

Assumptions of the **classical linear regression model**,
what to do when they are violated,
and **estimator properties**

Contents

1	Assumptions of the CLRM for inference	2
2	Estimators of $\{\beta, \sigma^2\}$ and their statistical properties	3
2.1	Estimator properties	3
2.2	ML estimator	4
2.3	OLS estimator	5
3	Post-estimation model diagnostics	6
4	How to deal with non-spherical errors	9
4.1	OLS: Sandwich estimators	9
	References	14
	Appendix A Misc.	15
A.1	Deriving the formula of the OLS estimator	15
A.2	Linear algebra — Positive-definite matrices	16
A.3	Kernel functions for non-parametric statistics	16

Disclaimer: Sections and lines in brown correspond to content which is very much ‘under construction’.

1 Assumptions of the CLRM for inference

The classical linear regression model (CLRM) consists of a set of assumptions that describe the data generating process (DGP). By decreasing order of importance (Gelman et al., 2020, Ch. 11):

Notation:	System of n equations	Matrix
Model:	$y_i = \mathbf{X}_i' \beta + e_i \quad (i = 1, \dots, n)$	$y = \mathbf{X} \beta + \mathbf{e}$
Assumptions		
A1. linearity	The model is linear in β	The model is linear in β
A2. identification	$\rho_{X_k, X_l} \approx 1$	$\mathbf{X}_{N \times K}$ has rank K
A3. exogeneity	$\mathbb{E}[e_i \mathbf{X}_i] = 0$	$\mathbb{E}[\mathbf{e} \mathbf{X}] = \mathbf{0}_{N \times 1}$
A4. spherical errors	$e_i \mathbf{X}_i \stackrel{\text{iid}}{\sim} (0, \sigma^2)$	$\mathbb{V}[\mathbf{e} \mathbf{X}] = \sigma^2 \mathbf{I}_N$
– independent errors	$\text{cov}[e_i, e_j \mathbf{X}] = 0$	
– homoskedastic errors	$\mathbb{V}[e_i \mathbf{X}_i] = \sigma^2 \quad \sigma_i^2$	
A5. normal errors	$e_i \mathbf{X}_i \sim \mathcal{N}(0, \sigma^2)$	$\mathbf{e} \mathbf{X} \sim \mathcal{N}(\mathbf{0}_{N \times 1}, \sigma^2 \mathbf{I}_N)$

(A1) **Linearity in the parameters** and correct model specification (notably an additive error term).
I.e., the linear functional form coincides with the actual DGP.

(A2) **Identification:** regressors are linearly independent (no perfect collinearity).
If this is violated, drop one regressor, or transform collinear regressors into a single X .

(A3) **Strict exogeneity of regressors:** all other factors that affect y are unrelated to X .
 $\mathbb{E}[\mathbf{e} | \mathbf{X}] = \mathbf{0}$ also implies $\mathbb{E}[\mathbf{e}] = \mathbf{0}$ and $\mathbb{E}[\mathbf{X}' \mathbf{e}] = \mathbf{0}$, leading to $\text{cov}[\mathbf{e}, \mathbf{X}] = \mathbf{0}$: X and \mathbf{e} are uncorrelated.

(A4) **Spherical errors**¹

- **Independent errors** \implies no autocorrelation: $\text{cov}[e_i, e_j | \mathbf{X}] = \mathbb{E}[e_i e_j | \mathbf{X}] = 0$

I.e., errors are randomly spread around the regression line.

If this is violated, e.g., by serial correlation (likely with time series data), try taking lags of regressors, or switch to an autoregressive or a moving average model...

- **Homoskedastic errors:** equal conditional variance $\mathbb{V}[e_i | \mathbf{X}] = \sigma^2$

The error variance is a measure of model uncertainty. Homoskedasticity means uncertainty, i.e., the spread of errors, is identical across the support of y .

If this is violated, $\hat{\beta}_{\text{OLS}}$ remains valid but isn't efficient — Weighted Least Squares has a lower variance. Look for omitted variables, remove outliers, perform a log-transformation...

(A5) **Normal errors**

This assumption is not required for estimating the regression but for making inferences, e.g., computing confidence intervals or p-values. Without (A5), t and F tests are invalid.

If this is violated, we have to appeal to asymptotics: the properties of $\hat{\beta}$ for large samples. Indeed: the one-sample t -test for β , which tests the null hypothesis that $\beta=0$, assumes that the sampling distribution of $\hat{\beta}$ is normal. If errors are not normal, then $\hat{\beta}$ isn't normal. However, when n is large enough, Laws of Large Numbers (LLNs) and Central Limit Theorems (CLTs) say that the asymptotic sampling distribution of $\hat{\beta}$ is normal, s.t. t and F tests are robust to departures from normality if n is large. If errors are highly non-normal (e.g., long tailed), appealing to an asymptotically normal approximation may be unreasonable, and one may want to consider an alternative (e.g., bootstrap).

¹Note that most statements are actually conditional statements. E.g., (A4) assumes *conditionally* homoskedastic errors.

2 Estimators of $\{\beta, \sigma^2\}$ and their statistical properties

Multiple estimators are often available to estimate some summary of the relationship between X and y . How one chooses between them (besides their ease of computation) is motivated by their **statistical properties**.

2.1 Estimator properties

Let $\hat{\theta}$ be an estimator for the population parameter θ , for a sample of size n . $\hat{\theta}$, as a function of the random sample, is a random variable. The various possible samples of size n would each lead to a different realization $\hat{\theta}_s$, which together make the estimator's distribution or pdf $f_{\hat{\theta}}$. $\hat{\theta}$ has:

- finite sample properties: characteristics of $f_{\hat{\theta}}$ for a finite n . *Ex: bias, efficiency (precision);*
- asymptotic properties: characteristics of $f_{\hat{\theta}}$ as $n \rightarrow \infty$. *Ex: consistency, asymptotic distribution.*

As we always deal with finite samples, finite sample properties may seem the most important. In effect, bias (concerned with the center of the pdf) and efficiency (its spread) are the most common selection criteria. But remember that they tell us nothing about the properties of the estimator *for our own sample*, they tell us only about the distribution of values from hypothetical samples.

Finite sample properties

- $\hat{\theta}$ is **unbiased** iff $\mathbb{E}[\hat{\theta}] = \theta$

The estimator is correct *in expectation* over all possible samples. I.e., its distribution is centered around the estimand. \triangle *But our estimate from our own sample could be anywhere within that distribution, e.g., far from its center.* The bias of $\hat{\theta}$ is $\mathbb{E}[\hat{\theta}] - \theta$.

- $\hat{\theta}$ is **efficient** or “best” iff it has the lowest possible variance of all estimators: $\mathbb{V}[\hat{\theta}] \leq \mathbb{V}[\tilde{\theta} \dots]$

Its distribution is condensed, thus though the realization $\hat{\theta}_s$ from any sample could be anywhere within that distribution, it will never be too far away from the mean (which, if the estimator is unbiased, is the true θ). For unbiased estimators, that variance is the Cramér-Rao lower bound.

Asymptotic properties

- $\hat{\theta}$ is **asymptotically unbiased** iff $\mathbb{E}[\hat{\theta}] \xrightarrow[n \rightarrow +\infty]{p} \theta$

- $\hat{\theta}$ is **asymptotically efficient** iff $\mathbb{V}[\hat{\theta}] \xrightarrow[n \rightarrow +\infty]{p} \text{asymptotic Cramér-Rao lower bound}$

- $\hat{\theta}$ is **consistent** iff $\hat{\theta} \xrightarrow[n \rightarrow +\infty]{p} \theta$

I.e., as we get enough data, then we know the truth. Sufficient conditions are that $\hat{\theta}$ be asymptotically unbiased, and its variance $\rightarrow 0$ as $n \rightarrow \infty$.

In frequentist statistics, the Maximum Likelihood (ML) and the Ordinary Least Squares (OLS) estimators are widely used. The next sections describe them and their properties.² A few preliminary remarks:

- OLS and ML are rooted in different mathematical disciplines: calculus for OLS, probabilities for ML. OLS makes no assumption on the probabilistic nature of the variables, it is deterministic.
- ML includes OLS as a special case: if $e_i|X_i \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \sigma^2)$, then $y_i|X_i \sim \text{MVN}(X_i'\beta, \sigma^2)$ and $\hat{\beta}_{\text{OLS}} = \hat{\beta}_{\text{ML}}$.

²For convenience, we will drop the notation $|X$, however all features of the distributions of the estimators $\hat{\beta}_{\text{ML}}$ and $\hat{\beta}_{\text{OLS}}$ presented are actually conditional on X .

2.2 ML estimator $\hat{\theta}_{\text{ML}} = \{\hat{\beta}_{\text{ML}}, \hat{\sigma}_{\text{ML}}^2\}$

Definition The likelihood function in a regression model is the probability density of the data given the parameters θ and predictors. The Maximum Likelihood (ML) estimator is then the value of the parameters θ s.t. under the assumed model, the observed data are most likely. Assuming iid observations, we have:

- the likelihood $\mathcal{L}(y|X, \theta) = f(X_1, \dots, X_n, \theta) = f(X_1, \theta) \times \dots \times f(X_n, \theta) = \prod_{i=1}^n f(X_i, \theta)$
- the log-likelihood $\log \mathcal{L}(y|X, \theta) = \sum_{i=1}^n \log f(X_i, \theta)$
- the ML estimator $\hat{\theta}_{\text{ML}} \equiv \underset{\theta}{\operatorname{argmax}} \mathcal{L}(y|X, \theta) = \underset{\theta}{\operatorname{argmax}} \log \mathcal{L}(y|X, \theta)$

Solution (A5) $\implies y|X \sim \text{MVN}(X\beta, \sigma^2 I_n)$. Therefore $\mathcal{L}(y|X, \beta, \sigma^2) = (2\pi\sigma^2)^{-\frac{n}{2}} e^{-\frac{(y-X\beta)'(y-X\beta)}{2\sigma^2}}$, and $\log \mathcal{L}(y|X, \beta, \sigma^2) = -\frac{n}{2} \ln(2\pi) - \frac{n}{2} \ln(\sigma^2) - \frac{(y-X\beta)'(y-X\beta)}{2\sigma^2}$. The two FOCs of the maximization problem give an exact closed-form solution:³

$$\begin{cases} \frac{\partial \log \mathcal{L}}{\partial \beta} = 0 \iff \frac{-1}{2\hat{\sigma}^2} (-2X'y + 2X'X\hat{\beta}) = 0 \iff \hat{\beta}_{\text{ML}} = (X'X)^{-1}X'y \\ \frac{\partial \log \mathcal{L}}{\partial \sigma^2} = 0 \iff \frac{-n}{2\hat{\sigma}^2} + \frac{(y-X\hat{\beta})'(y-X\hat{\beta})}{2\hat{\sigma}^4} = 0 \iff n\hat{\sigma}^2 = (y-X\hat{\beta})'(y-X\hat{\beta}) \iff \hat{\sigma}_{\text{ML}}^2 = \frac{\hat{e}'\hat{e}}{n} = \frac{r'r}{n} \end{cases}$$

Properties (assuming (A1)-(A5))

- Finite samples

$\hat{\beta}_{\text{ML}}$ is **unbiased** $\mathbb{E}[\hat{\beta}_{\text{ML}}|X] = \mathbb{E}[(X'X)^{-1}X'y|X] = \mathbb{E}[(X'X)^{-1}X'(X\hat{\beta} + e)|X]$
 $= \mathbb{E}[\hat{\beta}|X] + \mathbb{E}[(X'X)^{-1}X'e|X] = \beta$

efficient $\mathbb{V}[\hat{\beta}_{\text{ML}}|X] = \mathbb{E}[(\hat{\beta} - \mathbb{E}[\hat{\beta}]) (\hat{\beta} - \mathbb{E}[\hat{\beta}])' | X] = \mathbb{E}[(\hat{\beta} - \beta) (\hat{\beta} - \beta)' | X]$
 $= \mathbb{E}[(X'X)^{-1}X'e((X'X)^{-1}X'e)' | X]$
 $= (X'X)^{-1}X' \mathbb{E}[ee'|X] X(X'X)^{-1} = \sigma^2(X'X)^{-1} = \mathbb{V}[\hat{\beta}_{\text{ML}}|X]$

normally distributed $\hat{\beta}_{\text{ML}} = \beta + (X'X)^{-1}X'e \sim \mathcal{N}(\beta, \sigma^2(X'X)^{-1})$

$\hat{\sigma}_{\text{ML}}^2$ is *downward biased* $\mathbb{E}[\hat{\sigma}_{\text{ML}}^2|X] = \frac{1}{n} \mathbb{E}[r'_i r_i|X] = \dots = \frac{n-k}{n} \sigma^2 < \sigma^2$
The variance is underestimated. The size of the bias will decrease as the sample size gets larger. To overcome this problem, we can compute the sample variance s^2 instead of $\hat{\sigma}_{\text{ML}}^2$.

- Asymptotics

$\hat{\beta}_{\text{ML}}$ is **asymptotically unbiased** as is unbiased
asymptotically efficient as is efficient
consistent as 1. is asymptotically unbiased, and 2. $\mathbb{V}[\hat{\beta}_{\text{ML}}|X] = \dots \xrightarrow[n \rightarrow \infty]{p} 0$

$\hat{\sigma}_{\text{ML}}^2$ is **asymptotically unbiased** as $\lim_{n \rightarrow \infty} \mathbb{E}[\hat{\sigma}_{\text{ML}}^2|X] = \lim_{n \rightarrow \infty} (\sigma^2 - \frac{k}{n} \sigma^2) = \sigma^2$
asymptotically efficient as $\sqrt{n}(\hat{\sigma}_{\text{ML}}^2 - \sigma^2) \xrightarrow{d} \mathcal{N}(0, 2\sigma^4)$
consistent as 1. is asymptotically unbiased, and 2. $\mathbb{V}[\hat{\sigma}_{\text{ML}}^2|X] = \frac{2\sigma^4(n-k)}{n^2} \xrightarrow[n \rightarrow \infty]{p} 0$

³The likelihood function must be differentiable in order to apply the derivative test for determining maxima. In some cases, the FOCs can be solved explicitly (e.g., the OLS estimator maximizes the likelihood of the linear regression model). Under most circumstances, however, numerical methods will be necessary to find the maximum of the likelihood function.

2.3 OLS estimator $\hat{\theta}_{\text{OLS}} = \{\hat{\beta}_{\text{OLS}}, \hat{\sigma}_{\text{OLS}}^2\}$

Definition The fit of a model $y = g(X, \beta)$ to each data point is measured by its residual $r_i \equiv y_i - g(x_i, \beta)$. The Ordinary Least Squares (OLS) estimator computes, in the context of a model linear in the parameters $g(X, \beta) = \sum_{j=1}^k \beta_j h_j(X)$, the values of the parameters that minimize the sum of the squares of the residuals:

$$\hat{\beta}_{\text{OLS}} \equiv \underset{\beta}{\operatorname{argmin}} \sum_{i=1}^n r_i^2$$

Solution The FOC of the minimization problem gives an exact closed-form solution (which the SOC guarantees is a minimum iff the matrix $X'X$ is positive definite):

$$\hat{\beta}_{\text{OLS}} = (X'X)^{-1}X'y = (X'X)^{-1}X'(X\beta + e) = \beta + (X'X)^{-1}X'e$$

With the residuals $r_i \equiv \hat{e}_i$ from the fit, we compute as estimator of σ^2 the statistic $\hat{\sigma}_{\text{OLS}}^2 \equiv s^2 \equiv \frac{r'r}{n-k} = \frac{\sum r_i^2}{n-k}$.

Properties [assuming (A1)-(A3)]

- Finite samples

$$(A3) \implies \hat{\beta}_{\text{OLS}} \text{ unbiased}$$

$$(A4) \implies \hat{\beta}_{\text{OLS}} \text{ efficient among linear unbiased estimators}$$

Gauss-Markov Theorem: in the semi-parametric⁴ linear regression model, we cannot show that $\hat{\beta}_{\text{OLS}}$ is efficient, but we can show that it is the most efficient among linear⁵ unbiased estimators. It is the **Best Linear Unbiased Estimator (BLUE)**.

$$\mathbb{V}[\hat{\beta}|X] = \mathbb{E}[(\hat{\beta} - \mathbb{E}[\hat{\beta}])(\hat{\beta} - \mathbb{E}[\hat{\beta}])' | X] = \dots = \sigma^2(X'X)^{-1}$$

$$(A5) \implies \hat{\beta}_{\text{OLS}} \text{ efficient}$$

In the *parametric* linear *normal* regression model ($e_i \sim \mathcal{N}(0, \sigma^2)$), $\hat{\beta}_{\text{OLS}}$ is equal to $\hat{\beta}_{\text{ML}}$. Therefore it is efficient, it is the **Best Unbiased Estimator (BUE)**.

$$(A4) \implies \hat{\sigma}_{\text{OLS}}^2 \text{ unbiased}^6 \quad \mathbb{E}[s^2|X] = \frac{1}{n-k} \mathbb{E}[r'r|X] = \dots = \frac{1}{n-k} \sigma^2(n-k) = \sigma^2$$

- Asymptotics

$\hat{\beta}_{\text{OLS}}$ is **asymptotically unbiased** as is unbiased

asymptotically normally distributed by a CLT, $\sqrt{n}(\hat{\beta}_{\text{OLS}} - \beta) \xrightarrow{d} \mathcal{N}(0, M_{\text{XX}}^{-1} M_{\text{XX}} M_{\text{XX}}^{-1})$

asymptotically efficient as $\sigma^2(X'X)^{-1}$ is the smallest possible asymptotic variance

consistent as 1. is asymptotically unbiased, and 2. $\mathbb{V}[\hat{\beta}_{\text{OLS}}|X] = \dots \xrightarrow[n \rightarrow \infty]{p} 0$

$\hat{\sigma}_{\text{OLS}}^2$ is **asymptotically unbiased** as is unbiased

asymptotically efficient as $\sqrt{n}(\hat{\sigma}_{\text{OLS}}^2 - \sigma^2) \xrightarrow{d} \mathcal{N}(0, 2\sigma^4)$

consistent as 1. is asymptotically unbiased, and 2. $\mathbb{V}[\hat{\sigma}_{\text{OLS}}^2|X] = \frac{2\sigma^4}{n-k} \xrightarrow[n \rightarrow \infty]{p} 0$

⁴The distribution of e is not fully characterized.

⁵Here, linearity does not refer to the linearity of the model w.r.t the parameters, but to the linearity of $\hat{\beta}$ w.r.t. y , such that y enters the equation linearly: $\beta_j = \lambda_1 y_1 + \dots + \lambda_n y_n$. Indeed, $\hat{\beta}_{\text{OLS}} = (X'X)^{-1}X'y$ is linear in y .

⁶The residuals have $n-k$ degrees of freedom (k parameters $\hat{\beta}$ are estimated; the model has an intercept and $k-1$ regressors). We must hence divide by $n-k$ in order to bias-adjust any statistic that uses the residuals as proxy for the true errors.

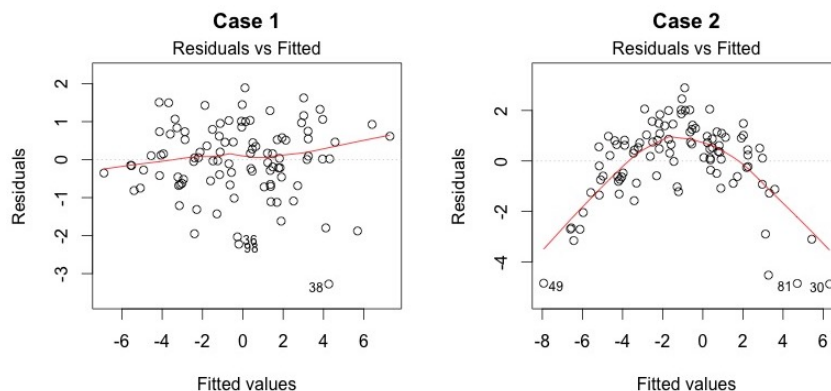
3 Post-estimation model diagnostics

After our statistical software has fit the model and spit out the estimates requested, we should check that the assumptions underlying these numbers hold. Several key assumptions can be diagnosed by looking at the residuals. Indeed, the residuals contain the variation in y and the patterns of the relationship between y and the explanatory variables that weren't explained by the model.

Four types of plots are presented below, with a description of the information that they provide w.r.t. an assumption of the CLRM.⁷ In each figure, the “Case 1” plot illustrates a case where the given assumption seems to be met relatively well, while the “Case 2” plot suggests the reverse.

1. “Residuals vs Fitted” plot — *Is there an unmodeled non-linear pattern?*

Residuals are plotted against fitted values. Are the residuals spread rather equally around a horizontal line, without distinct patterns? If so, that indicates that there are no non-linear relationships. If not, a non-linear relationship was not explained by the model and was therefore left out in the residuals. *Note: If y_i is discrete, such as in a logistic regression, then residuals are discrete. One shouldn't plot raw residuals, but binned residuals, i.e., divide the data equally into bins based on fitted values (s.t. each bin has the same number of points) and take the averages for each bin.*



2. “Scale-Location” plot — *Are the residuals homoscedastic?*

The square root of the absolute value of standardized residuals $\sqrt{|r_i|}$ is plotted against fitted values \hat{y}_i . Is the vertical spread of points uniform along x? A uniform spread indicates residuals have a uniform variance across the range of predicted values. The reverse suggests there is heteroscedasticity.

3. “Normal Q-Q” plot — *Are the residuals normally distributed?*

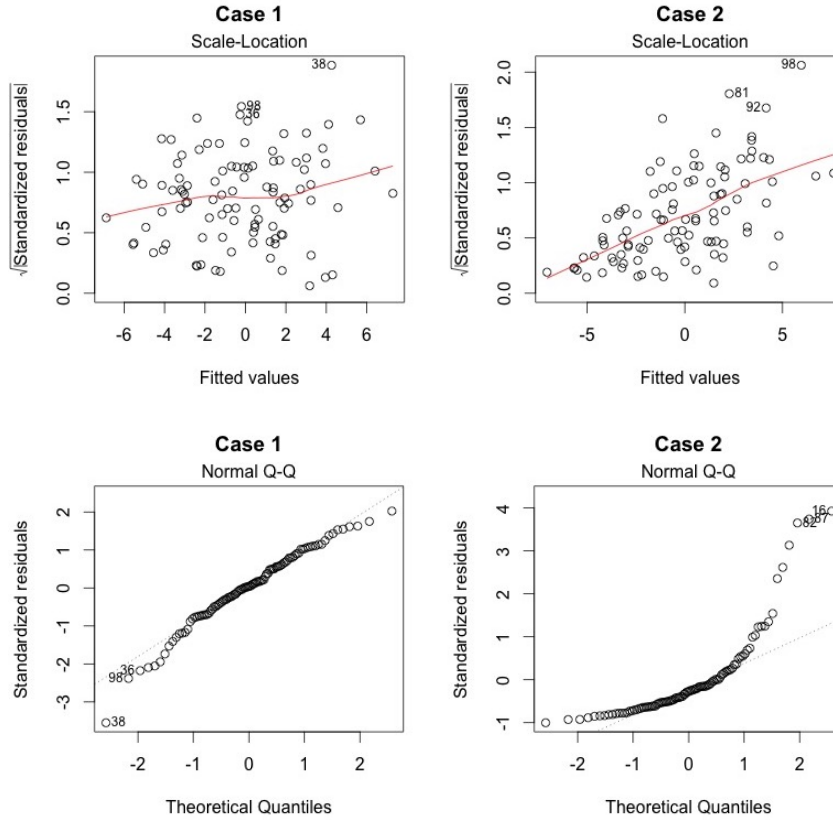
The quantiles of the residuals are plotted against the theoretical quantiles of the normal distribution. If the residuals are approximately normally distributed, we should see a roughly straight line.

4. “Residuals vs Leverage” plot — *Are there influential observations?*

First, let's distinguish outliers, high leverage points, and influential points:

- Outliers are observations with unusual outcome values y_i (i.e., that are considerably different from the rest of the data). They may not have a lot of influence on the regression line.
- High-leverage points are observations with unusual predictor values X_i . In linear regression, leverage measures how sensitive a fitted \hat{y}_i is to a change in the true y_i . High-leverage points will not have a lot of influence on the regression line if they lie close to it.

⁷These 4 particular types of plots are particularly easy to produce: they are built-in diagnostic plots for linear regression analysis in R (one need just run `plot()` to the fitted model object). The figures used here for illustrative purposes are taken from <https://data.library.virginia.edu/diagnostic-plots/>.

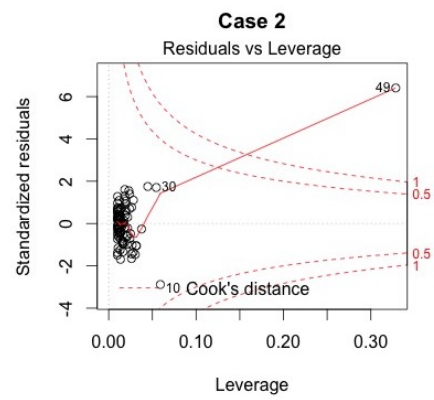
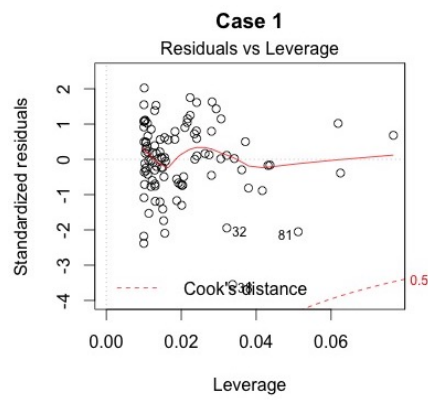


- Finally, influential points are observations whose removal from the data would cause a large change in the estimated regression line. I.e., they largely disagree with the trend.

A point has to have at least some leverage in order to be influential. To identify influential points, we can compute each observation's Cook's distance d_i , which measures the effect of omitting that observation on the combined parameter vector. Precisely, d_i has a component that reflects how well the model fits the i -th observation y_i and a component that measures how far that point is from the rest of the data. Points with $d_i > 1$ are generally considered to be influential.

In the "Residuals vs Leverage" plot, residuals are plotted against their leverage, and dotted red lines represent Cook's distances of 0.5 and 1. Points outside these lines have high Cook's distances, i.e., they have high influence.⁸

⁸If the Cook's distance lines aren't visible on the graph, it means that all points are well inside them — there are no influential points.



4 How to deal with non-spherical errors

4.1 OLS: Sandwich estimators

Assuming (A1)-(A3), by applying the CLT, we obtain the limit distribution of the rescaled $\hat{\beta}_{OLS}$:⁹

$$\sqrt{n}(\hat{\beta}_{OLS} - \beta) = \left(\frac{1}{n}X'X\right)^{-1} \frac{1}{\sqrt{n}}X'e = \underbrace{\left(\frac{1}{n}\sum_i x_i x_i'\right)^{-1}}_{\xrightarrow[n \rightarrow \infty]{p} M_{XX}} \underbrace{\frac{1}{\sqrt{n}}\sum_i x_i e_i}_{\xrightarrow[n \rightarrow \infty]{d} \mathcal{N}(0, M_{XX\Sigma X})} \xrightarrow[n \rightarrow \infty]{d} \mathcal{N}\left(0, M_{XX}^{-1} M_{XX\Sigma X} M_{XX}^{-1'}\right)$$

- where
- $M_{XX} \equiv \text{plim}\left(\frac{1}{n}X'X \mid X\right) \stackrel{!}{=} \lim\left(\mathbb{E}\left[\frac{1}{n}X'X \mid X\right]\right) = \lim\left(\frac{1}{n}X'X\right)$ is finite and $\neq 0$
 - $M_{XX\Sigma X} \equiv \text{plim}\left(\frac{1}{n}X'ee'X \mid X\right) = \lim\left(\mathbb{E}\left[\frac{1}{n}X'ee'X \mid X\right]\right) = \lim\left(\frac{1}{n}X' \mathbb{E}[ee' \mid X] X\right) \equiv \lim\left(\frac{1}{n}X'\Sigma X\right)$
 - Σ is the variance-covariance matrix of the error term: $\mathbb{E}[ee' \mid X]$

We talk of the limit distribution of $\sqrt{n}(\hat{\beta}_{OLS} - \beta)$, instead of $\hat{\beta}_{OLS}$, because $\hat{\beta}_{OLS}$ has a degenerate distribution with all mass at β . However, it would be more convenient to think of the distribution of $\hat{\beta}_{OLS}$ rather than carrying around $\sqrt{n}(\hat{\beta}_{OLS} - \beta)$. We do this by introducing the artifice of “asymptotic distribution”. We consider n large but not infinite, s.t. the asymptotics have kicked in, then we can drop the limits in the expressions (lim is dropped, plim becomes \mathbb{E}). We obtain $\hat{\beta}_{OLS}$ ’s asymptotic distribution:

$$\hat{\beta}_{OLS} \stackrel{a}{\sim} \mathcal{N}\left(\beta, \underbrace{\frac{1}{n}\left(\frac{1}{n}X'X\right)^{-1}\left(\frac{1}{n}X'\Sigma X\right)\left(\frac{1}{n}X'X\right)^{-1'}}_{\stackrel{a}{\mathbb{V}}[\hat{\beta}_{OLS}]}\right)$$

We need a **consistent** estimate of the asymptotic variance-covariance matrix $\stackrel{a}{\mathbb{V}}[\hat{\beta}_{OLS}]$ in order to do (sampling-based) statistical inference.¹¹ One approach is to use **sandwich estimators**.¹² The only unknown is Σ . We decompose the variance into its 3 $k \times k$ components: *bread*, *meat*, *bread*, and select a **consistent** estimator of the *meat* component that best represents our assumed error structure, to finally compute:

$$\hat{\mathbb{V}}[\hat{\beta}_{OLS}] \equiv \underbrace{(X'X)^{-1}}_{bread} \underbrace{X'\hat{\Sigma}X}_{meat} \underbrace{(X'X)^{-1'}}_{bread}$$

Error structure

► Spherical (A4)

$$\Sigma = \sigma^2 I, \text{ therefore } \stackrel{a}{\mathbb{V}} = (X'X)^{-1} X'\Sigma X (X'X)^{-1'} = \sigma^2 (X'X)^{-1} (X'X) (X'X)^{-1'} = \sigma^2 (X'X)^{-1}$$

⁹As in the entire document: (i) $\hat{\beta}_{OLS}$ refers to the vector of both the intercept and the slope coefficients, i.e., $k \geq 2$; (ii) for convenience, we drop the notation $|X$, however all features of $\hat{\beta}_{OLS}$ ’s distribution presented here are actually conditional on X .

¹⁰For a sample average \bar{Z}_n : by an LLN, $\text{plim } \bar{Z}_n = \lim \mathbb{E}[\bar{Z}_n]$.

¹¹The standard errors used in the t -test for $\hat{\beta}_{OLS}$ are none other than an estimate for $\sqrt{\stackrel{a}{\mathbb{V}}[\hat{\beta}_{OLS}]}$.

¹²All extremum estimators can actually be shown to be consistent and asymptotic normal, with an asymptotic variance matrix in the sandwich form: $\sqrt{n}(\hat{\theta} - \theta) \xrightarrow{d} \mathcal{N}(0, A(\theta)^{-1} B(\theta) A(\theta)^{-1})$. The sandwich algorithm presented here for OLS can be extended to all extremum estimators, e.g., ML and GMM. Which is not to say that it should be. [Freedman \(2006\)](#) points out that while White’s sandwich estimator often gives good results in OLS, the equivalent correction in ML does not necessarily make sense: “If the model is nearly correct, so are the usual standard errors, and robustification is unlikely to help much. On the other hand, if the model is seriously in error, the sandwich may help on the variance side, but the parameters being estimated by the ML are likely to be meaningless.” (If the specification—and hence the likelihood function—is incorrect, the parameter will be biased; why care about the variance of an estimator for the wrong parameter.)

We can consistently estimate the population variance σ^2 by the unbiased sample variance $s^2 = \frac{\sum_i r_i^2}{n-k}$, and hence $\hat{\mathbb{V}}^a$ by $\hat{\mathbb{V}}_s^a \equiv s^2 (X'X)^{-1}$. This expression is actually the Cramer-Rao lower bound, therefore $\hat{\beta}_{OLS}$ is **BLUE**.

► Not spherical

• Heteroskedastic

White (1980) proposes a non-parametric estimator for $\hat{\mathbb{V}}^a$, which provides heteroskedasticity-consistent (HC) or “robust” standard errors. They will be larger than those assuming homoskedasticity (which are downward-biased), as they account for the extra variation. HC SEs seem to have become best practice with large samples, as one can rarely assume homoskedastic errors.¹³¹⁴

The HC estimator uses $\hat{\Sigma}_H \equiv \frac{1}{n-k} \text{diag}[r_i^2]$, i.e., $X' \hat{\Sigma}_H X = \frac{1}{n-k} \sum_i r_i^2 x_i x_i'$. The resulting $\hat{\mathbb{V}}_H^a$ is consistent for $\hat{\mathbb{V}}^a$, even though r_i^2 is inconsistent for σ_i^2 .

• Autocorrelated

If errors are autocorrelated in any way (in time, space, both, by groups...), it means that the model is not capturing some feature of the DGP. There are essentially two ways to deal with this structure, in order to conduct proper inference:

1. Treat it as *substance*: incorporate the structure in the model. This also enables us to *study* it.
Ex: if errors are autocorrelated by group, model a multilevel data structure.
2. Treat it as *nuisance*: not incorporate the structure, but adjust for it after fitting the model.¹⁵
Ex: if errors are autocorrelated by group, cluster the standard errors.

The sandwich estimators below correspond to the second approach, and are built in a similar way:

- There is autocorrelation in e_i (in time, space, group...).
- The true covariance matrix of the errors Σ , and therefore the true asymptotic covariance matrix of $\hat{\beta}_{OLS}$, will contain these non-zero non-diagonal terms. We want our estimate of \mathbb{V} to be consistent, therefore we must produce a $\hat{\Sigma}$ that estimates these terms consistently.
- Two assumptions are made: 1. the process is **2nd order stationary**,¹⁶ 2. which covariance terms are potentially non-zero and the weights we give them (e.g., by using a kernel function).
- We estimate each of these terms by their sample equivalent. I.e., $\hat{\Sigma}$ is made of:
 - diagonal terms equal to White’s estimates;¹⁷
 - some non-zero non-diagonal terms that are the **sample autocovariances**.

¹³With a nonlinear conditional expectation function (CEF), the use of a linear model to approximate it should lead to heteroskedasticity (Angrist and Pischke, 2008, p.35). Indeed, as the quality of fit between the regression line and the CEF will vary with X , the residuals will be larger, on average, at values of X where the fit is poorer. The residual variance will increase with the square of the gap between the regression line $X\beta$ and the CEF $\mathbb{E}[y|X]$.

¹⁴Ideally, we would calculate an efficient estimator directly, instead of accepting an inefficient OLS and adjusting the SEs. The appropriate estimator is weighted least squares (WLS). However, its asymptotic efficiency rests on the correct specification of the pattern of heteroskedasticity. I.e., WLS is the better solution if we know the pattern, but we usually don’t.

¹⁵If we make no adjustments for this structure, default standard errors will generally overstate the estimator’s precision. Note that similarly as the note above, the first-best strategy would be to use generalized least squares (GLS), which produces an efficient estimator if we know the correct specification of the pattern of autocorrelation; but we usually don’t.

¹⁶A stationary process is “a stochastic process whose unconditional joint probability distribution does not change over the dimension of the process”. I.e., here:

- for a time series: the autocorrelation between 2 obs. that are m periods apart, is the same across the period;
- for a spatial process: the autocorrelation between 2 obs. that are apart by a distance d , is the same across the spatial field;
- for a process across groups: the autocorrelation between 2 obs. is fully determined by their group appartenance.

¹⁷These estimators are hence also heteroskedasticity-consistent.

* Serial (temporal) correlation

The dimension along which autocorrelation occurs is time.

Newey and West (1987) propose an estimator that accounts for serial correlation of unknown form in the errors of a single time series. It can be expanded to panel datasets, by estimating only correlations between lagged errors in the same cluster.

Consider a single time series $\{e_t\}$, and:

- Its autocovariance of lag l : $\gamma_e[t, t-l] \equiv \text{cov}[e_t, e_{t-l}] = \mathbb{E}[(X_t - \mu_t)(X_{t-l} - \mu_t)]$.
If the process is covariance-stationary, it is a function of the relative lag only: $\gamma_e[l]$
- Its bias-adjusted sample equivalent, for a sample $i = 1, \dots, T$: $g(l) = \frac{1}{T-k} \sum_{t=l+1}^T r_t r_{t-l}$.

The Newey-West estimator weights these covariance estimates using a triangular kernel function,¹⁸ s.t. the weight decreases linearly with the lag up to a chosen maximum lag L , and adds White's variance estimates:

$$\hat{\Sigma}_{\text{NW}} \equiv G(0) + \sum_{l=1}^L \left(1 - \frac{l}{L+1}\right) [G(l) + G(l)']$$

$$\begin{aligned} \text{Use } X' \hat{\Sigma}_{\text{NW}} X &\equiv \frac{1}{T-k} \sum_{t=1}^T r_t^2 \mathbf{x}_t \mathbf{x}_t' + \sum_{l=1}^L \left(1 - \frac{l}{L+1}\right) \left[\frac{1}{T-k} \sum_{t=l+1}^n r_t r_{t-l} \mathbf{x}_t \mathbf{x}_{t-l}' + \frac{1}{T-k} \sum_{t=l+1}^n r_t r_{t-l} \mathbf{x}_{t-l} \mathbf{x}_t' \right] \\ &= \frac{1}{T-k} \sum_{t=1}^T r_t^2 \mathbf{x}_t \mathbf{x}_t' + \frac{1}{T-k} \sum_{l=1}^L \left(1 - \frac{l}{L+1}\right) \sum_{t=l+1}^n r_t r_{t-l} (\mathbf{x}_t \mathbf{x}_{t-l}' + \mathbf{x}_{t-l} \mathbf{x}_t') \end{aligned}$$

Notes:

- $\hat{\Sigma}$ is consistent iff $L \rightarrow \infty$ and $\frac{L}{T^{1/4}} \rightarrow 0$ as $T \rightarrow \infty$, i.e., iff L grows slower than $T^{1/4}$. A common practice is hence to set L to the integer part of $T^{1/4}$.

* Spatial correlation

The dimension along which autocorrelation occurs is space.

This dimension is actually a dual dimension: while time or group appartenance are 1D, space is at least 2D. Conley (1999), under the supplementary assumption that the process is isotropic, proposes a consistent estimator for $\hat{\Sigma}$ that accounts for spatial correlation of unknown form in the errors, and for heteroskedasticity. It weights the sample covariances using a kernel function $k(s_i, s_j)$, where s_i is the location of observation i .

$$X' \hat{\Sigma}_{\text{Co}} X \equiv \frac{1}{n-k} \sum_{i=1}^n r_i^2 \mathbf{x}_i \mathbf{x}_i' + \frac{1}{n-k} \sum_{i=1}^n \sum_{j=1}^n k(s_i, s_j) e_i e_j \mathbf{x}_i \mathbf{x}_j'$$

Notes:

- Multiple choices of kernel are possible. Conley (2008) presents the uniform kernel but does not recommend it over another. $\hat{\Sigma}$ will be consistent if $\forall h, k(s, s+h) \rightarrow 1$ as $n \rightarrow \infty$, but slowly enough for the variance of $\hat{\Sigma}$ to collapse to zero. Assuming a stationary and isotropic process, $k(i, j)$ simplifies to a function of distance: $k(d_{ij})$. One can choose whichever distance metric fits one's context, e.g., a metric of economic distance.

¹⁸The modified Bartlett weights also ensure that $\hat{\Sigma}$ is positive semi-definite, which is required for the formation of asymptotic confidence interval and hypothesis testing.

- [Conley \(1999\)](#) shows that spatial dependence does not imply that SEs will necessarily increase. In his empirical example, 6 out of 9 spatial SE estimates are smaller than their iid counterparts.
- This estimator is very similar to the method of Kriging in geostatistics.

* Clustering

The dimension along which autocorrelation occurs is group appartenance.

Errors are correlated only within groups or “clusters”. I.e., $\mathbb{E}[e_i|x_i] = 0$, and $\mathbb{E}[e_i e_j | x_i, x_j] \neq 0$ iff $i, j \in$ same group g . The covariance matrix of the error term Σ has a block-diagonal structure.

Use $X'\hat{\Sigma}X = \frac{1}{n-k} \sum_{g=1}^G X_g' r_g r_g' X_g$. The resulting *cluster-robust* $\hat{\mathbb{V}}_c^a$ is consistent for \mathbb{V}^a .

Notes:

- This method is fully non-parametric, it allows for arbitrary dependence within a cluster.
- $\hat{\mathbb{V}}_c^a$ is also heteroskedasticity-consistent. It is typically $> \hat{\mathbb{V}}_H^a$ due to the addition of all non-diagonal terms within clusters.
- \triangle Do not use this estimator with too few clusters (rule of thumb: have > 40 clusters), as, like all sandwich estimators, it relies on asymptotics.¹⁹ [Cameron and Müller \(2015\)](#) recommends at least using critical values from the t_{G-1} distribution instead of the normal $\mathcal{N}(0,1)$. Note also that if clusters are unbalanced, the effective number of clusters is actually even lower.

\triangle Sandwich estimators are pointless in ML estimation These computations of adjusted SEs make sense only for the *linear* regression model estimated by OLS. In the case of a model that is nonlinear in the parameters (e.g., Logit and Probit models, which are usually estimated by ML), if for example the errors are heteroskedastic, then:

- the ML estimator of $\mathbb{V}^a[\hat{\beta}_{OLS}]$ is inconsistent (as in the linear model);
- but $\hat{\beta}_{ML}$ itself is also biased and inconsistent (unless the likelihood function is modified to correctly take into account the precise form of heteroskedasticity);

The reporting of robust standard errors in the context of nonlinear models such as Logit and Probit therefore doesn't make sense. What use is a consistent SE when the point estimate is wrong ([Freedman, 2006](#))? This is why it is important to test for model mis-specification (such as heteroskedasticity) when estimating models such as Logit, Probit, Tobit... Then, if need be, the model can be modified to take the heteroskedasticity into account before we estimate the parameters.

The reason why one can use a sandwich estimator in a linear model is because the coefficients and standard errors are determined separately. OLS coefficient estimates will be the same no matter what type of standard errors one chooses. In nonlinear models estimated by ML, the coefficients and standard errors can't be separated, they are jointly determined by maximizing the likelihood of finding y as it is given x . The problem applies to most of the standard models (binary, multinomial, ordered, and count) with the exception of GLS and poisson. Clustered standard errors will still correct the standard errors but they will now be attached to faulty coefficients. Essentially, you need to use something in the model to explain the clustering or you will bias your coefficients (and marginal effects/predicted probabilities) and not just your SEs. This is where fixed and random effects come back into play. By including either fixed effects or a random effect in the model you are using a variable or variables to directly model the problem. If done properly this can fix both

¹⁹The t -statistic $t_{\hat{\beta}} = \frac{\hat{\beta} - \beta}{\sqrt{\hat{\mathbb{V}}_c[\hat{\beta}]}} \underset{h_0}{\overset{a}{\rightsquigarrow}} \mathcal{N}(0,1)$. However, for finite G (and therefore, especially for small G), $t_{\hat{\beta}}$'s distribution

is unknown — even with normal errors. Intuitively, fewer clusters means there is less independent information in the sample (as the data are independent across clusters but not within). Using critical values from the standard normal distribution will downward-bias the variance estimate, leading to too narrow confidence intervals and over-rejection of the null.

the standard error issues and the biased coefficients. CCL: in linear models a sandwich estimator is good enough if you don't substantively care about group differences. In nonlinear models it can be a good aid to getting a better model but it will never be enough by itself.

References

- Angrist, J. and Pischke, J.-S. (2008). *Mostly harmless econometrics: An empiricist's companion*. Princeton University Press, Princeton, NJ, ISBN: [9781400829828](#), DOI: [10.1515/9781400829828](#).
- Cameron, A. C. and Miller, D. L. (2015). A practitioner's guide to cluster-robust inference. *The Journal of Human Resources*, 50(2):317–372, ISSN: 0022-166X.
- Conley, T. (1999). GMM estimation with cross sectional dependence. *Journal of Econometrics*, 92(1):1–45, DOI: [10.1016/S0304-4076\(98\)00084-0](#).
- Conley, T. G. (2008). Spatial Econometrics. In Durlauf, S. and Blume, L., editors, *The New Palgrave Dictionary of Economics*, volume 7, pages 741–747. 2nd edition, DOI: [10.1057/978-1-349-95121-5_2023-1](#).
- Freedman, D. A. (2006). On the So-Called “Huber Sandwich Estimator” and “Robust Standard Errors”. *The American Statistician*, 60(4):299–302, DOI: [10.1198/000313006X152207](#).
- Gelman, A., Hill, J., and Vehtari, A. (2020). *Regression and Other Stories*. Cambridge University Press, ISBN: [978-1-107-02398-7](#), DOI: [10.1017/9781139161879](#).
- Newey, W. K. and West, K. D. (1987). A simple, positive semi-definite, heteroskedasticity and autocorrelation consistent covariance matrix. *Econometrica*, 55(3):703–708, DOI: [10.2307/1913610](#).
- White, H. (1980). A heteroskedasticity-consistent covariance matrix estimator and a direct test for heteroskedasticity. *Econometrica*, 48(4):817–838, DOI: [10.2307/1912934](#).

A Misc.

A.1 Deriving the formula of the OLS estimator

Consider the multivariate linear regression model. We can write it as a system of n equations, or equivalently, in its matrix form:

$$y_i = X_i' \beta + e_i = \sum_{j=0}^k \beta_j x_{ij} + e_i, \quad e_i \stackrel{\text{iid}}{\sim} (0, \sigma^2), \quad \text{for } i = 1, \dots, n$$

$$y = X\beta + e, \quad e \sim (0, \sigma^2 I_n)$$

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} 1 & x_{11} & x_{12} & \cdots & x_{1k} \\ 1 & x_{21} & x_{22} & \cdots & x_{2k} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & x_{n2} & \cdots & x_{nk} \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_k \end{bmatrix} + \begin{bmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{bmatrix}$$

The OLS estimator is defined as the minimizer of the sum of squared residuals: $\hat{\beta}_{\text{OLS}} \equiv \underset{\beta}{\operatorname{argmin}} \text{SSR} \equiv \underset{\beta}{\operatorname{argmin}} \sum_{i=1}^n r_i^2 = \underset{\beta}{\operatorname{argmin}} \sum_{i=1}^n (y_i - X_i' \beta)^2$. We can solve for $\hat{\beta}_{\text{OLS}}$ using calculus:

- **Matrix form**

$$\hat{\beta}_{\text{OLS}} \equiv \underset{\beta}{\operatorname{argmin}} (y - X\beta)'(y - X\beta) = \underset{\beta}{\operatorname{argmin}} (y'y - \beta'X'y - y'X\beta + \beta'X'X\beta)$$

$$\begin{aligned} \text{FOC: } \frac{dS}{d\beta}(\hat{\beta}) = 0 &\iff \frac{d}{d\beta} (y'y - \beta'X'y - y'X\beta + \beta'X'X\beta) \Big|_{\beta=\hat{\beta}} = 0 \\ &\iff -X'y - (y'X)' + 2X'X\hat{\beta} \Big|_{\beta=\hat{\beta}} = 0^{20} \\ &\iff -2X'y + 2X'X\hat{\beta} = 0 \\ &\iff \hat{\beta} = (X'X)^{-1}X'y \end{aligned}$$

- **System of n equations**

$$\hat{\beta}_{\text{OLS}} \equiv \underset{\beta}{\operatorname{argmin}} \sum_i r_i^2 = \underset{\beta}{\operatorname{argmin}} \sum_i \left(y_i - \sum_k \beta_k x_{ik} \right)^2$$

$$\begin{aligned} \text{FOC: } \forall j, \frac{\partial \sum_i r_i^2}{\partial \beta_j} = 0 &\iff 2 \sum_i r_i \frac{\partial r_i}{\partial \beta_j} = 0 \\ &\iff \sum_i (y_i - \sum_k \beta_k x_{ik}) \frac{\partial (y_i - \sum_k \beta_k x_{ik})}{\partial \beta_j} = 0 \\ &\iff \sum_i (y_i - \sum_k \beta_k x_{ik}) (-x_{ij}) = 0 \\ &\iff \sum_i x_{ij} y_i = \sum_i x_{ij} \sum_k \beta_k x_{ik} \end{aligned}$$

²⁰Using denominator-layout notation, we have the following derivatives, or scalar-by-vector identities (where β and A are vectors): $\frac{d\beta'A}{d\beta} = \frac{dA'\beta}{d\beta} = A$, $\frac{d\beta'A\beta}{d\beta} = 2A\beta$.

For example, for the univariate regression model $y_i = \beta_0 + \beta_1 x_i + e_i$ (i.e., $x_{i0} = 1, x_{i1} = x_i$):

$$\begin{aligned}
 - \hat{\beta}_0: \quad \sum_i y_i &= \sum_i (\hat{\beta}_0 + \hat{\beta}_1 x_i) \iff \frac{1}{n} \sum_i y_i = \frac{1}{n} \sum_i \hat{\beta}_0 + \frac{1}{n} \sum_i \hat{\beta}_1 x_i \iff \bar{y} = \hat{\beta}_0 + \hat{\beta}_1 \bar{x} \\
 - \hat{\beta}_1: \quad \sum_i x_i y_i &= \sum_i x_i (\hat{\beta}_0 + \hat{\beta}_1 x_i) \iff \sum_i x_i y_i = \hat{\beta}_0 n \bar{x} + \hat{\beta}_1 \sum_i x_i^2 \\
 &\iff \frac{1}{n} \sum_i x_i y_i = (\bar{y} - \hat{\beta}_1 \bar{x}) \bar{x} + \hat{\beta}_1 \frac{1}{n} \sum_i x_i^2 \\
 &\iff \frac{1}{n} \sum_i (x_i y_i) - \bar{y} \bar{x} = \hat{\beta}_1 \left(\frac{1}{n} \sum_i (x_i^2) - \bar{x}^2 \right) \\
 &\iff \hat{\beta}_1 = \frac{\frac{1}{n} \sum_i (x_i y_i) - \bar{y} \bar{x}}{\frac{1}{n} \sum_i (x_i^2) - \bar{x}^2} = \dots = \frac{\sum_i (x_i - \bar{x})(y_i - \bar{y})}{\sum_i (x_i - \bar{x})^2}
 \end{aligned}$$

For the bivariate regression model $y_i = \beta_0 + \beta_1 x_i + \beta_2 z_i e_i$: add <https://scholar.princeton.edu/sites/default/files/bstewart/files/lecture6handout.pdf> slide 44...

A.2 Linear algebra — Positive-definite matrices

An $k \times k$ matrix A is **invertible** if there exists an $k \times k$ matrix B such that $AB = BA = I_k$.
A square matrix that is not invertible is called **singular or degenerate**.

The quasi-totality of square matrices are invertible.

Let M be an $k \times k$ symmetric real matrix, $\{\lambda_k\}$ its eigenvalues.

- M is **positive-definite** $\iff z' M z > 0$ for every vector $z \in \mathbb{R}^k \iff$ all $\{\lambda_k\}$ are > 0 .
- M is **positive semi-definite** $\iff z' M z \geq 0$ for every vector $z \in \mathbb{R}^k \iff$ all $\{\lambda_k\}$ are ≥ 0 .

Ex: The identity matrix I_k is positive-definite.

- Every positive definite matrix is invertible and its inverse is also positive definite.
- In statistics, the covariance matrix of a multivariate probability distribution is always symmetric and positive semi-definite; and it is positive definite unless one variable is an exact linear function of the others. Conversely, every positive semi-definite matrix is the covariance matrix of some multivariate distribution. Here we are talking about *population* covariance matrices. It is possible that the *sample* covariance matrix is singular, e.g., if there is exact collinearity, or when the number of observations is less than the number of variables.

A.3 Kernel functions for non-parametric statistics

A **kernel** is a non-negative real-valued integrable function $k()$, used as a **weighting function** in non-parametric estimation techniques. They are also called “window functions” (notably in time-series).

Some applications require the function to satisfy additional conditions, for instance:

- normalization: $\int_{-\infty}^{+\infty} k(u) du = 1$
In kernel density estimation, this ensures that the estimation produces a probability density function.
- symmetry: $\forall u, k(-u) = k(u)$
This ensures that the average of the corresponding distribution is equal to that of the sample used.

Examples of commonly used symmetric kernels, with the arbitrary bounded support $[-1, 1]$,²¹ i.e., $|u| \leq 1$:

²¹Such that $k(u) = 0$ for values of u outside the support.

- uniform: $k(u) = 1$
- triangular (Bartlett): $k(u) = 1 - |u|$
- parabolic (Epanechnikov): $k(u) = \frac{3}{4}(1 - u^2)$