

Covariance Matrices and Precision Matrices

The matrix of central second moments of a random vector \mathbf{b} is $\text{Var}(\mathbf{b})$. This is called the **covariance matrix**, **variance matrix**, or **variance-covariance matrix**,

$$\text{Var}(\mathbf{b}) \equiv \text{E}\left((\mathbf{b} - \text{E}(\mathbf{b}))(\mathbf{b} - \text{E}(\mathbf{b}))^\top\right). \quad (1)$$

When $\text{E}(\mathbf{b}) = \mathbf{0}$, $\text{Var}(\mathbf{b}) = \text{E}(\mathbf{b}\mathbf{b}^\top)$.

The i^{th} diagonal element of $\text{Var}(\mathbf{b})$ is $\text{Var}(b_i) = \text{E}(b_i - \text{E}(b_i))^2$.

The ij^{th} off-diagonal element of $\text{Var}(\mathbf{b})$ is

$$\text{Cov}(b_i, b_j) \equiv \text{E}\left((b_i - \text{E}(b_i))(b_j - \text{E}(b_j))\right). \quad (2)$$

If $i = j$, $\text{Cov}(b_i, b_j) = \text{Var}(b_i)$.

Since $\text{Cov}(b_i, b_j) = \text{Cov}(b_j, b_i)$, $\text{Var}(\mathbf{b})$ must be a symmetric matrix.

If b_i and b_j are statistically independent, then $\text{Cov}(b_i, b_j) = 0$. The converse is not true, however.

The inverse of a covariance matrix is a **precision matrix**.

The **correlation** between b_i and b_j is

$$\rho(b_i, b_j) \equiv \frac{\text{Cov}(b_i, b_j)}{(\text{Var}(b_i) \text{Var}(b_j))^{1/2}}, \quad \text{with } -1 \leq \rho(b_i, b_j) \leq 1. \quad (3)$$

Correlations can be arranged into a symmetric **correlation matrix** with all elements on the principal diagonal equal to 1.

$\text{Var}(\mathbf{b})$ must be **positive semidefinite** as well as symmetric.

In most cases, covariance matrices and correlation matrices are actually **positive definite**.

A $k \times k$ symmetric matrix \mathbf{A} is positive definite if, for all nonzero k -vectors \mathbf{x} , the matrix product $\mathbf{x}^\top \mathbf{A} \mathbf{x}$ is positive.

$\mathbf{x}^\top \mathbf{A} \mathbf{x}$ is called a **quadratic form**.

A quadratic form always involves a k -vector, in this case \mathbf{x} , and a $k \times k$ matrix, in this case \mathbf{A} .

$$\mathbf{x}^\top \mathbf{A} \mathbf{x} = \sum_{i=1}^k \sum_{j=1}^k x_i x_j A_{ij}. \quad (4)$$

If this can be zero but not negative, then \mathbf{A} is said to be positive semidefinite.

Any matrix of the form $\mathbf{B}^\top \mathbf{B}$ is positive semidefinite. $\mathbf{B}^\top \mathbf{B}$ is symmetric and, for any nonzero \mathbf{x} ,

$$\mathbf{x}^\top \mathbf{B}^\top \mathbf{B} \mathbf{x} = (\mathbf{B} \mathbf{x})^\top (\mathbf{B} \mathbf{x}) = \|\mathbf{B} \mathbf{x}\|^2 \geq 0. \quad (5)$$

This is positive unless $\mathbf{B} \mathbf{x} = \mathbf{0}$. In that case, since $\mathbf{x} \neq \mathbf{0}$, the columns of \mathbf{B} are linearly dependent, and \mathbf{B} does not have **full column rank**.

- B can have full rank but not full column rank if B has fewer rows than columns. If so, maximum possible rank is number of rows.
- When B does have full column rank, $B^\top B$ is positive definite.
- If A is positive definite, then any matrix $B^\top A B$ is positive definite if B has full column rank, positive semidefinite otherwise.
- A positive definite matrix cannot be singular, because, if A is singular, there must exist a nonzero x such that $Ax = \mathbf{0}$. But then $x^\top Ax = 0$ as well, which means that A is not positive definite.
- Diagonal elements of a positive definite matrix must all be positive. Suppose that A_{22} were negative. Then, if we chose x to be the vector e_2 , the quadratic form would just be $e_2^\top A e_2 = A_{22} < 0$.
- For a positive semidefinite matrix, diagonal elements may be 0.
- The off-diagonal elements of A may be of either sign.
- The inverse of a positive definite matrix always exists, and it is positive definite.

The identity matrix, \mathbf{I} , is a positive definite matrix, since

$$\mathbf{x}^\top \mathbf{I} \mathbf{x} = \sum_{i=1}^k x_i^2. \quad (6)$$

If the $k \times k$ matrix \mathbf{A} is symmetric and positive definite, then there always exists a full-rank $k \times k$ matrix \mathbf{B} such that $\mathbf{A} = \mathbf{B}^\top \mathbf{B}$.

Think of \mathbf{B} as a **square root matrix**. For any matrix \mathbf{A} , the corresponding matrix \mathbf{B} is not unique.

In particular, \mathbf{B} can be chosen to be symmetric, but it can also be chosen to be upper or lower triangular.

We can compute a triangular \mathbf{B} using the **Cholesky decomposition**.

Since a triangular matrix is extremely easy to invert, we can use the Cholesky decomposition to find $\mathbf{A}^{-1} = \mathbf{B}^{-1}(\mathbf{B}^\top)^{-1}$.

Finding a symmetric \mathbf{B} is a lot more work. It involves finding the eigenvectors of \mathbf{A} .

The OLS Covariance Matrix

The covariance matrix of $\hat{\beta}$ depends on the covariance matrix of the disturbances (error terms).

If the disturbances are IID, the covariance matrix of \mathbf{u} is the **scalar matrix** $\sigma^2 \mathbf{I}$:

$$\text{Var}(\mathbf{u}) = \text{E}(\mathbf{u}\mathbf{u}^\top) = \sigma^2 \mathbf{I}. \quad (7)$$

Note that (7) does not require the disturbances to be independent, or even to have the same distribution.

But they must all have the same variance, and the covariance of each pair of disturbances must be zero.

When every u_i has the same variance, then they are **homoskedastic**.

When the $\text{Var}(u_i)$ differ, then they are **heteroskedastic**.

In general, we denote the $N \times N$ error covariance matrix by $\mathbf{\Omega}$.

For time-series data, when Ω has nonzero off-diagonal elements, the disturbances are said to be **autocorrelated** or **serially correlated**.

If the observations of a sample characterize different locations in space, they may display **spatial autocorrelation**.

Another possibility is that the disturbances are **clustered**, correlated within each of G clusters but uncorrelated across them.

When observations are sorted by cluster, Ω is block-diagonal, with G diagonal blocks that correspond to the G clusters:

$$\Omega = \begin{bmatrix} \Omega_1 & \mathbf{O} & \dots & \mathbf{O} \\ \mathbf{O} & \Omega_2 & \dots & \mathbf{O} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{O} & \mathbf{O} & \dots & \Omega_G \end{bmatrix}. \quad (8)$$

Here Ω_g is the $N_g \times N_g$ covariance matrix for the observations belonging to the g^{th} cluster.

Disturbances that are autocorrelated or clustered may or may not also be heteroskedastic. The usual approach to cluster-robust inference also allows for heteroskedasticity of unknown form.

When the DGP belongs to the model we estimate,

$$\hat{\beta} - \beta_0 = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{u}. \quad (9)$$

When $\hat{\beta}$ is unbiased, $\text{Var}(\hat{\beta})$ is the expectation of the $k \times k$ matrix

$$(\hat{\beta} - \beta_0)(\hat{\beta} - \beta_0)^\top = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{u} \mathbf{u}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1}. \quad (10)$$

If we take this expectation, conditional on \mathbf{X} , we find that

$$\text{Var}(\hat{\beta}) = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \text{E}(\mathbf{u} \mathbf{u}^\top) \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \quad (11)$$

$$= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{\Omega} \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1}. \quad (12)$$

This form of covariance matrix is called a **sandwich covariance matrix**, because $\mathbf{X}^\top \mathbf{\Omega} \mathbf{X}$ is sandwiched between the two instances of $(\mathbf{X}^\top \mathbf{X})^{-1}$.

The diagonal elements of $\text{Var}(\hat{\beta})$ are particularly interesting. The square root of the k^{th} diagonal element is the standard error of $\hat{\beta}_k$.

In practice, we almost never know Ω , so we have to figure out how to estimate the matrix $\mathbf{X}^\top \Omega \mathbf{X}$.

That is what methods for **heteroskedasticity-consistent** and **cluster-robust** covariance matrix estimation do.

If $\Omega = \sigma_0^2 \mathbf{I}$, so that there is neither heteroskedasticity nor autocorrelation, then equation (12) simplifies greatly. It becomes

$$\text{Var}(\hat{\beta}) = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \sigma_0^2 \mathbf{I} \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \quad (13)$$

$$= \sigma_0^2 (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} = \sigma_0^2 (\mathbf{X}^\top \mathbf{X})^{-1}. \quad (14)$$

For the next few lectures, we will assume that (14) holds.

But this simplification is often not valid, and most modern empirical work uses sandwich covariance matrices.

Precision of the Least-Squares Estimates

When (14) holds, the precision matrix is

$$\frac{1}{\sigma_0^2} \mathbf{X}^\top \mathbf{X}. \quad (15)$$

This is inversely proportional to σ_0^2 . The more random variation there is in the disturbances, the less precise are the parameter estimates.

It is illuminating to rewrite (15) as

$$\frac{N}{\sigma_0^2} \left(\frac{1}{N} \mathbf{X}^\top \mathbf{X} \right). \quad (16)$$

Under standard assumptions, $N^{-1} \mathbf{X}^\top \mathbf{X}$ is $O(1)$.

Thus the precision matrix (16) must be $O(N)$. When \mathbf{X} is stochastic, it is $O_p(N)$.

If we double the sample size, the precision of $\hat{\beta}$ should roughly double, and the standard errors of the individual $\hat{\beta}_i$ will be, approximately, divided by $\sqrt{2}$.

Consider the model $\mathbf{y} = \beta_1 \boldsymbol{\iota} + \mathbf{u}$. Replacing \mathbf{X} with $\boldsymbol{\iota}$, we find that

$$\hat{\beta}_1 = (\boldsymbol{\iota}^\top \boldsymbol{\iota})^{-1} \boldsymbol{\iota}^\top \mathbf{y} = \frac{1}{N} \sum_{i=1}^N y_i, \text{ and} \quad (17)$$

$$\text{Var}(\hat{\beta}_1) = \sigma_0^2 (\boldsymbol{\iota}^\top \boldsymbol{\iota})^{-1} = \frac{1}{N} \sigma_0^2. \quad (18)$$

The precision of the sample mean is exactly proportional to N , since the variance is proportional to $1/N$.

We have seen that σ_0^2 and N affect the precision of $\hat{\beta}$. The third thing that does so is the matrix \mathbf{X} .

High-leverage observations affect it much more than average ones.

We can rewrite a linear regression model as

$$\mathbf{y} = \mathbf{x}_1\beta_1 + \mathbf{X}_2\boldsymbol{\beta}_2 + \mathbf{u}, \quad (19)$$

where \mathbf{X} has been partitioned into \mathbf{x}_1 and \mathbf{X}_2 .

By the FWL Theorem, (19) yields the same estimate of β_1 as the FWL regression

$$\mathbf{M}_2\mathbf{y} = \mathbf{M}_2\mathbf{x}_1\beta_1 + \text{residuals}, \quad (20)$$

where $\mathbf{M}_2 \equiv \mathbf{I} - \mathbf{X}_2(\mathbf{X}_2^\top\mathbf{X}_2)^{-1}\mathbf{X}_2^\top$.

This estimate is $\hat{\beta}_1 = \mathbf{x}_1^\top\mathbf{M}_2\mathbf{y} / \mathbf{x}_1^\top\mathbf{M}_2\mathbf{x}_1$, and its variance is $\sigma_0^2(\mathbf{x}_1^\top\mathbf{M}_2\mathbf{x}_1)^{-1}$. Thus the precision of $\hat{\beta}_1$ is

$$\frac{1}{\sigma_0^2} \mathbf{x}_1^\top\mathbf{M}_2\mathbf{x}_1. \quad (21)$$

How much information the sample gives us about β_1 is proportional to the squared Euclidean length of the vector $\mathbf{M}_2\mathbf{x}_1$.

When $\|M_2x_1\|$ is big, because N is large or at least some elements of M_2x_1 are large, $\hat{\beta}_1$ is relatively precise.

When $\|M_2x_1\|$ is small, because N is small or all the elements of M_2x_1 are small, $\hat{\beta}_1$ is relatively imprecise.

The squared Euclidean length of M_2x_1 is just the SSR from

$$x_1 = X_2c + \text{residuals}. \quad (22)$$

Thus the precision of $\hat{\beta}_1$ is proportional to the SSR from regression (22). It depends on X_2 just as much as it depends on x_1 .

The quadratic form $x_1^\top M_2x_1$ can be written as $\sum_{i=1}^N (M_2x_1)_i^2$. Each term in the sum is a squared residual. For high-leverage observations, these squared residuals are large.

When X_2 explains x_1 much better than a constant alone, the length of M_2x_1 is much less than the length of M_1x_1 .

In this case, x_1 is said to be **collinear** with some of the other regressors.

This should be called **approximate collinearity**, but it is often (wrongly) called **multicollinearity**.

- Collinearity can greatly reduce the precision of OLS estimates.
- Estimates can be imprecise even when N is very large.

When the disturbances are not independent and identically distributed, things get more complicated.

In general, for a sandwich covariance matrix, the precision of $\hat{\beta}_1$ is

$$\frac{(\mathbf{x}_1^\top \mathbf{M}_2 \mathbf{x}_1)^2}{\mathbf{x}_1^\top \mathbf{M}_2 \boldsymbol{\Omega} \mathbf{M}_2 \mathbf{x}_1}. \quad (23)$$

It is usual to assume that $\mathbf{x}_1^\top \mathbf{M}_2 \mathbf{x}_1 = O(N)$. But, when there is dependence, it may be that $\mathbf{x}_1^\top \mathbf{M}_2 \boldsymbol{\Omega} \mathbf{M}_2 \mathbf{x}_1 = O(N^*)$, with $N^* > N$.

In that case, the precision of $\hat{\beta}_1$ may be $O(N^2/N^*) < O(N)$. Thus information accumulates less rapidly than it does for independent observations.

Linear Functions of Parameter Estimates

Suppose we are interested in the variance of $\hat{\gamma}$, where $\gamma = \mathbf{w}^\top \boldsymbol{\beta}$, $\hat{\gamma} = \mathbf{w}^\top \hat{\boldsymbol{\beta}}$, and \mathbf{w} is a k -vector of known coefficients.

By choosing \mathbf{w} appropriately, we can make γ equal to any one of the β_i , or to the sum of the β_i , or to any linear combination of the β_i .

For example, if $\gamma = 3\beta_1 - \beta_4$, \mathbf{w} would be a vector with 3 as the first element, -1 as the fourth element, and 0 for all the other elements.

The variance of $\hat{\gamma}$ is just

$$\text{Var}(\hat{\gamma}) = \text{Var}(\mathbf{w}^\top \hat{\boldsymbol{\beta}}) = \text{E}(\mathbf{w}^\top (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0) (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0)^\top \mathbf{w}) \quad (24)$$

$$= \mathbf{w}^\top \text{E}((\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0) (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0)^\top) \mathbf{w} \quad (25)$$

$$= \mathbf{w}^\top \text{Var}(\hat{\boldsymbol{\beta}}) \mathbf{w}. \quad (26)$$

In general, $\text{Var}(\hat{\gamma})$ depends on every element of $\text{Var}(\hat{\boldsymbol{\beta}})$.

The **delta method**, discussed in Chapter 5, generalizes the result (26) to nonlinear functions of β .

Consider the special case in which $\gamma = 3\beta_1 - \beta_4$. In this case,

$$\text{Var}(\hat{\gamma}) = w_1^2 \text{Var}(\hat{\beta}_1) + w_4^2 \text{Var}(\hat{\beta}_4) + 2w_1w_4 \text{Cov}(\hat{\beta}_1, \hat{\beta}_4) \quad (27)$$

$$= 9 \text{Var}(\hat{\beta}_1) + \text{Var}(\hat{\beta}_4) - 6 \text{Cov}(\hat{\beta}_1, \hat{\beta}_4). \quad (28)$$

The variance of $\hat{\gamma}$ depends on $\text{Cov}(\hat{\beta}_1, \hat{\beta}_4)$ as well as on the variances of $\hat{\beta}_1$ and $\hat{\beta}_4$.

When that covariance is large and positive, $\text{Var}(\hat{\gamma})$ may be small, even when $\text{Var}(\hat{\beta}_1)$ and $\text{Var}(\hat{\beta}_4)$ are both large.

Just looking at standard errors of the $\hat{\beta}_i$ can be extremely misleading when we care about functions of the β_i .

Instead of using (28), we may be able to rearrange the regression so that γ is estimated directly. Or use `lincom` in Stata.

The Variance of Forecast Errors

Suppose we have computed $\hat{\beta}$ and wish to predict y_j , for j not in $i = 1, \dots, N$, using observed regressors \mathbf{X}_j .

The forecast of y_j is $\mathbf{X}_j\hat{\beta}$. The **prediction error**, or **forecast error**, has mean zero, and variance

$$E(y_j - \mathbf{X}_j\hat{\beta})^2 = E(\mathbf{X}_j\beta_0 + u_j - \mathbf{X}_j\hat{\beta})^2 \quad (29)$$

$$= E(u_j^2) + E(\mathbf{X}_j\beta_0 - \mathbf{X}_j\hat{\beta})^2 \quad (30)$$

$$= \sigma_0^2 + \text{Var}(\mathbf{X}_j\hat{\beta}). \quad (31)$$

- The first equality depends on correct specification.
- The second depends on disturbances being serially uncorrelated, which ensures that $E(u_j\mathbf{X}_j\hat{\beta}) = 0$.
- The third uses the fact that $\hat{\beta}$ is assumed to be unbiased.

Under classical assumptions,

$$\text{Var}(y_j - \mathbf{X}_j \hat{\boldsymbol{\beta}}) = \sigma_0^2 + \sigma_0^2 \mathbf{X}_j (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}_j^\top. \quad (32)$$

The first term is the variance of u_j . The second term arises because we use $\hat{\boldsymbol{\beta}}$ instead of $\boldsymbol{\beta}_0$.

Thus the variance of the forecast error is greater than the variance of the disturbance u_j .

In contrast, as we will see shortly, the variance of the residual \hat{u}_i is less than the variance of the disturbance u_i .

Suppose the y_i were not generated by $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{u}$ but by the DGP

$$\mathbf{y} = \mathbf{Z}\boldsymbol{\gamma}_0 + \mathbf{u}, \quad (33)$$

where some columns of \mathbf{Z} may belong to $\mathcal{S}(\mathbf{X})$. Thus the model we estimate is misspecified.

If the regressors are fixed,

$$\hat{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Z} \gamma_0 + (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{u} \quad (34)$$

$$= \beta_0 + (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{u}, \quad (35)$$

where the **pseudo-true** parameter vector β_0 is $(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Z} \gamma_0$.

The expected squared forecast error based on the false model $\mathbf{y} = \mathbf{X}\beta + \mathbf{u}$ is

$$\begin{aligned} \mathbb{E}(y_j - \mathbf{X}_j \hat{\beta})^2 &= \mathbb{E}(\mathbf{Z}_j \gamma_0 + u_j - \mathbf{X}_j \hat{\beta})^2 \\ &= \mathbb{E}(u_j^2) + \mathbb{E}(\mathbf{Z}_j \gamma_0 - \mathbf{X}_j \hat{\beta})^2 \end{aligned} \quad (36)$$

$$= \mathbb{E}(u_j^2) + \mathbb{E}(\mathbf{Z}_j \gamma_0 - \mathbf{X}_j \beta_0 + \mathbf{X}_j \beta_0 - \mathbf{X}_j \hat{\beta})^2 \quad (37)$$

$$= \sigma_0^2 + \mathbb{E}(\mathbf{Z}_j \gamma_0 - \mathbf{X}_j \beta_0)^\top (\mathbf{Z}_j \gamma_0 - \mathbf{X}_j \beta_0) + \text{Var}(\mathbf{X}_j \hat{\beta}). \quad (38)$$

The middle term in the last line is essentially a squared bias.

Under classical assumptions, $\text{Var}(\mathbf{X}_j \hat{\boldsymbol{\beta}}) = \sigma_0^2 \mathbf{X}_j (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}_j^\top$.

Thus the expected squared forecast error for the false model is equal to the one for the true model, given in (32), plus a squared bias term.

If $\mathbf{X}_j \boldsymbol{\beta}_0$ provides a good approximation to $\mathbf{Z}_j \boldsymbol{\gamma}_0$, then the middle term in (38) will be small.

But if it provides a poor approximation, the middle term may be large.

In that case, the expected squared forecast error may be much larger than expression (32) suggests.

We conclude that forecast errors will generally be larger than residuals for three reasons:

- $\text{Var}(u_i) > \text{Var}(\hat{u}_i)$ for all i (to be proved).
- $\text{Var}(\mathbf{X}_j \hat{\boldsymbol{\beta}}) > 0$ because $\hat{\boldsymbol{\beta}}$ is random.
- $\text{E}(\mathbf{Z}_j \boldsymbol{\gamma}_0 - \mathbf{X}_j \boldsymbol{\beta}_0)^\top (\mathbf{Z}_j \boldsymbol{\gamma}_0 - \mathbf{X}_j \boldsymbol{\beta}_0) > 0$.