

Multicollinearity

Multicollinearity (or collinearity) is the fact that an independent variable (IV) is a linear combination of other IVs.

The stronger this dependency, the larger the standard errors for estimators of this variable (i.e. estimators become unstable, confidence intervals larger and test on parameters for this IV are more likely not significant).

A first indication for multicollinearity would be high first stage (bivariate) correlations between independent variables (IVs).

But as it is known that this method is not sufficient (Allison, 2003) two approaches are commonly used.

The first approach measures directly the correlation between an IV and all other IVs.

Let $R_{X,Other}$ be the correlation between X and all other IVs, then "Tolerance" (TOL) is defined as $TOL := 1 - R^2_{X,Other}$

A small value of TOL indicates that X is highly correlated to the rest of IVs.

TOL, therefore measures the percentage of variation of an IV that could not be "explained" by all the other IVs.

As the standard error of the estimated parameter of X is depending on $1/TOL$, a small TOL value causes a large standard error, large confidence interval and likely a not significant test result for the affected parameter.

Mostly the "Variance inflation factor" $VIF = 1/TOL$ is used instead, as it could be interpreted more easily. It shows directly how much the standard error of the estimation is inflated by the multicollinearity. A VIF of 25 for a variable X , for example, means that the standard error for the parameter of X is 5 times higher (inflated) due to the correlation between X and all the other IVs (multicollinearity).

There is no test available, neither for TOL nor for VIF. In practise $VIF > 10$ (equivalently $TOL < 0.1$) would indicate a multicollinearity problem.

The second approach uses the eigenvalues of the model matrix and leaves it to the user to decide whether eigenvalue are extreme, indicating that the dimension of the problem could be (or should be) reduced. There are three measures "eigenvalues" "condition index" and "condition number".

In case of no collinearity all eigenvalues would be 1.

Eigenvalues smaller or larger than 1 would indicate departures from the ideal situation.

"Too" small or large eigenvalues would indicate multicollinearity problems.

While the condition index is the ratio between a specific eigenvalue and the maximum of all eigenvalues, the condition number is the root of largest eigenvalue divided by the smallest.

As an informal rule a condition index between 10 and 100 or condition numbers between 15 and 30 would indicate weak to serious problems. Unfortunately, as there are different methods how to scale the model matrix or using the correlation matrix instead, this second approach has some draw backs.