sigma)] = Var\left[\left\{I_n - \Sigma A'(A\Sigma A')^{-1}A\right\}\mathbf{x} + \Sigma A'(A\Sigma A')^{-1}b\right]$$

$$= Var\left[\left\{I_n - \Sigma A'(A\Sigma A')^{-1}A\right\}\mathbf{x}\right]$$

$$= \left\{I_n - \Sigma A'(A\Sigma A')^{-1}A\right\}Var[\mathbf{x}]\left\{I_n - \Sigma A'(A\Sigma A')^{-1}A\right\}'$$

$$= \left\{I_n - \Sigma A'(A\Sigma A')^{-1}A\right\}\Sigma\left\{I_n - A'(A\Sigma A')^{-1}A\Sigma\right\}$$

$$= \Sigma - \Sigma A'(A\Sigma A')^{-1}A\Sigma - \Sigma A'(A\Sigma A')^{-1}A\Sigma + \Sigma A'(A\Sigma A')^{-1}A\Sigma A'(A\Sigma A')^{-1}A\Sigma$$

$$= \Sigma - \Sigma A'(A\Sigma A')^{-1}A\Sigma - \Sigma A'(A\Sigma A')^{-1}A\Sigma + \Sigma A'(A\Sigma A')^{-1}A\Sigma$$

$$= \Sigma - \Sigma A'(A\Sigma A')^{-1}A\Sigma$$

Therefore, if $\mathbf{x} \sim (\mu, \Sigma)$, then $P(\mathbf{x}; \Omega, \Sigma) \sim \left(P(\mu; \Omega, \Sigma), \Sigma - \Sigma A'(A\Sigma A')^{-1}A\Sigma\right)$. As a check,

$E[A\mathbf{x}] = AP(\mu) = b$ and $Var[A\mathbf{x}] = A\left(\Sigma - \Sigma A'(A\Sigma A')^{-1}A\Sigma\right)A' = 0_{m\times m}$. This will prove

useful in Section 11.

## 7. PARAMETER ESTIMATION IN THE LINEAR STATISTICAL MODEL

Now let us apply our Gauss-Markov theorem to the linear statistical model. First, and as an easy start, we will apply it to derive the best linear unbiased estimator (BLUE) of the parameter $\beta$ in the model $\mathbf{y} = X_{t\times k}\beta_{k\times 1} + \mathbf{e}$, where $Var[\mathbf{e}] = \Sigma_{t\times t}$. $X$ is of full-column rank, and $\Sigma$ is positive-definite. The estimator is linear in $\mathbf{y}$, or $\hat{\beta} = W'\mathbf{y}$.[6] Because it is unbiased for all $\beta$, $E[\hat{\beta}] = E[W'\mathbf{y}] = W'X\beta = \beta$. So matrix $W'$ is constrained according to the equation $W'X = I_k$, which transposes as $X'W = I_k$. The best of the unbiased estimators minimizes $Var[\hat{\beta}] = Var[W'\mathbf{y}] = W'\Sigma W$. So the problem is to minimize

---

[6] For a reason immediately to become apparent, we use here the transpose of W.

$\Phi(W) = W'\Sigma W = W'(\Sigma^{-1})^{-1}W$ subject to $X'W = I_k$. The correspondences between the theorem and this model are:

| Theorem | $\leftarrow$ | Model |
|:---:|:---:|:---:|
| $\Sigma$ | $\leftarrow$ | $\Sigma^{-1}$ |
| $W$ | $\leftarrow$ | $W$ |
| $A$ | $\leftarrow$ | $X'$ |
| $B$ | $\leftarrow$ | $I_k$ |

$X'$ is of full-row rank, because $X$ is of full-column rank. Hence, according to the

theorem, $W^* = \Sigma^{-1}X''(X'\Sigma^{-1}X'')^{-1}I_k = \Sigma^{-1}X(X'\Sigma^{-1}X)^{-1}$. So $\hat{\boldsymbol{\beta}} = W^{*'}\mathbf{y} = (X'\Sigma^{-1}X)^{-1}X'\Sigma^{-1}\mathbf{y}$

and $Var[\hat{\boldsymbol{\beta}}] = W^{*'}\Sigma W^* = I_k'(X'\Sigma^{-1}X'')^{-1}I_k = (X'\Sigma^{-1}X)^{-1}$. Accordingly, $\hat{\boldsymbol{\beta}} = Var[\hat{\boldsymbol{\beta}}]X'\Sigma^{-1}\mathbf{y}$.


## 8. PREDICTION IN THE LINEAR STATISTICAL MODEL

The goal of most linear modeling is not to estimate the parameter $\beta$, but rather to estimate[7]

quantities which eventually will be observed. Although the model makes such quantities

dependent on the parameter, the parameter itself is usually hypothetical and never to be

observed. With partitioning between the observed $\mathbf{y}_1$ and the to-be-predicted $\mathbf{y}_2$ (hence,

containing missing values) the general form of the linear statistical model is:

$$\begin{bmatrix} \mathbf{y}_1 \\ \hline \mathbf{y}_2 \end{bmatrix} = \begin{bmatrix} X_{1\ (t_1 \times k)} \\ \hline X_{2\ (t_2 \times k)} \end{bmatrix} \beta_{k \times 1} + \begin{bmatrix} \mathbf{e}_1 \\ \hline \mathbf{e}_2 \end{bmatrix}, \text{ where } Var\begin{bmatrix} \mathbf{e}_1 \\ \hline \mathbf{e}_2 \end{bmatrix} = \begin{bmatrix} \Sigma_{11} & \vdots & \Sigma_{12} \\ \hline \Sigma_{21} & \vdots & \Sigma_{22} \end{bmatrix} = \Sigma$$

Not only must $\Sigma$ be symmetric and non-negative-definite, $\Sigma_{11}$ must be positive-definite,

and $X_1$ must be of full-column rank. We seek the best linear-in-$\mathbf{y}_1$, unbiased estimator

---

[7] More accurately, the goal is to predict – we seek the best linear unbiased prediction. But we will continue to call this BLUE, because BLUP already has a different technical meaning in statistics (*Wikipedia*, "Best linear unbiased prediction").

(BLUE) of $\mathbf{y}_2$, i.e., $\hat{\mathbf{y}}_2 = \mathbf{W}'\mathbf{y}_1$ for some matrix $\mathbf{W}'$, which depends only on the partitions of the design $\mathbf{X}$ and variance $\Sigma$ matrices. Because the estimator is unbiased for all $\beta$,

$$0 = E[\mathbf{y}_2 - \hat{\mathbf{y}}_2] = E[\mathbf{y}_2 - \mathbf{W}'\mathbf{y}_1] = (\mathbf{X}_2 - \mathbf{W}'\mathbf{X}_1)\beta.$$ Thus the estimator is unbiased if and only if $\mathbf{W}'\mathbf{X}_1 = \mathbf{X}_2$. By transposition, $\mathbf{X}_1'\mathbf{W} = \mathbf{X}_2'$, where W is $t_1 \times t_2$.

But now there is a complication in being "best." Predicting $\mathbf{y}_2$ as $\hat{\mathbf{y}}_2$, we will err by the amount $\mathbf{y}_2 - \hat{\mathbf{y}}_2$. So it is the prediction-error variance that we must minimize:

$$Var[\mathbf{y}_2 - \hat{\mathbf{y}}_2] = Var[\mathbf{y}_2 - \mathbf{W}'\mathbf{y}_1]$$

$$= Var\left[\begin{bmatrix} -\mathbf{W}' & I_{t_2} \end{bmatrix}\begin{bmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \end{bmatrix}\right]$$

$$= \begin{bmatrix} -\mathbf{W}' & I_{t_2} \end{bmatrix} Var\begin{bmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \end{bmatrix}\begin{bmatrix} -\mathbf{W}' & I_{t_2} \end{bmatrix}'$$

$$= \begin{bmatrix} -\mathbf{W}' & I_{t_2} \end{bmatrix}\begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}\begin{bmatrix} -\mathbf{W} \\ I_{t_2} \end{bmatrix}$$

$$= \mathbf{W}'\Sigma_{11}\mathbf{W} - \mathbf{W}'\Sigma_{12} - \Sigma_{21}\mathbf{W} + \Sigma_{22}$$

Although we can ignore the constant $\Sigma_{22}$ in the minimization, we cannot ignore the second and third terms, which are linear in W.

The key here is to apply the one-to-one transform $V \leftrightarrow W - \Sigma_{11}^{-1}\Sigma_{12}$. The transformation of the constraint set is $\{W : X_1'W = X_2'\} = \{V : X_1'V = X_1'W - X_1'\Sigma_{11}^{-1}\Sigma_{12} = X_2' - X_1'\Sigma_{11}^{-1}\Sigma_{12}\}$. So expressed in terms of V:

$$Var[\mathbf{y}_2 - \hat{\mathbf{y}}_2] = \mathbf{W}'\Sigma_{11}\mathbf{W} - \mathbf{W}'\Sigma_{12} - \Sigma_{21}\mathbf{W} + \Sigma_{22}$$

$$= \left(\mathbf{V} + \Sigma_{11}^{-1}\Sigma_{12}\right)'\Sigma_{11}\left(\mathbf{V} + \Sigma_{11}^{-1}\Sigma_{12}\right) - \left(\mathbf{V} + \Sigma_{11}^{-1}\Sigma_{12}\right)'\Sigma_{12} - \Sigma_{21}\left(\mathbf{V} + \Sigma_{11}^{-1}\Sigma_{12}\right) + \Sigma_{22}$$

$$= \left(\mathbf{V}' + \Sigma_{21}\Sigma_{11}^{-1}\right)\Sigma_{11}\left(\mathbf{V} + \Sigma_{11}^{-1}\Sigma_{12}\right) - \left(\mathbf{V}' + \Sigma_{21}\Sigma_{11}^{-1}\right)\Sigma_{12} - \Sigma_{21}\left(\mathbf{V} + \Sigma_{11}^{-1}\Sigma_{12}\right) + \Sigma_{22}$$

$$= \mathbf{V}'\Sigma_{11}\mathbf{V} + \mathbf{V}'\Sigma_{12} + \Sigma_{21}\mathbf{V} + \Sigma_{21}\Sigma_{11}^{-1}\Sigma_{12}$$

$$\quad - \mathbf{V}'\Sigma_{12} - \Sigma_{21}\Sigma_{11}^{-1}\Sigma_{12} - \Sigma_{21}\mathbf{V} - \Sigma_{21}\Sigma_{11}^{-1}\Sigma_{12} + \Sigma_{22}$$

$$= \mathbf{V}'\Sigma_{11}\mathbf{V} + \left(\Sigma_{22} - \Sigma_{21}\Sigma_{11}^{-1}\Sigma_{12}\right)$$

This transformation is a matrix version of completing the square. We can now apply the

Gauss-Markov theorem to the problem of minimizing $\mathbf{V}'\Sigma_{11}\mathbf{V} = \mathbf{V}'\left(\Sigma_{11}^{-1}\right)^{-1}\mathbf{V}$ subject to

$\mathbf{X}_1'\mathbf{V} = \mathbf{X}_2' - \mathbf{X}_1'\Sigma_{11}^{-1}\Sigma_{12}$. The correspondences are:

| Theorem | ← | Model |
|:---:|:---:|:---:|
| $\Sigma$ | ← | $\Sigma_{11}^{-1}$ |
| $\mathbf{W}$ | ← | $\mathbf{V}$ |
| $\mathbf{A}$ | ← | $\mathbf{X}_1'$ |
| $\mathbf{B}$ | ← | $\mathbf{X}_2' - \mathbf{X}_1'\Sigma_{11}^{-1}\Sigma_{12}$ |

As before, the conditions are met; $\mathbf{X}_1'$ is of full-row rank since $\mathbf{X}_1$ is of full-column rank.

Hence, $\mathbf{V}^* = \Sigma_{11}^{-1}\mathbf{X}_1\left(\mathbf{X}_1'\Sigma_{11}^{-1}\mathbf{X}_1\right)^{-1}\left(\mathbf{X}_2' - \mathbf{X}_1'\Sigma_{11}^{-1}\Sigma_{12}\right)$, and $\hat{\mathbf{y}}_2 = \mathbf{W}^*\mathbf{y}_1$, where:

$$\mathbf{W}^{*'} = \left(\mathbf{V}^* + \Sigma_{11}^{-1}\Sigma_{12}\right)'$$

$$= \mathbf{V}^{*'} + \Sigma_{21}\Sigma_{11}^{-1}$$

$$= \left(\mathbf{X}_2 - \Sigma_{21}\Sigma_{11}^{-1}\mathbf{X}_1\right)\left(\mathbf{X}_1'\Sigma_{11}^{-1}\mathbf{X}_1\right)^{-1}\mathbf{X}_1'\Sigma_{11}^{-1} + \Sigma_{21}\Sigma_{11}^{-1}$$

The minimized prediction-error variance is:

$$Var[\mathbf{y}_2 - \hat{\mathbf{y}}_2] = \mathbf{V}^{*'}\Sigma_{11}\mathbf{V}^* + \left(\Sigma_{22} - \Sigma_{21}\Sigma_{11}^{-1}\Sigma_{12}\right)$$

$$= \left(\mathbf{X}_2 - \Sigma_{21}\Sigma_{11}^{-1}\mathbf{X}_1\right)\left(\mathbf{X}_1'\Sigma_{11}^{-1}\mathbf{X}_1\right)^{-1}\left(\mathbf{X}_2 - \Sigma_{21}\Sigma_{11}^{-1}\mathbf{X}_1\right)' + \left(\Sigma_{22} - \Sigma_{21}\Sigma_{11}^{-1}\Sigma_{12}\right)$$

Introducing the estimator $\hat{\boldsymbol{\beta}} = \left(X_1'\Sigma_{11}^{-1}X_1\right)^{-1}X_1'\Sigma_{11}^{-1}\mathbf{y}_1$ and its variance $Var\left[\hat{\boldsymbol{\beta}}\right] = \left(X_1'\Sigma_{11}^{-1}X_1\right)^{-1}$

allows us to simplify:

$$\hat{\mathbf{y}}_2 = W^{*'}\mathbf{y}_1$$

$$= \left(X_2 - \Sigma_{21}\Sigma_{11}^{-1}X_1\right)\left(X_1'\Sigma_{11}^{-1}X_1\right)^{-1}X_1'\Sigma_{11}^{-1}\mathbf{y}_1 + \Sigma_{21}\Sigma_{11}^{-1}\mathbf{y}_1$$

$$= \left(X_2 - \Sigma_{21}\Sigma_{11}^{-1}X_1\right)\hat{\boldsymbol{\beta}} + \Sigma_{21}\Sigma_{11}^{-1}\mathbf{y}_1$$

$$= X_2\hat{\boldsymbol{\beta}} + \Sigma_{21}\Sigma_{11}^{-1}\left(\mathbf{y}_1 - X_1\hat{\boldsymbol{\beta}}\right)$$

$$Var\left[\mathbf{y}_2 - \hat{\mathbf{y}}_2\right] = \left(X_2 - \Sigma_{21}\Sigma_{11}^{-1}X_1\right)Var\left[\hat{\boldsymbol{\beta}}\right]\left(X_2 - \Sigma_{21}\Sigma_{11}^{-1}X_1\right)' + \left(\Sigma_{22} - \Sigma_{21}\Sigma_{11}^{-1}\Sigma_{12}\right)$$

Combining this completing-the-square technique with the Gauss-Markov theorem makes for

a cleaner and more elegant proof than that in Halliwell [1997, Appendix C, 328-330].

## 9. LINEAR STATISTICAL MODELS WITH PARAMETER CONSTRAINTS

Here we will impose upon the model of Section 8 a constraint on β, viz., $R_{j \times k}\beta_{k \times 1} = r_{j \times 1}$.

The rows of $R$ must be linearly independent, i.e., $R$ must be of full-row rank. The

constraint set $\left\{\beta \in \Re^k : R\beta = r\right\}$ is non-empty because right inverses of R exist, most

obviously $R'\left(RR'\right)^{-1}$. Hence, $\beta = R'\left(RR'\right)^{-1}r$ exists and satisfies the constraint.

Two procedures are commonly employed to solve β-constrained linear models. The first is

to reduce the parameter dimension according to the equation $\beta = \beta_0 + S_{k \times (k-j)}\gamma_{(k-j) \times 1}$, for

some matrix S (of full-column rank) such that $RS = 0_{j \times (k-j)}$, as done in Halliwell [1997,

Appendix B, 321-324]. This is the purist approach to the problem, but it requires an

understanding of eigen-decomposition, cannot be performed in Excel without add-ins, and

may suffer from the numerical-analysis problem of deciding when small eigenvalues should

be zeroed. The second procedure is to employ the Lagrange multiplier (Judge [1988, §6.2, 235-237]) to minimize $\Lambda(\beta, \lambda_{j\times 1}) = (\mathbf{y} - X\beta)' \Sigma^{-1}(\mathbf{y} - X\beta) + 2\lambda'(R\beta - \mathbf{r})$. But a third procedure (Halliwell [1998, Appendix C]) to us is the most convincing.

This procedure is to treat the $\beta$ constraint as the limit of $\mathbf{r} = R\beta + \boldsymbol{\eta}$ as $Var[\boldsymbol{\eta}] \to 0_{j\times j}$. We could have specified the variance as $\sigma^2 I_j$, and the limit as $\sigma^2 \to 0$; but for the sake of generality we will let $Var[\boldsymbol{\eta}] = H$ be any positive-definite matrix. So we can form the following augmented linear model, which satisfies the conditions of Section 8:

$$
\begin{bmatrix} \mathbf{y}_1 \\ \mathbf{r} \\ \hline \mathbf{y}_2 \end{bmatrix} = \begin{bmatrix} X_1 \\ R \\ \hline X_2 \end{bmatrix} \beta + \begin{bmatrix} \mathbf{e}_1 \\ \boldsymbol{\eta} \\ \hline \mathbf{e}_2 \end{bmatrix}, \text{ where } Var \begin{bmatrix} \mathbf{e}_1 \\ \boldsymbol{\eta} \\ \hline \mathbf{e}_2 \end{bmatrix} = \begin{bmatrix} \Sigma_{11} & 0 & \vdots & \Sigma_{12} \\ 0 & H & \vdots & 0 \\ \hline \Sigma_{21} & 0 & \vdots & \Sigma_{22} \end{bmatrix}
$$

The parameter estimator, which depends on H, is:

$$
\begin{aligned}
\hat{\boldsymbol{\beta}}(H) &= \left( \begin{bmatrix} X_1' & R' \end{bmatrix} \begin{bmatrix} \Sigma_{11} & 0 \\ 0 & H \end{bmatrix}^{-1} \begin{bmatrix} X_1 \\ R \end{bmatrix} \right)^{-1} \begin{bmatrix} X_1' & R' \end{bmatrix} \begin{bmatrix} \Sigma_{11} & 0 \\ 0 & H \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{y}_1 \\ \mathbf{r} \end{bmatrix} \\
&= \left( \begin{bmatrix} X_1' & R' \end{bmatrix} \begin{bmatrix} \Sigma_{11}^{-1} & 0 \\ 0 & H^{-1} \end{bmatrix} \begin{bmatrix} X_1 \\ R \end{bmatrix} \right)^{-1} \begin{bmatrix} X_1' & R' \end{bmatrix} \begin{bmatrix} \Sigma_{11}^{-1} & 0 \\ 0 & H^{-1} \end{bmatrix} \begin{bmatrix} \mathbf{y}_1 \\ \mathbf{r} \end{bmatrix} \\
&= \left( X_1' \Sigma_{11}^{-1} X_1 + R' H^{-1} R \right)^{-1} \left( X_1' \Sigma_{11}^{-1} \mathbf{y}_1 + R' H^{-1} \mathbf{r} \right) \\
&= Var[\hat{\boldsymbol{\beta}}(H)] \left( X_1' \Sigma_{11}^{-1} \mathbf{y}_1 + R' H^{-1} \mathbf{r} \right)
\end{aligned}
$$

Therefore, according to the formulas of the previous section, the predictor is:

$$
\begin{aligned}
\hat{\mathbf{y}}_2(H) &= X_2 \hat{\boldsymbol{\beta}}(H) + \begin{bmatrix} \Sigma_{21} & 0 \end{bmatrix} \begin{bmatrix} \Sigma_{11} & 0 \\ 0 & H \end{bmatrix}^{-1} \left( \begin{bmatrix} \mathbf{y}_1 \\ \mathbf{r} \end{bmatrix} - \begin{bmatrix} X_1 \\ R \end{bmatrix} \hat{\boldsymbol{\beta}}(H) \right) \\
&= X_2 \hat{\boldsymbol{\beta}}(H) + \Sigma_{21} \Sigma_{11}^{-1} \left( \mathbf{y}_1 - X_1 \hat{\boldsymbol{\beta}}(H) \right)
\end{aligned}
$$

And the variance of the prediction error is:

$$Var[\mathbf{y}_2 - \hat{\mathbf{y}}_2(\mathrm{H})] = \left( X_2 - \begin{bmatrix} \Sigma_{21} & 0 \end{bmatrix} \begin{bmatrix} \Sigma_{11} & 0 \\ 0 & \mathrm{H} \end{bmatrix}^{-1} \begin{bmatrix} X_1 \\ R \end{bmatrix} \right) Var[\hat{\boldsymbol{\beta}}(\mathrm{H})] \left( X_2 - \begin{bmatrix} \Sigma_{21} & 0 \end{bmatrix} \begin{bmatrix} \Sigma_{11} & 0 \\ 0 & \mathrm{H} \end{bmatrix}^{-1} \begin{bmatrix} X_1 \\ R \end{bmatrix} \right)'$$

$$+ \left( \Sigma_{22} - \begin{bmatrix} \Sigma_{21} & 0 \end{bmatrix} \begin{bmatrix} \Sigma_{11} & 0 \\ 0 & \mathrm{H} \end{bmatrix}^{-1} \begin{bmatrix} \Sigma_{12} \\ 0 \end{bmatrix} \right)$$

$$= \left( X_2 - \Sigma_{21}\Sigma_{11}^{-1}X_1 \right) Var[\hat{\boldsymbol{\beta}}(\mathrm{H})] \left( X_2 - \Sigma_{21}\Sigma_{11}^{-1}X_1 \right)' + \left( \Sigma_{22} - \Sigma_{21}\Sigma_{11}^{-1}\Sigma_{12} \right)$$

These two formulas depend on H only insofar as $\hat{\boldsymbol{\beta}}$ depends on H. Therefore, it remains

for us to determine $\boldsymbol{\beta}^* = \lim\limits_{\mathrm{H}\to 0} \hat{\boldsymbol{\beta}}(\mathrm{H})$.

We start with $Var[\boldsymbol{\beta}^*] = \lim\limits_{\mathrm{H}\to 0} \left( X_1'\Sigma_{11}^{-1}X_1 + R'\mathrm{H}^{-1}R \right)^{-1}$. The following proof makes use of the

theorem $(A + BDC)^{-1} = A^{-1} - A^{-1}B(D^{-1} + CA^{-1}B)^{-1}CA^{-1}$ (cf. Judge [1988, A.7, 938]; the

inverses must exist, as they do here): [8]

$$Var[\boldsymbol{\beta}^*] = \lim\limits_{\mathrm{H}\to 0} \left( X_1'\Sigma_{11}^{-1}X_1 + R'\mathrm{H}^{-1}R \right)^{-1}$$

$$= \lim\limits_{\mathrm{H}\to 0} \left\{ \left( X_1'\Sigma_{11}^{-1}X_1 \right)^{-1} - \left( X_1'\Sigma_{11}^{-1}X_1 \right)^{-1} R' \left[ \mathrm{H} + R\left( X_1'\Sigma_{11}^{-1}X_1 \right)^{-1} R' \right]^{-1} R\left( X_1'\Sigma_{11}^{-1}X_1 \right)^{-1} \right\}$$

$$= \left( X_1'\Sigma_{11}^{-1}X_1 \right)^{-1} - \left( X_1'\Sigma_{11}^{-1}X_1 \right)^{-1} R' \left[ 0 + R\left( X_1'\Sigma_{11}^{-1}X_1 \right)^{-1} R' \right]^{-1} R\left( X_1'\Sigma_{11}^{-1}X_1 \right)^{-1}$$

$$= \left( X_1'\Sigma_{11}^{-1}X_1 \right)^{-1} - \left( X_1'\Sigma_{11}^{-1}X_1 \right)^{-1} R' \left[ R\left( X_1'\Sigma_{11}^{-1}X_1 \right)^{-1} R' \right]^{-1} R\left( X_1'\Sigma_{11}^{-1}X_1 \right)^{-1}$$

$$= Var[\hat{\boldsymbol{\beta}}] - Var[\hat{\boldsymbol{\beta}}]R' \left[ R Var[\hat{\boldsymbol{\beta}}]R' \right]^{-1} R Var[\hat{\boldsymbol{\beta}}]$$

The variance of the constrained estimator is neatly expressed in terms of the variance of the

unconstrained estimator $Var[\hat{\boldsymbol{\beta}}] = \left( X_1'\Sigma_{11}^{-1}X_1 \right)^{-1}$. As a check:

$$Var[R\boldsymbol{\beta}^*] = R Var[\boldsymbol{\beta}^*]R' = R Var[\hat{\boldsymbol{\beta}}]R' - R Var[\hat{\boldsymbol{\beta}}]R' \left[ R Var[\hat{\boldsymbol{\beta}}]R' \right]^{-1} R Var[\hat{\boldsymbol{\beta}}]R' = 0$$

---

[8] In the following formulas the existence of the inverse of $R\left( X_1'\Sigma_{11}^{-1}X_1 \right)^{-1} R' = R Var[\hat{\boldsymbol{\beta}}]R'$ is guaranteed, since the variance matrix is positive-definite and R is of full-row rank.

In order to take the limit of $\hat{\boldsymbol{\beta}}(H)$ we need the following intermediate result:

$$Var\left[\hat{\boldsymbol{\beta}}(H)\right]R'H^{-1} = \left(X'\Sigma_{11}^{-1}X\right)^{-1}\left\{I_k - R'\left[H + R\left(X'\Sigma_{11}^{-1}X\right)^{-1}R'\right]^{-1}R\left(X'\Sigma_{11}^{-1}X\right)^{-1}\right\}R'H^{-1}$$

$$= \left(X'\Sigma_{11}^{-1}X\right)^{-1}\left\{R' - R'\left[H + R\left(X'\Sigma_{11}^{-1}X\right)^{-1}R'\right]^{-1}R\left(X'\Sigma_{11}^{-1}X\right)^{-1}R'\right\}H^{-1}$$

$$= \left(X'\Sigma_{11}^{-1}X\right)^{-1}R'\left\{I_j - \left[H + R\left(X'\Sigma_{11}^{-1}X\right)^{-1}R'\right]^{-1}R\left(X'\Sigma_{11}^{-1}X\right)^{-1}R'\right\}H^{-1}$$

$$= \left(X'\Sigma_{11}^{-1}X\right)^{-1}R'\left\{I_j - \left[H + R\left(X'\Sigma_{11}^{-1}X\right)^{-1}R'\right]^{-1}\left[H + R\left(X'\Sigma_{11}^{-1}X\right)^{-1}R' - H\right]\right\}H^{-1}$$

$$= \left(X'\Sigma_{11}^{-1}X\right)^{-1}R'\left\{I_j - I_j + \left[H + R\left(X'\Sigma_{11}^{-1}X\right)^{-1}R'\right]^{-1}H\right\}H^{-1}$$

$$= \left(X'\Sigma_{11}^{-1}X\right)^{-1}R'\left\{\left[H + R\left(X'\Sigma_{11}^{-1}X\right)^{-1}R'\right]^{-1}H\right\}H^{-1}$$

$$= \left(X'\Sigma_{11}^{-1}X\right)^{-1}R'\left[H + R\left(X'\Sigma_{11}^{-1}X\right)^{-1}R'\right]^{-1}$$

$$= Var\left[\hat{\boldsymbol{\beta}}\right]R'\left[H + R\,Var\left[\hat{\boldsymbol{\beta}}\right]R'\right]^{-1}$$

Therefore:

$$\boldsymbol{\beta}^* = \lim_{H \to 0}\hat{\boldsymbol{\beta}}(H)$$

$$= \lim_{H \to 0}Var\left[\hat{\boldsymbol{\beta}}(H)\right]\left(X_1'\Sigma_{11}^{-1}\mathbf{y}_1 + R'H^{-1}\mathbf{r}\right)$$

$$= \lim_{H \to 0}Var\left[\hat{\boldsymbol{\beta}}(H)\right]X_1'\Sigma_{11}^{-1}\mathbf{y}_1 + \lim_{H \to 0}Var\left[\hat{\boldsymbol{\beta}}(H)\right]R'H^{-1}\mathbf{r}$$

$$= Var\left[\boldsymbol{\beta}^*\right]X_1'\Sigma_{11}^{-1}\mathbf{y}_1 + \left(\lim_{H \to 0}Var\left[\hat{\boldsymbol{\beta}}(H)\right]R'H^{-1}\right)\left(\plim_{H \to 0}\mathbf{r}\right)$$

$$= Var\left[\boldsymbol{\beta}^*\right]X_1'\Sigma_{11}^{-1}\mathbf{y}_1 + \left(Var\left[\hat{\boldsymbol{\beta}}\right]R'\left[R\,Var\left[\hat{\boldsymbol{\beta}}\right]R'\right]^{-1}\right)\left(\plim_{Var[\boldsymbol{\eta}] \to 0}\left[R\boldsymbol{\beta} + \boldsymbol{\eta}\right]\right)$$

$$= Var\left[\boldsymbol{\beta}^*\right]X_1'\Sigma_{11}^{-1}\mathbf{y}_1 + \left(Var\left[\hat{\boldsymbol{\beta}}\right]R'\left[R\,Var\left[\hat{\boldsymbol{\beta}}\right]R'\right]^{-1}\right)R\boldsymbol{\beta}$$

$$= Var\left[\boldsymbol{\beta}^*\right]X_1'\Sigma_{11}^{-1}\mathbf{y}_1 + Var\left[\hat{\boldsymbol{\beta}}\right]R'\left[R\,Var\left[\hat{\boldsymbol{\beta}}\right]R'\right]^{-1}\mathbf{r}$$

$$= \left\{Var\left[\hat{\boldsymbol{\beta}}\right] - Var\left[\hat{\boldsymbol{\beta}}\right]R'\left[R\,Var\left[\hat{\boldsymbol{\beta}}\right]R'\right]^{-1}R\,Var\left[\hat{\boldsymbol{\beta}}\right]\right\}X_1'\Sigma_{11}^{-1}\mathbf{y}_1 + Var\left[\hat{\boldsymbol{\beta}}\right]R'\left[R\,Var\left[\hat{\boldsymbol{\beta}}\right]R'\right]^{-1}\mathbf{r}$$

$$= \left\{I_k - Var\left[\hat{\boldsymbol{\beta}}\right]R'\left[R\,Var\left[\hat{\boldsymbol{\beta}}\right]R'\right]^{-1}R\right\}Var\left[\hat{\boldsymbol{\beta}}\right]X_1'\Sigma_{11}^{-1}\mathbf{y}_1 + Var\left[\hat{\boldsymbol{\beta}}\right]R'\left[R\,Var\left[\hat{\boldsymbol{\beta}}\right]R'\right]^{-1}\mathbf{r}$$

$$= \left\{I_k - Var\left[\hat{\boldsymbol{\beta}}\right]R'\left[R\,Var\left[\hat{\boldsymbol{\beta}}\right]R'\right]^{-1}R\right\}\hat{\boldsymbol{\beta}} + Var\left[\hat{\boldsymbol{\beta}}\right]R'\left[R\,Var\left[\hat{\boldsymbol{\beta}}\right]R'\right]^{-1}\mathbf{r}$$

As a check:

$$\mathbf{R}\boldsymbol{\beta}^* = \mathbf{R}\left\{\mathbf{I}_k - Var[\hat{\boldsymbol{\beta}}]\mathbf{R}'[\mathbf{R}Var[\hat{\boldsymbol{\beta}}]\mathbf{R}']^{-1}\mathbf{R}\right\}\hat{\boldsymbol{\beta}} + \mathbf{R}Var[\hat{\boldsymbol{\beta}}]\mathbf{R}'[\mathbf{R}Var[\hat{\boldsymbol{\beta}}]\mathbf{R}']^{-1}\mathbf{r}$$

$$= \left\{\mathbf{R} - \mathbf{R}\,Var[\hat{\boldsymbol{\beta}}]\mathbf{R}'[\mathbf{R}Var[\hat{\boldsymbol{\beta}}]\mathbf{R}']^{-1}\mathbf{R}\right\}\hat{\boldsymbol{\beta}} + \mathbf{R}Var[\hat{\boldsymbol{\beta}}]\mathbf{R}'[\mathbf{R}Var[\hat{\boldsymbol{\beta}}]\mathbf{R}']^{-1}\mathbf{r}$$

$$= \left\{\mathbf{R} - \mathbf{I}_j\mathbf{R}\right\}\hat{\boldsymbol{\beta}} + \mathbf{I}_j\mathbf{r}$$

$$= \mathbf{r}$$

And so, to summarize, the solution of a β-constrained linear model is the solution of the unconstrained model with the substitution of $\boldsymbol{\beta}^*$ for $\hat{\boldsymbol{\beta}}$, where:

$$\boldsymbol{\beta}^* = \left\{\mathbf{I}_k - Var[\hat{\boldsymbol{\beta}}]\mathbf{R}'[\mathbf{R}Var[\hat{\boldsymbol{\beta}}]\mathbf{R}']^{-1}\mathbf{R}\right\}\hat{\boldsymbol{\beta}} + Var[\hat{\boldsymbol{\beta}}]\mathbf{R}'[\mathbf{R}Var[\hat{\boldsymbol{\beta}}]\mathbf{R}']^{-1}\mathbf{r}$$

$$Var[\boldsymbol{\beta}^*] = Var[\hat{\boldsymbol{\beta}}] - Var[\hat{\boldsymbol{\beta}}]\mathbf{R}'[\mathbf{R}Var[\hat{\boldsymbol{\beta}}]\mathbf{R}']^{-1}\mathbf{R}Var[\hat{\boldsymbol{\beta}}]$$

## 10. PARAMETER CONSTRAINTS AS PROJECTIONS

The formulas above for $\boldsymbol{\beta}^*$ and $Var[\boldsymbol{\beta}^*]$ may seem cumbersome, perhaps even repugnant. However, they become perspicuous when interpreted as a projection. From Section 6 we take the projection formula $P(\mathbf{x}; \{A\mathbf{x} = \mathbf{b}\}, \Sigma) = \left\{I_n - \Sigma A'(A\Sigma A')^{-1}A\right\}\mathbf{x} + \Sigma A'(A\Sigma A')^{-1}\mathbf{b}$, where $\Sigma = Var[\mathbf{x}]$. But now let the constraint space $\Omega$ be $\{\beta \in \Re^k : \mathbf{R}\beta = \mathbf{r}\}$. In this case:

$$\boldsymbol{\beta}^* = \left\{\mathbf{I}_k - Var[\hat{\boldsymbol{\beta}}]\mathbf{R}'[\mathbf{R}Var[\hat{\boldsymbol{\beta}}]\mathbf{R}']^{-1}\mathbf{R}\right\}\hat{\boldsymbol{\beta}} + Var[\hat{\boldsymbol{\beta}}]\mathbf{R}'[\mathbf{R}Var[\hat{\boldsymbol{\beta}}]\mathbf{R}']^{-1}\mathbf{r}$$

$$= P(\hat{\boldsymbol{\beta}}; \{\mathbf{R}\beta = \mathbf{r}\}, Var[\hat{\boldsymbol{\beta}}])$$

Hence, the constrained parameter estimator is the projection of the unconstrained estimator according to the metric of the variance of the unconstrained estimator. Just to corroborate, we see that the variance of the constrained estimator,

$Var[\boldsymbol{\beta}^*] = Var[\hat{\boldsymbol{\beta}}] - Var[\hat{\boldsymbol{\beta}}]R'[R\,Var[\hat{\boldsymbol{\beta}}]R']^{-1}R\,Var[\hat{\boldsymbol{\beta}}]$ accords with the projection variance

$Var[P(\mathbf{x}; \Omega, \Sigma)] = \Sigma - \Sigma A'(A\Sigma A')^{-1} A\Sigma$ .

Similarly to how we argued in Section 4, $Var[\hat{\boldsymbol{\beta}}]$ can be Cholesky-decomposed as

$Var[\hat{\boldsymbol{\beta}}] = QQ'$ for some non-singular $Q_{k \times k}$ . So the constrained variance can be factored as

$Var[\boldsymbol{\beta}^*] = Q\{I_k - Q'R'(RQQ'R')^{-1}RQ\}Q' = QMQ'$ , and its rank is that of the idempotent

matrix $M$ , whose rank equals its trace. Again, to continue as in Section 4:

$$Tr(M) = Tr(I_k) - Tr(Q'R'(RQQ'R')^{-1}RQ) = k - Tr((RQQ'R')^{-1}RQQ'R') = k - j$$

Therefore, $rank\left(Var[\boldsymbol{\beta}^*]\right) = k - j = rank\left(Var[\hat{\boldsymbol{\beta}}]\right) - j$ . The parameter constraint reduces

the parameter variance by $j$ degrees of freedom. In words, $R\boldsymbol{\beta}^*$ is a degenerate random

variable, or a constant. Certainly it is, since by the constraint $R\boldsymbol{\beta}^* = r$ .

All this shows that the solution of a parameter-constrained model is equivalent to the projection of the solution of an unconstrained model. There seems to be a certain commutativity between constraining/projecting and solving.

## 11. INFORMATION AS PROJECTION

We start with the equation of Section 7: $\mathbf{y} = X_{t \times k}\beta_{k \times 1} + \mathbf{e}$ , where $Var[\mathbf{e}] = \Sigma_{t \times t}$ . However,

let us suppose that $\beta$ is known and needs no estimation. Our best prediction of $\mathbf{y}$ is $X\beta$,

whose prediction-error variance is $Var[\mathbf{y} - \hat{\mathbf{y}}] = Var[\mathbf{y} - X\beta] = Var[\mathbf{e}] = \Sigma$ . At this stage we

are saying nothing more than $\mathbf{y} \sim (X\beta, \Sigma)$. But furthermore suppose that we have observed

$A_{m \times t}\mathbf{y}$, where A is of full-row rank. Name the observed value $\mathbf{b}_{m \times 1}$. The problem is to

predict $\mathbf{y}$ after the observation.

Since $\mathbf{y} = X\beta + \mathbf{e}$, $\mathbf{b} = A\mathbf{y} = AX\beta + A\mathbf{e}$, where $Var[A\mathbf{e}] = A\Sigma A'$. Since $\mathbf{y}$ is the same in the

observation as in the prediction, the observation covaries with the prediction; in fact,

$Cov[A\mathbf{y}, \mathbf{y}] = Cov[A\mathbf{e}, \mathbf{e}] = A\Sigma$. We can predict $\mathbf{y}$ according to the parameter-constrained

model of Section 9:

$$
\begin{bmatrix} \mathbf{b} \\ \hline \beta_0 \\ \hline \mathbf{y} \end{bmatrix} = \begin{bmatrix} AX \\ \hline I_k \\ \hline X \end{bmatrix} \beta + \begin{bmatrix} A\mathbf{e} \\ \hline 0 \\ \hline \mathbf{e} \end{bmatrix}, \text{ where } Var\begin{bmatrix} A\mathbf{e} \\ \hline 0 \\ \hline \mathbf{e} \end{bmatrix} = \begin{bmatrix} A\Sigma A' & 0 & A\Sigma \\ 0 & 0 & 0 \\ \hline \Sigma A' & 0 & \Sigma \end{bmatrix}
$$

Although this will work, a simpler and more appealing model can be constructed if one

allows for zero-dimensional matrices.[9] Because all $m \times 0$ and $0 \times n$ matrices are of rank zero,

$A_{m \times 0}B_{0 \times n} = 0_{m \times n}$. This is nothing more than the nullity of the empty summation operator,

i.e., $(AB)_{ij} = \sum_{k=1}^{0} (A)_{ik}(B)_{kj} = 0$. The simpler model is:

$$
\begin{bmatrix} \mathbf{b} - AX\beta \\ \hline \mathbf{y} - X\beta \end{bmatrix} = \begin{bmatrix} X_1 \;_{(m \times 0)} \\ \hline X_2 \;_{(t \times 0)} \end{bmatrix} \gamma_{0 \times 1} + \begin{bmatrix} A\mathbf{e} \\ \hline \mathbf{e} \end{bmatrix}, \text{ where } Var\begin{bmatrix} A\mathbf{e} \\ \hline \mathbf{e} \end{bmatrix} = \begin{bmatrix} A\Sigma A' & A\Sigma \\ \hline \Sigma A' & \Sigma \end{bmatrix}
$$

Its solution begins with:

---

[9] It is a windfall for a matrix language to allow for zeros in the dimensions of its arrays, as do APL, J, and R. SAS/IML does not; at least it did not in the late 1990s (version 7), when the author last used it.

$$\hat{\gamma}_{0\times1} = Var[\hat{\gamma}]X_1'(A\Sigma A')^{-1}(\mathbf{b} - AX\beta)$$
$$= \left(X_1'(A\Sigma A')^{-1}X_1\right)^{-1}X_1'(A\Sigma A')^{-1}(\mathbf{b} - AX\beta)$$
$$= (0\times0)^{-1}(0\times m)\cdot(m\times m)^{-1}(m\times1)$$
$$= (0\times0)^{-1}(0\times1)$$

The only thing to give pause here is the inverse of the $0\times0$ matrix. But the space of real 0-vectors, $\Re^0$, contains just one element, viz., the origin. It is closed under addition and multiplication ($0+0 = 0\times0 = 0$), and 0 serves as its identity element. So in $\Re^0$, $0^{-1} = 0$. Hence, $(0\times0)^{-1} = (0\times0)$. Therefore, $Var[\hat{\gamma}] = (0\times0)$ and $\hat{\gamma}_{0\times1} = (0\times1)$.[10] Finally:

$$\hat{\mathbf{y}} = X\beta + \widehat{\mathbf{y} - X\beta}$$
$$= X\beta + X_2\hat{\gamma} + \Sigma A'(A\Sigma A')^{-1}(\{\mathbf{b} - AX\beta\} - X_1\hat{\gamma})$$
$$= X\beta + 0_{t\times1} + \Sigma A'(A\Sigma A')^{-1}(\{\mathbf{b} - AX\beta\} - 0_{m\times1})$$
$$= X\beta + \Sigma A'(A\Sigma A')^{-1}(\mathbf{b} - AX\beta)$$

The variance of its prediction error is:

$$Var[\mathbf{y} - \hat{\mathbf{y}}] = Var[(\mathbf{y} - X\beta) - (\hat{\mathbf{y}} - X\beta)]$$
$$= \left(X_2 - \Sigma A'(A\Sigma A')^{-1}X_1\right)Var[\hat{\gamma}]\left(X_2 - \Sigma A'(A\Sigma A')^{-1}X_1\right)' + \Sigma - \Sigma A'(A\Sigma A')^{-1}A\Sigma$$
$$= (t\times0)(0\times0)(0\times t) + \Sigma - \Sigma A'(A\Sigma A')^{-1}A\Sigma$$
$$= 0_{t\times t} + \Sigma - \Sigma A'(A\Sigma A')^{-1}A\Sigma$$
$$= \Sigma - \Sigma A'(A\Sigma A')^{-1}A\Sigma$$

Except for slight notational differences, this solution is the same as that of Section 6. As long as no parameter needs to be estimated (or the parameter dimension is $0\times1$), the linear statistical model treats "$m$ dimensions" of prior information as a projection into a subspace of $t - m$ dimensions.

---

[10] To elaborate on the previous footnote, we have verified that APL, J, and R yield these results. Therefore, they correctly treat $(0\times0)^{-1}$ as $(0\times0)$.

## 12. COMBINING ESTIMATES

It is not uncommon for an actuary linearly to combine two or more unbiased estimators of the same quantity. Of course, it is desirable for the combination to be best. In the simplest situation of independent scalar estimators, the best combination uses weights inversely proportional to the variances of the estimators. But with Gauss-Markov theorem one can determine the best linear combination of vector estimators, even if they are not independent.

To frame the problem, suppose that we have $n$ unbiased estimators $\hat{\mathbf{y}}_i$ of the same $t{\times}1$ vector $\mathbf{y}$, as well as their $t{\times}t$ prediction-error variances $\Sigma_{ii} = Var[\mathbf{y} - \hat{\mathbf{y}}_i]$. Suppose also that we have the $t{\times}t$ prediction-error covariances $\Sigma_{ij} = Cov[\mathbf{y} - \hat{\mathbf{y}}_i, \mathbf{y} - \hat{\mathbf{y}}_j]$. Frequently the covariances are $0_{t \times t}$, but there are realistic exceptions. Stack the estimators and block their (co)variances:

$$\hat{\mathbf{Y}}_{nt \times 1} = \begin{bmatrix} \hat{\mathbf{y}}_1 \\ \vdots \\ \hat{\mathbf{y}}_n \end{bmatrix}, \quad Var[\mathbf{Y} - \hat{\mathbf{Y}}]_{nt \times nt} = \begin{bmatrix} \Sigma_{11} & \cdots & \Sigma_{1n} \\ \vdots & \ddots & \vdots \\ \Sigma_{n1} & \cdots & \Sigma_{nn} \end{bmatrix}$$

The variance matrix must be non-negative-definite; but we will assume it to be positive definite, hence invertible. If the weight given to $\hat{\mathbf{y}}_i$ is the $t{\times}t$ matrix $\mathbf{W}_i'$, the combined estimator will be $\hat{\mathbf{y}} = \sum_{i=1}^{n} \mathbf{W}_i' \hat{\mathbf{y}}_i = \begin{bmatrix} \mathbf{W}_1' & \cdots & \mathbf{W}_n' \end{bmatrix}_{t \times nt} \begin{bmatrix} \hat{\mathbf{y}}_1 \\ \vdots \\ \hat{\mathbf{y}}_n \end{bmatrix} = \mathbf{W}'\hat{\mathbf{Y}}$. In order for this estimator to be unbiased, $\begin{bmatrix} \mathbf{W}_1' & \cdots & \mathbf{W}_n' \end{bmatrix} \begin{bmatrix} \mathbf{I}_t \\ \vdots \\ \mathbf{I}_t \end{bmatrix} = \mathbf{W}' \begin{bmatrix} \mathbf{I}_t \\ \vdots \\ \mathbf{I}_t \end{bmatrix} = \mathbf{I}_t$. The transpose of this constraint is

$\begin{bmatrix} \mathbf{I}_t & \cdots & \mathbf{I}_t \end{bmatrix}_{t \times nt} \mathbf{W}_{nt \times t} = \mathbf{I}_t$. The best combination will minimize the combined prediction-error variance $\mathbf{W}' Var\begin{bmatrix} \mathbf{Y} - \hat{\mathbf{Y}} \end{bmatrix} \mathbf{W}$. Posing the problem in the proper form, we seek to minimize $\mathbf{W}' Var\begin{bmatrix} \mathbf{Y} - \hat{\mathbf{Y}} \end{bmatrix} \mathbf{W} = \mathbf{W}' \left( Var^{-1} \begin{bmatrix} \mathbf{Y} - \hat{\mathbf{Y}} \end{bmatrix} \right)^{-1} \mathbf{W}$ subject to $\begin{bmatrix} \mathbf{I}_t & \cdots & \mathbf{I}_t \end{bmatrix}_{t \times nt} \mathbf{W}_{nt \times t} = \mathbf{I}_t$.

According to the Gauss-Markov theorem:

$$\mathbf{W}^* = Var^{-1} \begin{bmatrix} \mathbf{Y} - \hat{\mathbf{Y}} \end{bmatrix} \begin{bmatrix} \mathbf{I}_t \\ \vdots \\ \mathbf{I}_t \end{bmatrix} \left( \begin{bmatrix} \mathbf{I}_t & \cdots & \mathbf{I}_t \end{bmatrix} Var^{-1} \begin{bmatrix} \mathbf{Y} - \hat{\mathbf{Y}} \end{bmatrix} \begin{bmatrix} \mathbf{I}_t \\ \vdots \\ \mathbf{I}_t \end{bmatrix} \right)^{-1} \mathbf{I}_t$$

$$= Var^{-1} \begin{bmatrix} \mathbf{Y} - \hat{\mathbf{Y}} \end{bmatrix} \begin{bmatrix} \mathbf{I}_t \\ \vdots \\ \mathbf{I}_t \end{bmatrix} \left( \begin{bmatrix} \mathbf{I}_t & \cdots & \mathbf{I}_t \end{bmatrix} Var^{-1} \begin{bmatrix} \mathbf{Y} - \hat{\mathbf{Y}} \end{bmatrix} \begin{bmatrix} \mathbf{I}_t \\ \vdots \\ \mathbf{I}_t \end{bmatrix} \right)^{-1}$$

If the covariances $\Sigma_{i \neq j}$ are zero, this simplifies to:

$$\mathbf{W}^* = \begin{bmatrix} \Sigma_{11}^{-1} \\ \vdots \\ \Sigma_{nn}^{-1} \end{bmatrix} \left( \sum_{i=1}^{n} \Sigma_{ii}^{-1} \right)^{-1}$$

It is recognizable as the matrix version of the well known rule of weighting independent scalar estimates inversely proportionally to their variances.[11] Appendix B will provide a simple example of covarying estimates, and will outline its importance to conjoint modeling, or to modeling in which ultimate paid and incurred losses must be equal.

---

[11] Unlike scalar weighting, a matrix-weighted average can fall outside its extremes, e.g.:

$$\begin{bmatrix} 0.50 & 0.25 \\ 0.25 & 0.40 \end{bmatrix} \begin{bmatrix} 400 \\ 420 \end{bmatrix} + \begin{bmatrix} 0.50 & -0.25 \\ -0.25 & 0.60 \end{bmatrix} \begin{bmatrix} 440 \\ 400 \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 425 \\ 398 \end{bmatrix}, \text{ yet } 398 \notin [400, 420].$$ This is due to non-zero off-diagonal amounts ($\pm$ 0.25) in the weighting matrices. In practice, such amounts are relatively small, and the matrix-weighted averages lie within their extremes. Cf. Judge [1988, 287].

## 13. CONCLUSION

The Gauss-Markov theorem is truly profound. It provides a lucid basis for solving a wide range of modeling and estimation problems, even within the rudimentary matrix functionality of Excel. As the many sections of this paper have demonstrated, it deserves to be liberated from being an appendage to the least-squares approach to linear statistical modeling.[12]

---

[12] For a brief history of the least-squares method and the true relation of the Gauss-Markov theorem to it see Appendix C.

# REFERENCES

[1.] Halliwell, Leigh, "Conjoint Prediction of Paid and Incurred Losses," *1997 Loss Reserving Discussion Papers*, Casualty Actuarial Society, 1997, 241-379, www.casact.org/pubs/forum/97sforum/97sf1241.pdf.

[2.] Halliwell, Leigh, "Statistical Models and Credibility," *CAS Forum* (Winter 1998), www.casact.org/pubs/forum/98wforum/98wf061.pdf.

[3.] Halliwell, Leigh, "Chain-Ladder Bias: Its Reason and Meaning," *Variance*, 1:2, 2007, 214-247, www.variancejournal.org/issues/01-02/214.pdf.

[4.] Judge, George G., Hill, R. C., *et al.*, *Introduction to the Theory and Practice of Econometrics* (Second Edition), New York, John Wiley & Sons, 1988.

[5.] Wikipedia contributors, "Best linear unbiased prediction," *Wikipedia*, http://en.wikipedia.org/wiki/Best_linear_unbiased_prediction (accessed September 2015).

[6.] Wikipedia contributors, "Gauss–Markov theorem," *Wikipedia*, http://en.wikipedia.org/wiki/Gauss-Markov_theorem (accessed September 2015).

[7.] Wikipedia contributors, "Least squares," *Wikipedia,* http://en.wikipedia.org/wiki/Least_squares (accessed September 2015)

[8.] Wikipedia contributors, "Mahalanobis distance," *Wikipedia*, http://en.wikipedia.org/wiki/Mahalanobis_distance (accessed September 2015).

## APPENDIX A

## GEOMETRIC MATTERS CONCERNING VECTORS IN $\Re^n$

In this appendix we will interpret the vector dot product and prove the triangle inequality.

For $x, y \in \Re^n$ the dot product $x \bullet y = x'y = \sum_{i=1}^{n} x_i y_i$. However, this is easily generalized

with a $\Sigma$ metric as $x'\Sigma^{-1}y$. The $\Sigma$-metric triangle inequality is:

$$\sqrt{(x+y)'\Sigma^{-1}(x+y)} \leq \sqrt{x'\Sigma^{-1}x} + \sqrt{y'\Sigma^{-1}y}$$

As for the dot product, let $\hat{y}$ be the $\Sigma$ projection of y onto x. If $x_{n\times 1} \neq 0$, then x is of full-

column rank, $x'\Sigma^{-1}x$ is 1×1 positive-definite, and $(x'\Sigma^{-1}x)^{-1}$ exists. So vector y will $\Sigma$-

project as some multiple of x, or $\hat{y} = x\beta$. From Section 7, $\beta = (x'\Sigma^{-1}x)^{-1}x'\Sigma^{-1}y$. Hence,

$\hat{y} = x(x'\Sigma^{-1}x)^{-1}x'\Sigma^{-1}y$. Accordingly, $x'\Sigma^{-1}y = x'\Sigma^{-1}x(x'\Sigma^{-1}x)^{-1}x'\Sigma^{-1}y = x'\Sigma^{-1}\hat{y}$. So the $\Sigma$-

metric dot product of two vectors is equal to $\Sigma$-metric dot product of one vector and the $\Sigma$

projection of the other onto it. Although $(x'\Sigma^{-1}x)^{-1}$ does not exist if $x = 0$, we know that

the projection of any vector onto 0 is 0. Hence, our geometric interpretation of the dot

product is valid for all x and y. For the Euclidean metric $\Sigma = I_n$, the projection is the

perpendicular, and $|\hat{y}| = |y|cos\theta$, where $\theta$ is the angle between the two vectors with vertex at

0. From this follows the well-known formula $x \bullet y = x'y = |x||y|cos\theta$.

In preparation for the triangle inequality, since $\Sigma$ is positive-definite, the following 2×2 symmetric matrix is non-negative-definite:

$$\begin{bmatrix} x_{n\times 1} & y_{n\times 1} \end{bmatrix}' \Sigma^{-1} \begin{bmatrix} x & y \end{bmatrix} = \begin{bmatrix} x' \\ y' \end{bmatrix} \Sigma^{-1} \begin{bmatrix} x & y \end{bmatrix} = \begin{bmatrix} x'\Sigma^{-1}x & x'\Sigma^{-1}y \\ y'\Sigma^{-1}x & y'\Sigma^{-1}y \end{bmatrix}$$

It is a theorem that the determinant of a non-negative-definite matrix is non-negative; but we can readily prove it here for the 2×2 case. Such a matrix can be Cholesky factored as

$$\begin{bmatrix} a & 0 \\ b & d \end{bmatrix}\begin{bmatrix} a & b \\ 0 & d \end{bmatrix}$$, for real numbers $a$, $b$, and $c$. This equals $\begin{bmatrix} a^2 & ab \\ ba & b^2+d^2 \end{bmatrix}$, whose

determinant is $a^2(b^2+d^2) - a^2 b^2 = a^2 d^2$, which must be greater than or equal to zero.

Now let '~' stand for the relationship in the triangle inequality:

$$\sqrt{(x+y)'\Sigma^{-1}(x+y)} \sim \sqrt{x'\Sigma^{-1}x} + \sqrt{y'\Sigma^{-1}y}$$

Because the quantities under all the radical signs are non-negative, the following transformations will not affect the relationship:

$$(x+y)'\Sigma^{-1}(x+y) \sim x'\Sigma^{-1}x + 2\sqrt{x'\Sigma^{-1}x}\sqrt{y'\Sigma^{-1}y} + y'\Sigma^{-1}y$$

$$x'\Sigma^{-1}x + x'\Sigma^{-1}y + y'\Sigma^{-1}x + y'\Sigma^{-1}y \sim x'\Sigma^{-1}x + 2\sqrt{x'\Sigma^{-1}x}\sqrt{y'\Sigma^{-1}y} + y'\Sigma^{-1}y$$

$$x'\Sigma^{-1}y + y'\Sigma^{-1}x \sim 2\sqrt{x'\Sigma^{-1}x}\sqrt{y'\Sigma^{-1}y}$$

$$2x'\Sigma^{-1}y \sim 2\sqrt{x'\Sigma^{-1}x}\sqrt{y'\Sigma^{-1}y}$$

$$x'\Sigma^{-1}y \sim \sqrt{x'\Sigma^{-1}x}\sqrt{y'\Sigma^{-1}y}$$

$$(x'\Sigma^{-1}y)^2 \sim x'\Sigma^{-1}x \cdot y'\Sigma^{-1}y$$

$$x'\Sigma^{-1}y \cdot y'\Sigma^{-1}x \sim x'\Sigma^{-1}x \cdot y'\Sigma^{-1}y$$

$$0 \sim x'\Sigma^{-1}x \cdot y'\Sigma^{-1}y - x'\Sigma^{-1}y \cdot y'\Sigma^{-1}x$$

But the expression on the right of the last line is the determinant of our 2×2 non-negative-definite matrix. Therefore '~' is '≤'. Thus have we proven the triangle inequality in $\mathfrak{R}^n$ for every valid $\Sigma$ metric.

## APPENDIX B

## COVARYING ESTIMATORS AND CONJOINT PREDICTION

This appendix furnishes a simple, but not too contrived, example of combining estimators that are not independent. Let $X_i \sim [\mu, \sigma^2]$ be independent random variables. Our task will be to estimate the mean $\mu$. However, we must estimate it from two known statistics, $Y_1 = (X_1 + X_2 + X_3)/3$ and $Y_2 = (X_3 + X_4)/2$. Four $X$ variables have been melded into two $Y$ variables: $Y_1 \sim [\mu, \sigma^2/3]$ and $Y_2 \sim [\mu, \sigma^2/2]$. But since $X_3$ is common to both, they are not independent; rather, $Cov[Y_1, Y_2] = Cov[X_3/3, X_3/2] = \sigma^2/6$. So the first two moments of the **y** vector are:

$$\mathbf{y} = \begin{bmatrix} Y_1 \\ Y_2 \end{bmatrix} \sim \left( \begin{bmatrix} \mu \\ \mu \end{bmatrix}, \sigma^2 \begin{bmatrix} 1/3 & 1/6 \\ 1/6 & 1/2 \end{bmatrix} \right)$$

The problem is to minimize $\mathbf{W}'\sigma^2 \begin{bmatrix} 1/3 & 1/6 \\ 1/6 & 1/2 \end{bmatrix} \mathbf{W}$ subject to $\begin{bmatrix} 1 & 1 \end{bmatrix} \mathbf{W}_{2\times 1} = \mathbf{I}_1$. By the Gauss-Markov theorem (the $\sigma^2$ cancels, so it's omitted):

$$\mathbf{W}^* = \begin{bmatrix} 1/3 & 1/6 \\ 1/6 & 1/2 \end{bmatrix}^{-1} \begin{bmatrix} 1 \\ 1 \end{bmatrix} \left( \begin{bmatrix} 1 & 1 \end{bmatrix} \begin{bmatrix} 1/3 & 1/6 \\ 1/6 & 1/2 \end{bmatrix}^{-1} \begin{bmatrix} 1 \\ 1 \end{bmatrix} \right)^{-1} \mathbf{I}_1$$

$$= \frac{1}{5/36} \begin{bmatrix} 1/2 & -1/6 \\ -1/6 & 1/3 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} \left( \begin{bmatrix} 1 & 1 \end{bmatrix} \frac{1}{5/36} \begin{bmatrix} 1/2 & -1/6 \\ -1/6 & 1/3 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} \right)^{-1}$$

$$= \begin{bmatrix} 1/2 & -1/6 \\ -1/6 & 1/3 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} \left( \begin{bmatrix} 1 & 1 \end{bmatrix} \begin{bmatrix} 1/2 & -1/6 \\ -1/6 & 1/3 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} \right)^{-1} = \begin{bmatrix} 1/3 \\ 1/6 \end{bmatrix} (1/2)^{-1} = \begin{bmatrix} 2/3 \\ 1/3 \end{bmatrix}$$

So the minimal variance results from combining in a 2:1 = 10:5 ratio. One who ignored the covariance would weight them in a 3:2 = 9:6 ratio, underweighting the first and overweighting the second. The minimal variance itself is:

$$\begin{bmatrix} 2/3 & 1/3 \end{bmatrix} \sigma^2 \begin{bmatrix} 1/3 & 1/6 \\ 1/6 & 1/2 \end{bmatrix} \begin{bmatrix} 2/3 \\ 1/3 \end{bmatrix} = \sigma^2 \frac{5}{18} = \sigma^2 \cdot 0.2\overline{7}$$

Since $\sigma^2/4 < \sigma^2 \cdot 0.2\overline{7} =< \sigma^2/3$, the informational value of the two $Y$ statistics lies in between the informational values of three and four independent $X$ statistics.

As for conjoint prediction, the following model combines submodels *a* and *b*:

$$\begin{bmatrix} \mathbf{y}_{a1} \\ \mathbf{y}_{b1} \\ \hline \mathbf{y}_{a2} \\ \mathbf{y}_{b2} \end{bmatrix} = \begin{bmatrix} \mathbf{X}_{a1} & 0 \\ 0 & \mathbf{X}_{b1} \\ \hline \mathbf{X}_{a2} & 0 \\ 0 & \mathbf{X}_{b2} \end{bmatrix} \begin{bmatrix} \beta_a \\ \beta_b \end{bmatrix} + \begin{bmatrix} \mathbf{e}_{a1} \\ \mathbf{e}_{b1} \\ \hline \mathbf{e}_{a2} \\ \mathbf{e}_{b2} \end{bmatrix}, \text{where } Var \begin{bmatrix} \mathbf{e}_{a1} \\ \mathbf{e}_{b1} \\ \hline \mathbf{e}_{a2} \\ \mathbf{e}_{b2} \end{bmatrix} = \begin{bmatrix} \Sigma_{a11} & 0 & \Sigma_{a12} & 0 \\ 0 & \Sigma_{b11} & 0 & \Sigma_{b12} \\ \hline \Sigma_{a21} & 0 & \Sigma_{a22} & 0 \\ 0 & \Sigma_{b21} & 0 & \Sigma_{b22} \end{bmatrix}$$

One who works through the formulas of Section 8 will find that the solution of the combination is identical to the combination of the separate solutions (Halliwell [1998, Appendix A]). Were it not for this good fortune, one would have to model everything in order to model anything. So this combination is trivial; although the submodels are written down together, they interact neither in the design matrix nor in the variance structure. But conjoint prediction (Halliwell [1997]) makes use of the fact that paid losses (model *a*) and incurred losses (model *b*) must ultimately be equal by exposure period. This constrains the variance matrix; the sums of the paid and the incurred errors of each exposure period must be equal. But additionally, it imposes a restriction on the parameters. The *a priori*, or prior-to-any-observation, expected values are $E[\mathbf{y}_a] = \mathbf{X}_a \beta_a$ and $E[\mathbf{y}_b] = \mathbf{X}_b \beta_b$. The exposure-period sums of these paid and incurred vectors must also be equal. A "semi-conjoint"

model adds the appropriate β constraint to the trivial combination:

$$R\beta = \left( Q \begin{bmatrix} X_a & 0 \\ 0 & X_b \end{bmatrix} \right) \cdot \begin{bmatrix} \beta_a \\ \beta_b \end{bmatrix} = 0 = r \; ; \text{ consequently, } \hat{\beta}_a \text{ and } \hat{\beta}_b \text{ will covary. This parameter}$$

covariance will introduce covariance, or off-diagonal blocks, into $Var \begin{bmatrix} \mathbf{y}_{a2} - \hat{\mathbf{y}}_{a2} \\ \mathbf{y}_{b2} - \hat{\mathbf{y}}_{b2} \end{bmatrix}$.

According to Section 12, one may best combine the semi-constrained solutions $\begin{bmatrix} \hat{\mathbf{y}}_{a2} \\ \hat{\mathbf{y}}_{b2} \end{bmatrix}$ and

$Var \begin{bmatrix} \mathbf{y}_{a2} - \hat{\mathbf{y}}_{a2} \\ \mathbf{y}_{b2} - \hat{\mathbf{y}}_{b2} \end{bmatrix}$ according to the linear constraint that exposure-period sums of paid and

incurred losses are equal. Equivalently, in terms of Section 11, one can project the semi-constrained solution into the subspace of the constraint. Although a proof of this has so far eluded us, it works with examples. So the Gauss-Markov theorem seems to allow modeling temporarily to ignore variance restrictions in order to arrive at a tentative solution that can rather easily be collapsed by the hitherto ignored restrictions into the desired solution. This is the meaning of the sentence at the end of Section 10: "There seems to be a certain commutativity between constraining/projecting and solving." Conjoint prediction by collapsing a semi-conjoint model is much easier than fully conjoint prediction; it requires no eigen-decomposition, and is amenable to a spreadsheet solution.

## APPENDIX C

## LEAST-SQUARES VERSUS GAUSS-MARKOV

Many, probably most, actuaries think in terms of linear regression, rather than in terms of linear modeling. The standard linear-regression problem begins with $t$ observed quantities $y_j$. Each observation is associated with a $k$-tuple of known variables $(x_{j1},\ldots,x_{jk})$, on which the observation is believed linearly to depend, i.e., $y_j = x_{j1}\beta_1 + \ldots + x_{jk}\beta_k$. Of course, if $t = k$ and the $k$-tuples are linearly independent, one is merely solving simultaneous equations for the $\beta_j$. The regression problem arises when $t > k$, and the equations are approximate: $y_j \approx x_{j1}\beta_1 + \ldots + x_{jk}\beta_k$. One then needs to find the values of $\beta_j$ that make $x_{j1}\beta_1 + \ldots + x_{jk}\beta_k$ most closely approximate the $y_j$. A reasonable method, called "least squares," is to find the $\beta_j$ that minimize the sum of the squared errors, i.e., to minimize

$$f(\beta_1,\ldots,\beta_k) = \sum_{j=1}^{t}(y_j - x_{j1}\beta_1 + \ldots + x_{jk}\beta_k)^2 .$$ This is a problem well within the capability of a first-year calculus student.

The least-squares criterion for fitting, or "regressing," the best line to data first appeared in print in 1805, when Legendre published his Nouvelles méthodes pour la détermination des orbites des comètes. Earlier, in 1801, Gauss had applied the method to predict the reappearance of Ceres, which had just been discovered and then lost. However, he did not publish the method until 1806 in his Theoria Motus Corporum Coelestium in sectionibus conicis solem ambientium. Apparently, he did not refer to Legendre; and in the ensuing controversy over priority Gauss insisted that he had worked out the method at least as early

as 1795, when at the age of eighteen he entered the University of Göttingen.[13]  The relevant point of this interesting story is that in this early astronomical setting, the least-squares method was not statistical modeling.  It was applied to deterministically moving objects (comets and the newly discovered asteroids).  All uncertainty stemmed from the imprecision of the astronomers.  But for the first time it was realized that many economical but "fuzzy" observations could be more useful than one costly but "sharp" observation.

Gradually the approximate equations were turned into exact ones with random error terms: $\boldsymbol{y}_j = x_{j1}\beta_1 + \ldots + x_{jk}\beta_k + \boldsymbol{e}_j$.  Gauss himself in 1822 stated the optimality of the least-squares method, an early form of BLUE.  So today we talk of the "Gauss-Markov" theorem because Gauss started it.  But the linear algebra and statistical theory that developed after his death in 1855 culminated in the work of Andrey Markov (1856-1922).  Even today it is common for students to be introduced into linear modeling by way of least squares; many texts still refer to the matrix formula $\boldsymbol{\beta} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$ as the OLS ("Ordinary Least Squares") estimator.

How does the least-squares method differ from our version of the Gauss-Markov theorem? To put it in modern terms, both deal with estimating the β parameter in the model of Section 7: $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}$, where $Var[\mathbf{e}] = \Sigma$.  But instead of finding the $t \times k$ matrix W that will make estimator $\hat{\boldsymbol{\beta}} = \mathbf{W}'\mathbf{y}$ unbiased and of minimal variance, as per the Gauss-Markov theorem, the least-squares method seeks the value of β for which Xβ most closely

---

[13] Wikipedia "Least squares" gives an excellent account of this history, which is also recounted in many histories of mathematics.  However, there is slight disagreement about some of the dates.  Most historians cede the priority to Gauss.

approximates $\mathbf{y}$. "Closeness" here requires distance as measured by the $\Sigma$ metric. So the least-squares problem is to minimize $f(\beta) = (\mathbf{y} - \mathbf{X}\beta)' \Sigma^{-1} (\mathbf{y} - \mathbf{X}\beta)$. As with the Gauss-Markov theorem, X must be of full-column rank and $\Sigma$ must be positive-definite. But the two approaches are not logically equivalent. Although they yield the same answer, BLUE is *a posteriori* to the least-squares answer, whereas it is *a priori* to the Gauss-Markov.

The usual approach to the minimization is by means of multivariate calculus:

$$\frac{\partial f}{\partial \beta} = -2\mathbf{X}'\Sigma^{-1}(\mathbf{y} - \mathbf{X}\beta)$$

$$\frac{\partial^2 f}{\partial \beta \partial \beta'} = 2\mathbf{X}'\Sigma^{-1}\mathbf{X}$$

Setting the first derivative to $\mathbf{0}_{k\times 1}$, we derive $\tilde{\boldsymbol{\beta}} = \left(\mathbf{X}'\Sigma^{-1}\mathbf{X}\right)^{-1}\mathbf{X}'\Sigma^{-1}\mathbf{y}$. Since the second derivative is positive-definite, the critical value $\tilde{\boldsymbol{\beta}}$ is a minimum, as desired. However, since vector differentiation is unfamiliar to many (cf. Judge [1988, Appendix A.16]), we will solve the problem algebraically:

$$
\begin{aligned}
f(\beta) &= (\mathbf{y} - \mathbf{X}\beta)' \Sigma^{-1} (\mathbf{y} - \mathbf{X}\beta) \\
&= \left(\mathbf{y} - \mathbf{X}\tilde{\boldsymbol{\beta}} + \mathbf{X}\left[\tilde{\boldsymbol{\beta}} - \beta\right]\right)' \Sigma^{-1} \left(\mathbf{y} - \mathbf{X}\tilde{\boldsymbol{\beta}} + \mathbf{X}\left[\tilde{\boldsymbol{\beta}} - \beta\right]\right) \\
&= (\mathbf{y} - \mathbf{X}\tilde{\boldsymbol{\beta}})' \Sigma^{-1} (\mathbf{y} - \mathbf{X}\tilde{\boldsymbol{\beta}}) + (\mathbf{y} - \mathbf{X}\tilde{\boldsymbol{\beta}})' \Sigma^{-1} \mathbf{X}\left[\tilde{\boldsymbol{\beta}} - \beta\right] + \left[\tilde{\boldsymbol{\beta}} - \beta\right]' \mathbf{X}'\Sigma^{-1}(\mathbf{y} - \mathbf{X}\tilde{\boldsymbol{\beta}}) + \left[\tilde{\boldsymbol{\beta}} - \beta\right]' \mathbf{X}'\Sigma^{-1}\mathbf{X}\left[\tilde{\boldsymbol{\beta}} - \beta\right] \\
&= f(\tilde{\boldsymbol{\beta}}) + 2\left[\tilde{\boldsymbol{\beta}} - \beta\right]' \mathbf{X}'\Sigma^{-1}(\mathbf{y} - \mathbf{X}\tilde{\boldsymbol{\beta}}) + \left[\tilde{\boldsymbol{\beta}} - \beta\right]' \mathbf{X}'\Sigma^{-1}\mathbf{X}\left[\tilde{\boldsymbol{\beta}} - \beta\right] \\
&= f(\tilde{\boldsymbol{\beta}}) + 2\left[\tilde{\boldsymbol{\beta}} - \beta\right]' \left\{\mathbf{X}'\Sigma^{-1}\mathbf{y} - \mathbf{X}'\Sigma^{-1}\mathbf{X}\tilde{\boldsymbol{\beta}}\right\} + \left[\tilde{\boldsymbol{\beta}} - \beta\right]' \mathbf{X}'\Sigma^{-1}\mathbf{X}\left[\tilde{\boldsymbol{\beta}} - \beta\right] \\
&= f(\tilde{\boldsymbol{\beta}}) + 2\left[\tilde{\boldsymbol{\beta}} - \beta\right]' \left\{\mathbf{X}'\Sigma^{-1}\mathbf{y} - \mathbf{X}'\Sigma^{-1}\mathbf{X}\left(\mathbf{X}'\Sigma^{-1}\mathbf{X}\right)^{-1}\mathbf{X}'\Sigma^{-1}\mathbf{y}\right\} + \left[\tilde{\boldsymbol{\beta}} - \beta\right]' \mathbf{X}'\Sigma^{-1}\mathbf{X}\left[\tilde{\boldsymbol{\beta}} - \beta\right] \\
&= f(\tilde{\boldsymbol{\beta}}) + 2\left[\tilde{\boldsymbol{\beta}} - \beta\right]' \left\{\mathbf{X}'\Sigma^{-1}\mathbf{y} - \mathbf{X}'\Sigma^{-1}\mathbf{y}\right\} + \left[\tilde{\boldsymbol{\beta}} - \beta\right]' \mathbf{X}'\Sigma^{-1}\mathbf{X}\left[\tilde{\boldsymbol{\beta}} - \beta\right] \\
&= f(\tilde{\boldsymbol{\beta}}) + \left[\tilde{\boldsymbol{\beta}} - \beta\right]' \mathbf{X}'\Sigma^{-1}\mathbf{X}\left[\tilde{\boldsymbol{\beta}} - \beta\right] \\
&\geq f(\tilde{\boldsymbol{\beta}})
\end{aligned}
$$

As in Section 2, the last line is to be taken in a matrix-definite sense. Moreover, since $X'\Sigma^{-1}X$ is positive-definite, the inequality is strict except for $\beta = \tilde{\boldsymbol{\beta}}$. Therefore, $\tilde{\boldsymbol{\beta}}$ uniquely minimizes $f(\beta) = (\mathbf{y} - X\beta)' \Sigma^{-1}(\mathbf{y} - X\beta)$. Geometrically, the least-squares method drops a $\Sigma$ perpendicular from $\mathbf{y}$ to the linear subspace swept by $X\beta$.

But now that we have a "least-squares" estimator, we must check its "BLUE-ness." This is the meaning of the sentence above, that BLUE is *a posteriori* to the least-squares answer. Of course, from our *a priori* Gauss-Markov approach, we already know it to be BLUE, since it is identical to the Section 7 formula $\hat{\boldsymbol{\beta}} = \left(X'\Sigma^{-1}X\right)^{-1} X'\Sigma^{-1}\mathbf{y}$. If $\tilde{\boldsymbol{\beta}}$ were not identical to $\hat{\boldsymbol{\beta}}$, $\tilde{\boldsymbol{\beta}}$ would be either biased or not as good as $\hat{\boldsymbol{\beta}}$; it would lack either the 'B' or the 'U' of BLUE.

Finally, despite the historical development from least squares to Gauss-Markov, this is neither a "distinction without a difference" nor a matter of taste. Developing the theory of linear statistical modeling from our Gauss-Markov theorem allows us cleanly to solve problems that the least-squares approach can solve only with difficulty, if at all – such problems as predicting (Section 8), constraining (Section 9), projecting (Section 10), incorporating prior information (Section 11), and combining estimates (Section 12).