

Heteroskedasticity-Robust Inference

Consider the linear regression model with exogenous regressors,

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{u}, \quad \mathbb{E}(\mathbf{u}) = \mathbf{0}, \quad \mathbb{E}(\mathbf{u}\mathbf{u}^\top) = \boldsymbol{\Omega}, \quad (1)$$

where $\boldsymbol{\Omega}$ is an $N \times N$ matrix with i^{th} diagonal element equal to $\omega_i^2 > 0$ and all the off-diagonal elements equal to 0.

Since \mathbf{X} is assumed to be exogenous, the expectations in (1) can be treated as conditional on \mathbf{X} .

We would get the same asymptotic results if, instead of treating \mathbf{X} as exogenous, we assumed that $\mathbb{E}(u_i | \mathbf{X}_i) = 0$.

The disturbances in (1) are uncorrelated and have mean 0, but their variances differ. They are said to be **heteroskedastic**.

We assume that the investigator knows nothing about the ω_i^2 . In other words, the form of the heteroskedasticity is completely unknown.

Whatever the form of $\mathbf{\Omega}$, the covariance matrix of $\hat{\boldsymbol{\beta}}$ is

$$\begin{aligned} \mathbb{E}((\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0)(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0)^\top) &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbb{E}(\mathbf{u}\mathbf{u}^\top) \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \\ &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{\Omega} \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1}. \end{aligned} \quad (2)$$

This is often called a **sandwich covariance matrix**, for obvious reasons.

If we knew the ω_i^2 , we could evaluate (2). In fact, we could do better and obtain efficient estimates of $\boldsymbol{\beta}$.

Observations with low variance convey more information than ones with high variance, and so the former should be given greater weight.

But it is assumed that we do not know the ω_i^2 . We cannot hope to estimate them consistently without making additional assumptions, because there are N of them.

For the purposes of asymptotic theory, we wish to consider the covariance matrix of $N^{1/2}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0)$, the limit of N times the matrix (2).

The asymptotic covariance matrix of $N^{1/2}(\hat{\beta} - \beta_0)$ is

$$\left(\lim_{N \rightarrow \infty} \frac{1}{N} \mathbf{X}^\top \mathbf{X} \right)^{-1} \left(\lim_{N \rightarrow \infty} \frac{1}{N} \mathbf{X}^\top \boldsymbol{\Omega} \mathbf{X} \right) \left(\lim_{N \rightarrow \infty} \frac{1}{N} \mathbf{X}^\top \mathbf{X} \right)^{-1}. \quad (3)$$

Under standard assumptions, the factor $(\lim N^{-1} \mathbf{X}^\top \mathbf{X})^{-1}$ tends to the positive definite matrix $\mathbf{S}_{\mathbf{X}^\top \mathbf{X}}^{-1}$.

To estimate $\mathbf{S}_{\mathbf{X}^\top \mathbf{X}}^{-1}$, we can simply use the matrix $(N^{-1} \mathbf{X}^\top \mathbf{X})^{-1}$ itself.

In a very famous paper, White (1980) showed that, under certain conditions, the middle matrix can be estimated consistently by

$$\frac{1}{N} \mathbf{X}^\top \hat{\boldsymbol{\Omega}} \mathbf{X}, \quad (4)$$

where $\hat{\boldsymbol{\Omega}}$ is an *inconsistent* estimator of $\boldsymbol{\Omega}$.

The simplest version of $\hat{\boldsymbol{\Omega}}$ is a diagonal matrix with i^{th} diagonal element equal to \hat{u}_i^2 , the i^{th} squared OLS residual.

The matrix $\lim(N^{-1}\mathbf{X}^\top\mathbf{\Omega}\mathbf{X})$ is a $k \times k$ symmetric matrix. Therefore, it has exactly $(k^2 + k)/2$ distinct elements.

Since this number is independent of the sample size, the matrix can be estimated consistently. Its jl^{th} element is

$$\lim_{N \rightarrow \infty} \left(\frac{1}{N} \sum_{i=1}^N \omega_i^2 x_{ij} x_{il} \right). \quad (5)$$

This is estimated by the jl^{th} element of (4). For the simplest version of $\hat{\mathbf{\Omega}}$, the estimator is

$$\frac{1}{N} \sum_{i=1}^N \hat{u}_i^2 x_{ij} x_{il}. \quad (6)$$

Because $\hat{\beta}$ is consistent for β_0 , \hat{u}_i must be consistent for u_i , and \hat{u}_i^2 is therefore consistent for u_i^2 .

However, \hat{u}_i^2 does not estimate ω_i^2 consistently.

Asymptotically, expression (6) is equal to

$$\frac{1}{N} \sum_{i=1}^N u_i^2 x_{ij} x_{il} = \frac{1}{N} \sum_{i=1}^N (\omega_i^2 + v_i) x_{ij} x_{il} \quad (7)$$

$$= \frac{1}{N} \sum_{i=1}^N \omega_i^2 x_{ij} x_{il} + \frac{1}{N} \sum_{i=1}^N v_i x_{ij} x_{il}, \quad (8)$$

where v_i is defined to equal u_i^2 minus its mean of ω_i^2 .

Under suitable assumptions about the x_{ij} and the ω_i^2 , we can apply a law of large numbers to the second term in (8).

Since $E(v_i) = 0$, this term converges to 0, while the first term converges to expression (5).

Because $N^{-1} \sum_{i=1}^N \hat{u}_i^2 x_{ij} x_{il} \stackrel{a}{=} N^{-1} \sum_{i=1}^N u_i^2 x_{ij} x_{il}$, these arguments imply that (6) consistently estimates (5).

In practice, of course, we omit the factors of $1/N$. We simply use the matrix

$$\widehat{\text{Var}}_h(\hat{\beta}) \equiv (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \hat{\Omega} \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \quad (9)$$

directly to estimate the covariance matrix of $\hat{\beta}$.

A more revealing way to write (9) is

$$\widehat{\text{Var}}_h(\hat{\beta}) \equiv (\mathbf{X}^\top \mathbf{X})^{-1} \left(\sum_{i=1}^N \hat{u}_i^2 \mathbf{X}_i^\top \mathbf{X}_i \right) (\mathbf{X}^\top \mathbf{X})^{-1}, \quad (10)$$

where \mathbf{X}_i is the i^{th} row of \mathbf{X} . This makes it clear that the $N \times N$ matrix $\hat{\Omega}$ is never used.

The sandwich estimator (10) is a **heteroskedasticity-consistent covariance matrix estimator**, or **HCCME**. It is valid for heteroskedasticity of unknown form.

By taking square roots of the diagonal elements of (10), we can obtain **heteroskedasticity-robust standard errors**.

Asymptotic Theory for OLS

With heteroskedasticity of unknown form, Theorem 4.3 needs to be replaced by

$$N^{1/2}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0) \xrightarrow{d} N\left(\mathbf{0}, \mathbf{S}_{\mathbf{X}^\top \mathbf{X}}^{-1} \left(\lim_{N \rightarrow \infty} \frac{1}{N} \mathbf{X}^\top \boldsymbol{\Omega} \mathbf{X} \right) \mathbf{S}_{\mathbf{X}^\top \mathbf{X}}^{-1}\right) \quad (11)$$

and

$$\lim_{N \rightarrow \infty} \frac{1}{N} \widehat{\text{Var}}_h(\hat{\boldsymbol{\beta}}) = \mathbf{S}_{\mathbf{X}^\top \mathbf{X}}^{-1} \left(\lim_{N \rightarrow \infty} \frac{1}{N} \mathbf{X}^\top \boldsymbol{\Omega} \mathbf{X} \right) \mathbf{S}_{\mathbf{X}^\top \mathbf{X}}^{-1}. \quad (12)$$

We conclude that $\hat{\boldsymbol{\beta}}$ is root- N consistent and asymptotically normal, with (12) providing a consistent estimator of its covariance matrix.

Of course, all the factors of N are omitted when we actually make inferences about $\hat{\boldsymbol{\beta}}$.

Alternative Forms of HCCME

The original HCCME (10) of White (1980), often called HC_0 , uses squared residuals to estimate the diagonal elements Ω .

But least-squares residuals tend to be too small. Better estimators inflate the squared residuals (MacKinnon and White, 1985).

HC_1 : Use \hat{u}_i^2 in $\hat{\Omega}$ and then multiply the entire matrix by the scalar $N/(N - k)$, for a standard degrees-of-freedom correction.

HC_2 : Use $\hat{u}_i^2 / (1 - h_i)$ in $\hat{\Omega}$, where

$$h_i \equiv \mathbf{X}_i(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}_i^\top \quad (13)$$

is the i^{th} diagonal element of the “hat” matrix \mathbf{P}_X .

Recall the result that, when $\text{Var}(u_i) = \sigma^2$ for all i , $E(\hat{u}_i^2) = \sigma^2(1 - h_i)$. Therefore, the ratio of \hat{u}_i^2 to $1 - h_i$ would have expectation σ^2 if the disturbances were homoskedastic.

HC₃: Use $\hat{u}_i^2 / (1 - h_i)^2$ in $\hat{\Omega}$. This is a computationally efficient approximation to a **jackknife estimator**.

A jackknife estimator omits one observation at a time when obtaining estimates and fitted values.

The usual way to define a jackknife covariance matrix is

$$\widehat{\text{Var}}_{\text{JK}}(\hat{\beta}) = \frac{N-1}{N} \sum_{i=1}^N (\hat{\beta}^{(i)} - \hat{\beta})(\hat{\beta}^{(i)} - \hat{\beta})^\top, \quad (14)$$

where $\hat{\beta}^{(i)}$ is the vector of estimates when the i^{th} observation is omitted. We can also omit groups of observations.

It is not at all obvious that (14) is numerically equal to

$$\frac{N-1}{N} (\mathbf{X}^\top \mathbf{X})^{-1} \left(\sum_{i=1}^N \frac{\hat{u}_i^2}{(1-h_i)^2} \mathbf{X}_i^\top \mathbf{X}_i \right) (\mathbf{X}^\top \mathbf{X})^{-1}, \quad (15)$$

but it can be verified numerically.

MacKinnon and White (1985) actually proposed a slightly more complicated version of (15), based on a different version of (14),

$$\widehat{\text{Var}}_{\text{JK}}(\hat{\beta}) = \frac{N-1}{N} \sum_{i=1}^N (\hat{\beta}^{(i)} - \bar{\beta})(\hat{\beta}^{(i)} - \bar{\beta})^{\top}, \quad (16)$$

where $\bar{\beta} = N^{-1} \sum_{i=1}^N \hat{\beta}^{(i)}$, the mean of the omit-one estimates.

Dividing by $(1 - h_i)^2$ actually seems to overcorrect the residuals.

But observations with large variances often tend to have residuals that are very much too small. Thus, HC_3 may be attractive if large variances are associated with large values of h_i .

Inferences based on any HCCME, especially HC_0 and HC_1 , may be seriously inaccurate even when the sample size is moderately large if some observations have much higher leverage than others.

By default, Stata uses HC_1 with the “robust” and “vec(robust)” options. But “vce(hc2)” and “vce(hc3)” provide HC_2 and HC_3 .

When Does Heteroskedasticity Matter?

Even when the disturbances are heteroskedastic, we do not necessarily have to use an HCCME.

Consider the jl^{th} element of $N^{-1}\mathbf{X}^{\top}\mathbf{\Omega}\mathbf{X}$, which is

$$\frac{1}{N} \sum_{i=1}^N \omega_i^2 x_{ij} x_{il}. \quad (17)$$

If the limit as $N \rightarrow \infty$ of the average of the ω_i^2 exists and is denoted σ^2 , then expression (17) can be rewritten as

$$\sigma^2 \frac{1}{N} \sum_{i=1}^N x_{ij} x_{il} + \frac{1}{N} \sum_{i=1}^N (\omega_i^2 - \sigma^2) x_{ij} x_{il}. \quad (18)$$

The first term here is just the jl^{th} element of $\sigma^2 N^{-1}\mathbf{X}^{\top}\mathbf{X}$.

If it happens that

$$\lim_{N \rightarrow \infty} \frac{1}{N} \sum_{i=1}^N (\omega_i^2 - \sigma^2) x_{ij} x_{il} = 0 \quad (19)$$

for $j, l = 1, \dots, k$, then we find that

$$\lim_{N \rightarrow \infty} \frac{1}{N} \mathbf{X}^\top \boldsymbol{\Omega} \mathbf{X} = \sigma^2 \lim_{N \rightarrow \infty} \frac{1}{N} \mathbf{X}^\top \mathbf{X}. \quad (20)$$

If so, the asymptotic covariance matrix of $N^{1/2}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0)$ is just

$$\left(\lim_{N \rightarrow \infty} \frac{1}{N} \mathbf{X}^\top \mathbf{X} \right)^{-1} \sigma^2 \left(\lim_{N \rightarrow \infty} \frac{1}{N} \mathbf{X}^\top \mathbf{X} \right) \left(\lim_{N \rightarrow \infty} \frac{1}{N} \mathbf{X}^\top \mathbf{X} \right)^{-1} = \sigma^2 \mathbf{S}_{\mathbf{X}^\top \mathbf{X}}^{-1}. \quad (21)$$

The usual OLS estimate of σ^2 is $s^2 = (1/(N - k)) \sum_{i=1}^N \hat{u}_i^2$.

If we assume that we can apply a law of large numbers, the probability limit of $N^{-1} \sum_{i=1}^N \hat{u}_i^2$ is

$$\lim_{N \rightarrow \infty} \frac{1}{N} \sum_{i=1}^N \omega_i^2 = \sigma^2. \quad (22)$$

In this special case, the usual OLS covariance matrix estimator $s^2(\mathbf{X}^\top \mathbf{X})^{-1}$ is valid asymptotically.

If we are estimating a sample mean, then $\mathbf{X} = \mathbf{1}$, and

$$\frac{1}{N} \sum_{i=1}^N \omega_i^2 x_{ij} x_{il} = \frac{1}{N} \sum_{i=1}^N \omega_i^2 t_i^2 = \frac{1}{N} \sum_{i=1}^N \omega_i^2 \rightarrow \sigma^2 \text{ as } N \rightarrow \infty. \quad (23)$$

Thus (19) holds, and we do not have to worry about heteroskedasticity.

Only heteroskedasticity related to the squares and cross-products of the x_{ij} affects the validity of the usual OLS covariance matrix estimator.

HAC Covariance Matrix Estimation

The assumption that the matrix Ω is diagonal is what makes it possible to estimate $N^{-1}\mathbf{X}^\top\Omega\mathbf{X}$ consistently and obtain an HCCME, even though Ω itself cannot be estimated consistently.

The matrix $N^{-1}\mathbf{X}^\top\Omega\mathbf{X}$ can sometimes be estimated consistently for a model that uses time-series data when the disturbances are correlated across time periods.

Observations that are close to each other may be strongly correlated, but observations that are far apart may be uncorrelated or nearly so.

If so, only the elements of Ω that are on or close to the principal diagonal are large.

We may be able to obtain an estimate of the covariance matrix of the parameter estimates that is **heteroskedasticity and autocorrelation consistent**, or **HAC**.

Cluster-Robust Inference

Data are often collected at the individual level, but each observation is associated with a higher-level entity, such as a city, state, province, or country, a classroom or school, a hospital, or perhaps a time period.

Thus each observation belongs to a **cluster**, and the regression disturbances may be correlated within the clusters.

It seems natural to allow for any form of correlation within each of G clusters, while assuming no correlation across clusters.

The resulting covariance matrix is called a **cluster-robust variance estimator** or **CRVE**.

Instead of presenting an elementary exposition based on ETM2, I will present a more detailed one based on:

James G. MacKinnon, Morten Ø. Nielsen, and Matthew D. Webb, “Cluster-robust inference: A guide to empirical practice.” *Journal of Econometrics*, 2023, **232**, 272–299.