

appendix

2023-01-22

Appendix

Effect size

Derivation

Effect size is usually defined as the difference of a parameter or a statistic derived from a sample. Centred on a conventional statistical test, we usually can deduce the effect size from the test statistic by substituting the null parameter value. When considering the diagnostics of residual departures, there exist many possibilities of test statistics for a variety of model assumptions. Meanwhile, diagnostic plots such as the residual plot have no general agreement on measuring how strong a model violation pattern is. To build a bridge between various residual-based tests, and the visual test, we focus on the shared information embedded in the testing procedures, which is the distribution of residuals. When comes to comparison of distribution, Kullback-Leibler divergence is a classical way to represent the information loss or entropy increase caused by the approximation to the true distribution, which in our case, the inefficiency due to the use of false model assumptions.

Following the terminology introduced by @kullback1951information, P represents the measured probability distribution, and Q represents the assumed probability distribution. The Kullback-Leibler divergence is defined as $\int_{-\infty}^{\infty} \log(p(x)/q(x))p(x)dx$, where $p(\cdot)$ and $q(\cdot)$ denote probability densities of P and Q .

Let $\mathbf{X}_a = (\mathbf{1}, \mathbf{X})$ denotes the p regressors with n observations, $\mathbf{R}_a = \mathbf{I} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$ denotes the residual operator, and let $\boldsymbol{\varepsilon} \sim N(\mathbf{0}, \sigma^2\mathbf{I})$ denotes the error. Using the Frisch-Waugh-Lovell theorem, residuals $\mathbf{e} = \mathbf{R}_a\boldsymbol{\varepsilon}$. Because $\text{rank}(\mathbf{R}_a) = n - p < n$, \mathbf{e} follows a degenerate multivariate normal distribution and does not have a density. Since the Kullback-Leibler divergence requires a proper density function, we need to simplify the covariance matrix of \mathbf{e} by setting all the off-diagonal elements to 0. Then, the residuals will assumed to follow $N(\mathbf{0}, \text{diag}(\mathbf{R}_a\sigma^2))$ under the null hypothesis that the model is correctly specified. If the model is however misspecified due to omitted variables \mathbf{Z} , or a non-constant variance \mathbf{V} , the distribution of residuals can be derived as $N(\mathbf{R}_a\mathbf{Z}\boldsymbol{\beta}_z, \text{diag}(\mathbf{R}_a\sigma^2))$ and $N(\mathbf{0}, \text{diag}(\mathbf{R}_a\mathbf{V}\mathbf{R}_a'))$ respectively.

By assuming both P and Q are multivariate normal density functions, the Kullback-Leibler divergence can be rewritten as

$$KL = \frac{1}{2} \left(\log \frac{|\Sigma_p|}{|\Sigma_q|} - n + \text{tr}(\Sigma_p^{-1}\Sigma_q) + (\mu_p - \mu_q)' \Sigma_p^{-1} (\mu_p - \mu_q) \right).$$

Then, we can combine the two residual departures into one formula

$$KL = \frac{1}{2} \left(\log \frac{|\text{diag}(\mathbf{R}_a\mathbf{V}\mathbf{R}_a')|}{|\text{diag}(\mathbf{R}_a\sigma^2)|} - n + \text{tr}(\text{diag}(\mathbf{R}_a\mathbf{V}\mathbf{R}_a')^{-1}\text{diag}(\mathbf{R}_a\sigma^2)) + \boldsymbol{\mu}_z^T (\mathbf{R}_a\mathbf{V}\mathbf{R}_a')^{-1} \boldsymbol{\mu}_z \right). \quad (1)$$

When there are omitted variables but constant error variance, the formula can be reduced to

$$KL = \frac{1}{2} (\boldsymbol{\mu}_z^T (\text{diag}(\mathbf{R}_a\sigma^2))^{-1} \boldsymbol{\mu}_z).$$

And when the model equation is correctly specified but the error variance is non-constant, the formula can be reduced to

$$KL = \frac{1}{2} \left(\log \frac{|diag(\mathbf{R}_a \mathbf{V} \mathbf{R}_a')|}{|diag(\mathbf{R}_a \sigma^2)|} - n + tr(diag(\mathbf{R}_a \mathbf{V} \mathbf{R}_a')^{-1} diag(\mathbf{R}_a \sigma^2)) \right).$$

Since we assume $\sigma = 1$ for the heteroskedasticity model, the final form of the formula is

$$KL = \frac{1}{2} \left(\log \frac{|diag(\mathbf{R}_a \mathbf{V} \mathbf{R}_a')|}{|diag(\mathbf{R}_a)|} - n + tr(diag(\mathbf{R}_a \mathbf{V} \mathbf{R}_a')^{-1} diag(\mathbf{R}_a)) \right).$$

Caculation

The effect size can be calculated using Equation 1 for each lineup by plugging in the data generated from the simulated model. However, since the term \mathbf{R}_a will differ for each sample, the effect size will be different even if the parameter values are the same. To overcome this undesired property, for the same set of parameter values, we simulate a sufficient large number of samples, and take the average of the effect size as the final effect size.

Experiment setup

Mapping of subjects to experimental factors

Mapping of subjects to experimental factors

Demographics

Data collection interface

Batch effect size

Sensitivity analysis for α

Other results