

Linear Models

Carlo Favero

Econometric Modelling of Financial Returns

- Financial data are mostly observational data: they are not generated by well-designed experiment to test hypothesis, they are given to the econometrician.
- These data can be used to construct non-causal predictive models and to evaluate treatment effects.
- The second exercise involves a deeper understanding of causation while the implementation of non-causal predictive modelling requires understanding conditional expectations.

Non experimental data

- What do we do with non-experimental data ?
- When you try and explain returns you do not control how the data are generated
- Think of estimating a CAPM model

Non experimental data: CAPM

$$\left(r_t^i - r_t^{rf} \right) = \beta_{0,i} + \beta_{1,i} \left(r_t^m - r_t^{rf} \right) + u_{i,t}$$

$$\left(r_t^m - r_t^{rf} \right) = \mu_m + u_{m,t}$$

$$u_{i,t} \sim n.i.d. \left(0, \sigma_i^2 \right)$$

$$\begin{pmatrix} u_{i,t} \\ u_{m,t} \end{pmatrix} \sim n.i.d. \left[\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} \sigma_{ii} & 0 \\ 0 & \sigma_{mm} \end{pmatrix} \right]$$

- you have a time-series on excess returns of a given stock on the safe asset and on the market on the safe asset (you cannot control these data)
- you specify a linear model in which the excess returns on asset i are function of the excess returns on the market, and the excess returns on the market follow a CER.
- you introduce hypothesis on the error terms, to estimate parameters, to implement tests, and to use the model for

Econometric Modelling of Financial Returns

- Econometric models of financial returns specify the distribution of a vector of variables \mathbf{y}_t conditional upon other variables \mathbf{z}_t that are helpful in predicting them.
- The mapping between \mathbf{y}_t and \mathbf{z}_t is determined by some functional relation and some unknown parameters.
- All the relevant variables are stochastic and they are therefore characterized by a density function.
- Linear Econometric Models specify conditional means of the \mathbf{y}_t as linear functions of the \mathbf{z}_t .

Econometric Modelling of Financial Returns

- the data

$$D(\mathbf{y}_t, \mathbf{z}_t, \mathbf{w}_t \mid \mathbf{Y}_{t-1}, \mathbf{Z}_{t-1}, \mathbf{W}_{t-1}, \boldsymbol{\theta})$$

- a general multivariate model

$$D(\mathbf{y}_t, \mathbf{z}_t \mid \mathbf{Y}_{t-1}, \mathbf{Z}_{t-1}, \boldsymbol{\beta})$$

- decomposing a multivariate into conditional and marginal

$$D(\mathbf{y}_t \mid \mathbf{z}_t, \mathbf{Y}_{t-1}, \mathbf{Z}_{t-1}, \boldsymbol{\beta}_1) D(\mathbf{z}_t \mid \mathbf{Y}_{t-1}, \mathbf{Z}_{t-1}, \boldsymbol{\beta}_2)$$

- a general linear univariate conditional model

$$\begin{aligned} y_t &= \boldsymbol{\beta}'_1 \mathbf{z}_t + u_{1t} \\ \mathbf{z}_t &= \boldsymbol{\beta}_2 \mathbf{z}_{t-1} + u_{2t} \end{aligned}$$

- There are many ways in which the CAPM can go wrong:
 - other factors beyond the market are relevant in determining excess returns on asset i
 - the excess returns on the market do depend on excess returns on asset i
 - the model is non-linear
 - the residuals are non-normal and their variance-covariance matrix does not fit the assumptions made

- Conditional densities are best interpreted as the outcome of a reduction process that allows a simplified representation of reality.
- Of course such a simplified representation omits an enormous amount of information.
- The validity of the model adopted is crucially affected by the importance of the omitted information in determining the density of y_t .

Reduction Process

- To understand the reduction process partition the set of all variables into three types of variables:

$$\mathbf{x}_t = (\mathbf{w}_t, \mathbf{y}_t, \mathbf{z}_t),$$

\mathbf{w}_t identifies variables which are ignored in the specification of the econometric model.

- Exclusion is obtained by factorizing the joint density and integrating it with respect to \mathbf{w}_t . In formal terms, we have no information loss only if

$$D(\mathbf{y}_t, \mathbf{z}_t \mid \mathbf{Y}_{t-1}, \mathbf{Z}_{t-1}, \boldsymbol{\beta}) = D(\mathbf{y}_t, \mathbf{z}_t, \mathbf{w}_t \mid \mathbf{Y}_{t-1}, \mathbf{Z}_{t-1}, \mathbf{W}_{t-1}, \boldsymbol{\theta}).$$

This is the statistical model considered by the econometrician, this is technically called i.e. the reduced form of the structure of interest. In general this reduced form is a more general model than the one estimated. It is constructed by parameterizing $E(\mathbf{y}_t, \mathbf{z}_t \mid \mathbf{Y}_{t-1}, \mathbf{Z}_{t-1}, \boldsymbol{\beta})$ and by deriving a vector of innovations from the difference between the vector of observed variables and the vector of their means.

The CAPM Reduced Form

In the case of the CAPM the general specification of the reduced form is the following one:

$$\begin{aligned} \begin{pmatrix} r_t^i - r_t^{rf} \end{pmatrix} &= \mu_i + \beta_i u_{m,t} + u_{i,t} \\ \begin{pmatrix} r_t^m - r_t^{rf} \end{pmatrix} &= \mu_m + u_{m,t} \\ \begin{pmatrix} u_{i,t} \\ u_{m,t} \end{pmatrix} &\sim n.i.d. \left[\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} \sigma_{ii} & \sigma_{im} \\ \sigma_{im} & \sigma_{mm} \end{pmatrix} \right] \end{aligned}$$

From the Statistical Model to the Conditional Model: the CAPM

Statistical model

$$\begin{aligned}\left(r_t^i - r_t^{rf}\right) &= \mu_i + \beta_i u_{m,t} + u_{i,t} \\ \left(r_t^m - r_t^{rf}\right) &= \mu_m + u_{m,t} \\ \begin{pmatrix} u_{i,t} \\ u_{m,t} \end{pmatrix} &\sim n.i.d. \left[\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} \sigma_{ii} & \sigma_{im} \\ \sigma_{im} & \sigma_{mm} \end{pmatrix} \right]\end{aligned}$$

Estimated Equation

$$E \left(\left(r_t^i - r_t^{rf} \right) \mid \left(r_t^m - r_t^{rf} \right), \mathbf{Y}_{t-1}, \mathbf{Z}_{t-1}, \beta_i \right) = \alpha_i + \beta_i \left(r_t^m - r_t^{rf} \right)$$

if $\sigma_{im} = 0$, then the estimated equation is a valid approximation to the statistical model for the estimation of β_i

To illustrate how estimation can be performed to derive conditional expectations, consider the following general representation of the model of interest:

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon},$$

$$\mathbf{y} = \begin{pmatrix} y_1 \\ \cdot \\ \cdot \\ \cdot \\ y_N \end{pmatrix}, \quad \mathbf{X} = \begin{pmatrix} x_{11} & x_{12} & \cdot & \cdot & x_{1k} \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ x_{N1} & x_{N2} & \cdot & \cdot & x_{Nk} \end{pmatrix},$$

$$\boldsymbol{\beta} = \begin{pmatrix} \beta_1 \\ \cdot \\ \cdot \\ \cdot \\ \beta_k \end{pmatrix}, \quad \boldsymbol{\epsilon} = \begin{pmatrix} \epsilon_1 \\ \cdot \\ \cdot \\ \cdot \\ \epsilon_N \end{pmatrix}.$$

- The simplest way to derive estimates of the parameters of interest is the ordinary least squares (OLS) method.
- Such a method chooses values for the unknown parameters to minimize the magnitude of the non-observable components.
- In our simple bivariate case this amounts to choosing a line that goes through the scatterplot of excess returns on each asset on the market excess returns such that it provides the best fit.
- The best fit is obtained by minimizing the sum of squared vertical deviations of the data points from the fitted line.

Define the following quantity:

$$\mathbf{e}(\boldsymbol{\beta}) = \mathbf{y} - \mathbf{X}\boldsymbol{\beta},$$

where $\mathbf{e}(\boldsymbol{\beta})$ is a $(n \times 1)$ vector. If we treat $\mathbf{X}\boldsymbol{\beta}$, as a (conditional) prediction for \mathbf{y} , then we can consider $\mathbf{e}(\boldsymbol{\beta})$ as a forecasting error. The sum of the squared errors is then

$$\mathbf{S}(\boldsymbol{\beta}) = \mathbf{e}(\boldsymbol{\beta})' \mathbf{e}(\boldsymbol{\beta}).$$

The OLS method produces an estimator of β , $\hat{\beta}$, defined as follows:

$$\mathbf{S}(\hat{\beta}) = \min_{\beta} \mathbf{e}(\beta)' \mathbf{e}(\beta).$$

Given $\hat{\beta}$, we can define an associated vector of residual $\hat{\mathbf{e}}$ as $\hat{\mathbf{e}} = \mathbf{y} - \mathbf{X}\hat{\beta}$. The OLS estimator is derived by considering the necessary and sufficient conditions for $\hat{\beta}$ to be a unique minimum for \mathbf{S} :

- 1 $\mathbf{X}'\hat{\mathbf{e}} = 0$;
- 2 $\text{rank}(\mathbf{X}) = k$.

Condition 1 imposes orthogonality between the right-hand side variables on the OLS residuals, and ensures that residuals have an average of zero when a constant is included among the regressors. Condition 2 requires that the columns of the \mathbf{X} matrix are linearly independent: no variable in \mathbf{X} can be expressed as a linear combination of the other variables in \mathbf{X} .

From 1 we derive an expression for the OLS estimates:

$$\begin{aligned}\mathbf{X}'\hat{\boldsymbol{\varepsilon}} &= \mathbf{X}'(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}) = \mathbf{X}'\mathbf{y} - \mathbf{X}'\mathbf{X}\hat{\boldsymbol{\beta}} = 0, \\ \hat{\boldsymbol{\beta}} &= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}.\end{aligned}$$

Properties of the OLS estimates

We have derived the OLS estimator without any assumption on the statistical structure of the data. However, the statistical structure of the data is needed to define the properties of the estimator. To illustrate them, we refer to the basic concepts of mean and variance of vector variables.

Given a generic vector of variables, \mathbf{x} , and its mean vector $E(\mathbf{x})$

$$\mathbf{x} = \begin{pmatrix} x_1 \\ \cdot \\ \cdot \\ \cdot \\ x_n \end{pmatrix}, E(\mathbf{x}) = \begin{pmatrix} E(x_1) \\ \cdot \\ \cdot \\ \cdot \\ E(x_n) \end{pmatrix}$$

Properties of the OLS estimates

We define the mean matrix of outer products $E(\mathbf{xx}')$ as:

$$E(\mathbf{xx}') = \begin{pmatrix} E(x_1^2) & E(x_1x_2) & \cdot & \cdot & E(x_1x_n) \\ \cdot & E(x_2^2) & \cdot & \cdot & E(x_2x_n) \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ E(x_nx_1) & E(x_nx_2) & \cdot & \cdot & E(x_n^2) \end{pmatrix}$$

The variance-covariance matrix of \mathbf{x} is defined as:

$$\begin{aligned} \text{var}(\mathbf{x}) &= E(\mathbf{x} - E(\mathbf{x}))E(\mathbf{x} - E(\mathbf{x}))' \\ &= E(\mathbf{xx}') - E(\mathbf{x})E(\mathbf{x})'. \end{aligned}$$

The variance-covariance matrix is symmetric and positive definite, by construction. Given an arbitrary \mathbf{A} vector of dimension n , we have:

$$\text{var}(\mathbf{A}'\mathbf{x}) = \mathbf{A}'\text{var}(\mathbf{x})\mathbf{A}.$$

The first hypothesis

The first relevant hypothesis for the derivation of the statistical properties of OLS regards the relationship between disturbances and regressors in the estimated equation. This hypothesis is constructed in two parts: first we assume that $E(\mathbf{y}_i | \mathbf{x}_i) = \mathbf{x}_i' \boldsymbol{\beta}$, ruling out the contemporaneous correlation between residuals and regressors (note that assuming the validity of this hypothesis implies that there are no omitted variables correlated with the regressors), second we assume that the components of the available sample are independently drawn. The second assumption guarantees the equivalence between $E(\mathbf{y}_i | \mathbf{x}_i) = \mathbf{x}_i' \boldsymbol{\beta}$ and $E(\mathbf{y}_i | \mathbf{x}_1, \dots, \mathbf{x}_i, \dots, \mathbf{x}_n) = \mathbf{x}_i' \boldsymbol{\beta}$. Using vector notation, we have:

$$E(\mathbf{y} | \mathbf{X}) = \mathbf{X}\boldsymbol{\beta},$$

which is equivalent to

$$E(\boldsymbol{\epsilon} | \mathbf{X}) = \mathbf{0}. \quad (1)$$

The first hypothesis

Note that hypothesis (1) is very demanding. It implies that

$$E(\epsilon_i \mid \mathbf{x}_1, \dots, \mathbf{x}_i, \dots, \mathbf{x}_n) = 0 \quad (i = 1, \dots, n).$$

The conditional mean is, in general, a non-linear function of $(\mathbf{x}_1, \dots, \mathbf{x}_i, \dots, \mathbf{x}_n)$ and (1) requires that such a function is a constant of zero. Note that (1) requires that each regressor is orthogonal not only to the error term associated with the same observation ($E(x_{ik}\epsilon_i) = 0$ for all k), but also to the error term associated with each other observation ($E(x_{jk}\epsilon_i) = 0$ for all $j \neq k$). This statement is proved by using the properties of conditional expectations.

The first hypothesis

Since $E(\boldsymbol{\epsilon} \mid \mathbf{X}) = \mathbf{0}$ implies, from the law of iterated expectations, that $E(\boldsymbol{\epsilon}) = \mathbf{0}$, we have

$$E(\epsilon_i \mid x_{jk}) = E[E(\epsilon_i \mid \mathbf{x}) \mid x_{jk}] = 0. \quad (2)$$

Then

$$\begin{aligned} E(\epsilon_i x_{jk}) &= E[E(\epsilon_i x_{jk} \mid x_{jk})] \\ &= E[x_{jk} E(\epsilon_i \mid x_{jk})] \\ &= 0. \end{aligned}$$

The second and third hypothesis

The second hypothesis defines the constancy of the conditional variance of shocks:

$$E(\boldsymbol{\epsilon}'\boldsymbol{\epsilon} \mid \mathbf{X}) = \sigma^2 I, \quad (3)$$

where σ^2 is a constant independent from \mathbf{X} . In the case of our data, this is a strong assumption unlikely to be met in practice.

The third hypothesis is the one already introduced, which guarantees that the OLS estimator can be derived:

$$\text{rank}(\mathbf{X}) = k. \quad (4)$$

Under hypotheses (1) – (4) we can derive the properties of the OLS estimator.

Property 1: unbiasedness

The conditional expectation (with respect to \mathbf{X}) of the OLS estimates is the vector of unknown parameters $\boldsymbol{\beta}$:

$$\begin{aligned}\hat{\boldsymbol{\beta}} &= (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}' (\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}) \\ &= \boldsymbol{\beta} + (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\boldsymbol{\epsilon} \\ E(\hat{\boldsymbol{\beta}} | \mathbf{X}) &= \boldsymbol{\beta} + (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'E(\boldsymbol{\epsilon} | \mathbf{X}) \\ &= \boldsymbol{\beta},\end{aligned}$$

by hypothesis (1).

Property 2: variance of OLS

The conditional variance of the OLS estimator is $\sigma^2 (\mathbf{X}'\mathbf{X})^{-1}$:

$$\begin{aligned} \text{var}(\hat{\boldsymbol{\beta}} | \mathbf{X}) &= E\left(\left(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}\right)\left(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}\right)' | \mathbf{X}\right) \\ &= E\left(\left(\mathbf{X}'\mathbf{X}\right)^{-1}\mathbf{X}'\boldsymbol{\epsilon}\boldsymbol{\epsilon}'\mathbf{X}\left(\mathbf{X}'\mathbf{X}\right)^{-1} | \mathbf{X}\right) \\ &= \left(\mathbf{X}'\mathbf{X}\right)^{-1}\mathbf{X}'E\left(\boldsymbol{\epsilon}\boldsymbol{\epsilon}' | \mathbf{X}\right)\mathbf{X}\left(\mathbf{X}'\mathbf{X}\right)^{-1} \\ &= \left(\mathbf{X}'\mathbf{X}\right)^{-1}\mathbf{X}'\sigma^2\mathbf{I}\mathbf{X}\left(\mathbf{X}'\mathbf{X}\right)^{-1} \\ &= \sigma^2\left(\mathbf{X}'\mathbf{X}\right)^{-1}. \end{aligned}$$

Property 3: Gauss-Markov theorem

The OLS estimator is the most efficient in the class of linear unbiased estimators. Consider the class of linear estimators:

$$\beta_L = \mathbf{L}y.$$

This class is defined by the set of matrices ($k \times n$) \mathbf{L} , which are fixed when conditioning upon \mathbf{X} . \mathbf{L} does not depend on y . Therefore we have:

$$E(\beta_L | \mathbf{X}) = E(\mathbf{L}\mathbf{X}\beta + \mathbf{L}\varepsilon | \mathbf{X}) = \mathbf{L}\mathbf{X}\beta,$$

and $\mathbf{L}\mathbf{X}\beta = \beta$ only if $\mathbf{L}\mathbf{X} = \mathbf{I}_k$. Such a condition is obviously satisfied by the OLS estimator, which is obtained by setting $\mathbf{L} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$. The variance of the general estimator in the class of linear unbiased estimators is readily obtained as:

$$\text{var}(\beta_L | \mathbf{X}) = E(\mathbf{L}\varepsilon\varepsilon'\mathbf{L}' | \mathbf{X}) = \sigma^2\mathbf{L}\mathbf{L}'.$$

To show that the OLS estimator is the most efficient within this class we have to show that the variance of the OLS estimator differs from the variance of the generic estimator in the class by a positive semidefinite matrix.

To this aim define $\mathbf{D} = \mathbf{L} - (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'$; $\mathbf{LX} = \mathbf{I}$ requires $\mathbf{DX} = \mathbf{0}$.

$$\begin{aligned} \mathbf{LL}' &= \left((\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}' + \mathbf{D} \right) \left(\mathbf{X} (\mathbf{X}'\mathbf{X})^{-1} + \mathbf{D}' \right) \\ &= (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{X} (\mathbf{X}'\mathbf{X})^{-1} + (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{D}' + \\ &\quad + \mathbf{DX} (\mathbf{X}'\mathbf{X})^{-1} + \mathbf{DD}' \\ &= (\mathbf{X}'\mathbf{X})^{-1} + \mathbf{DD}', \end{aligned}$$

from which we have that

$$\text{var}(\boldsymbol{\beta}_L | \mathbf{X}) = \text{var}(\widehat{\boldsymbol{\beta}} | \mathbf{X}) + \sigma^2 \mathbf{DD}',$$

which proves the point. For any given matrix \mathbf{D} , (not necessarily square), the symmetric matrix \mathbf{DD}' is positive semidefinite.

Residual Analysis

Consider the following representation:

$$\begin{aligned}\hat{\boldsymbol{\epsilon}} &= \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}} \\ &= \mathbf{y} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} = \mathbf{M}\mathbf{y},\end{aligned}$$

where $\mathbf{M} = \mathbf{I}_n - \mathbf{Q}$, and $\mathbf{Q} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$. The $(n \times n)$ matrices \mathbf{M} and \mathbf{Q} , have the following properties:

- 1 they are symmetric: $\mathbf{M}' = \mathbf{M}, \mathbf{Q}' = \mathbf{Q}$;
- 2 they are idempotent: $\mathbf{Q}\mathbf{Q} = \mathbf{Q}, \mathbf{M}\mathbf{M} = \mathbf{M}$;
- 3 $\mathbf{M}\mathbf{X} = \mathbf{0}, \mathbf{M}\mathbf{Q} = \mathbf{0}, \mathbf{Q}\mathbf{X} = \mathbf{X}$.

Residual Analysis

Note that the OLS projection for \mathbf{y} can be written as $\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{Q}\mathbf{y}$, and that $\hat{\boldsymbol{\epsilon}} = \mathbf{M}\mathbf{y}$, from which we have the known result of orthogonality between the OLS residuals and regressors. We also have $\mathbf{M}\mathbf{y} = \mathbf{M}\mathbf{X}\boldsymbol{\beta} + \mathbf{M}\boldsymbol{\epsilon} = \mathbf{M}\boldsymbol{\epsilon}$, given that $\mathbf{M}\mathbf{X} = \mathbf{0}$. Therefore we have a very well-specified relation between the OLS residuals and the errors in the model $\hat{\boldsymbol{\epsilon}} = \mathbf{M}\boldsymbol{\epsilon}$, which cannot be used to derive the errors given the residuals, since the \mathbf{M} matrix is not invertible.

We can re-write the sum of squared residuals as:

$$S(\hat{\boldsymbol{\beta}}) = \hat{\boldsymbol{\epsilon}}'\hat{\boldsymbol{\epsilon}} = \boldsymbol{\epsilon}'\mathbf{M}'\mathbf{M}\boldsymbol{\epsilon} = \boldsymbol{\epsilon}'\mathbf{M}\boldsymbol{\epsilon}.$$

$S(\hat{\boldsymbol{\beta}})$ is an obvious candidate for the construction of an estimate for σ^2 .

To derive an estimate of σ^2 from $S(\hat{\beta})$, we introduce the concept of trace. The trace of a square matrix is the sum of all elements on its principal diagonal. The following properties are relevant:

- 1 given any two square matrices \mathbf{A} and \mathbf{B} , $tr(\mathbf{A} + \mathbf{B}) = tr\mathbf{A} + tr\mathbf{B}$;
- 2 given any two matrices \mathbf{A} and \mathbf{B} , $tr(\mathbf{AB}) = tr(\mathbf{BA})$;
- 3 the rank of an idempotent matrix is equal to its trace.

Residual Analysis

Using property 2 together with the fact that a scalar coincides with its trace, we have:

$$\boldsymbol{\epsilon}'\mathbf{M}\boldsymbol{\epsilon} = \text{tr}(\boldsymbol{\epsilon}'\mathbf{M}\boldsymbol{\epsilon}) = \text{tr}(\mathbf{M}\boldsymbol{\epsilon}\boldsymbol{\epsilon}').$$

Now we analyse the expected value of $S(\hat{\boldsymbol{\beta}})$, conditional upon \mathbf{X} :

$$\begin{aligned} E\left(S(\hat{\boldsymbol{\beta}}) \mid \mathbf{X}\right) &= E(\text{tr}\mathbf{M}\boldsymbol{\epsilon}\boldsymbol{\epsilon}' \mid \mathbf{X}) \\ &= \text{tr}E(\mathbf{M}\boldsymbol{\epsilon}\boldsymbol{\epsilon}' \mid \mathbf{X}) \\ &= \text{tr}\mathbf{M}(E\boldsymbol{\epsilon}\boldsymbol{\epsilon}' \mid \mathbf{X}) \\ &= \sigma^2 \text{tr}\mathbf{M}. \end{aligned}$$

From properties 1 and 2 we have:

$$\begin{aligned} \text{tr}\mathbf{M} &= \text{tr}\mathbf{I}_n - \text{tr}\left(\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\right) \\ &= n - \text{tr}\left(\mathbf{X}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\right) \\ &= n - k. \end{aligned}$$

Therefore, an unbiased estimate of σ^2 is given by $s^2 = S(\hat{\boldsymbol{\beta}}) / (n - k)$.

Residual Analysis: the R-squared

Using the result of orthogonality between the OLS projections and residuals, we can write:

$$\text{var}(\mathbf{y}) = \text{var}(\hat{\mathbf{y}}) + \text{var}(\hat{\boldsymbol{\epsilon}}),$$

from which we can derive the following residual-based indicator of the goodness of fit:

$$R^2 = \frac{\text{var}(\hat{\mathbf{y}})}{\text{var}(\mathbf{y})} = 1 - \frac{\text{var}(\hat{\boldsymbol{\epsilon}})}{\text{var}(\mathbf{y})}.$$

The information contained in R^2 is associated with the information contained in the standard error of the regression, which is the square root of the estimated variance of OLS residuals.