

OVERIDENTIFICATION

Authored by
mohammad looti

October 14, 2025

RECOMMENDED CITATION

mohammad looti (2025). *OVERIDENTIFICATION*. PSYCHOLOGICAL SCALES. Retrieved from <https://scales.arabpsychology.com/?p=48333>

OVERIDENTIFICATION

Primary Disciplinary Field(s): Econometrics, Statistics, Structural Equation Modeling (SEM)

1. Core Definition

The term **overidentification** refers to a state within a statistical or econometric model where the available information, typically derived from observed data moments or constraints, exceeds the minimum necessary number of parameters required to uniquely and consistently estimate the model structure. Fundamentally, an overidentified model possesses more constraints or degrees of freedom than unknown parameters. This concept is crucial in fields relying on simultaneous equations, such as econometrics, and in psychometric modeling using techniques like **Structural Equation Modeling** (SEM).

The distinction between identification statuses is central to statistical inference. A model is considered **just-identified** (or exactly identified) if the number of known pieces of information precisely matches the number of parameters to be estimated, resulting in unique solutions but precluding formal tests of model validity based on internal constraints. Conversely, a model is **underidentified** if there are fewer pieces of information than parameters, making unique estimation impossible. **Overidentification**, however, is often considered desirable because the excess information imposes testable restrictions on the model. These excess constraints allow researchers to perform statistical tests to evaluate whether the model's assumptions align with the observed data structure.

In practical terms, the existence of more parameters than are strictly required to dictate the model accurately, as described in the initial definition, implies redundancy of information. This redundancy is harnessed to verify the internal consistency of the model specification itself. If the model is correctly specified, the redundant information should not contradict the estimated parameter values. The ability to test these restrictions forms the basis of crucial specification tests, such as the Sargan or Hansen J tests, which are essential for validating instrumental variable estimation methods.

2. Identification in the General Linear Model Context

While the most critical applications of overidentification lie outside simple Ordinary Least Squares (OLS), the fundamental principle is rooted in the broader context of the **General Linear Model** (GLM). In standard regression settings, identification typically involves ensuring that the design matrix is of full column rank, preventing perfect multicollinearity. However, in advanced GLM extensions, particularly those involving systems of equations or latent variables, identification becomes complex and relates to the ability to recover structural parameters from reduced-form

estimates.

In simultaneous equation models, where endogenous variables appear as predictors in other equations, OLS estimation is inconsistent due to endogeneity bias. To achieve consistent estimation, researchers employ techniques like Instrumental Variables (IV) or Two-Stage Least Squares (2SLS). It is within these instrumental variable frameworks that **overidentification** gains profound significance. The identification status determines not just whether consistent estimates can be obtained, but also whether the underlying assumptions about the instruments are empirically plausible.

The mathematical foundation for overidentification involves the rank condition and the order condition. The **order condition** states that for an equation to be identified, the number of excluded exogenous variables (instruments) must be greater than or equal to the number of endogenous variables included in the equation. When the number of excluded instruments is strictly greater than the number of endogenous variables, the equation is **overidentified**. This surplus of valid instruments provides the necessary statistical leverage to test their exogeneity assumptions.

3. Instrumental Variables and Overidentifying Restrictions

In the context of **Instrumental Variables** (IV) estimation, a key assumption is that the instruments are exogenous--meaning they are uncorrelated with the structural error term. When a model is overidentified, the researcher has more instruments than necessary to simply estimate the parameters. The redundant instruments imply **overidentifying restrictions** (OIRs). These restrictions state that all valid instruments should produce consistent estimates of the structural parameters, regardless of which subset of valid instruments is used.

The primary benefit of **overidentification** in IV estimation is the capacity to test the validity of the instruments themselves. Since the consistency of IV estimators hinges on the assumption of instrument exogeneity, researchers must be able to verify this condition, especially if the theoretical justification for instrument validity is weak or debatable. When the model is only just-identified, this crucial assumption cannot be statistically verified using internal data constraints.

If the model is severely overidentified--that is, possessing many more instruments than endogenous variables--it suggests a high degree of redundancy. While beneficial for testing, this can also imply that the instruments are weak or that the linear combination of instruments used in the first stage of 2SLS might introduce efficiency issues if many are only weakly correlated with the endogenous variables. Therefore, careful selection and validation of instruments are paramount, leveraging the constraints imposed by the **overidentification** status.

4. Overidentification in Structural Equation Modeling (SEM)

In structural equation modeling (SEM), the concept of identification is related to the relationship between the observed covariance matrix (Σ) and the parameters (θ) of the hypothesized structural model. SEM aims to determine if the hypothesized model structure can uniquely reproduce the population covariance matrix. An SEM is **overidentified** when the number of non-redundant elements in the observed covariance matrix (i.e., the unique variances and covariances) is greater than the number of free parameters that must be estimated.

The degrees of freedom (df) in an SEM are calculated as the number of data points (unique elements in Σ) minus the number of free parameters. If $df > 0$, the model is **overidentified**. This positive degree of freedom is essential because it is the numerator of the model chi-square test statistic. This statistic assesses the overall goodness-of-fit by testing the null hypothesis (H_0) that the population covariance matrix is equal to the covariance matrix implied by the model ($\Sigma = \Sigma(\theta)$).

Unlike in IV estimation where overidentification primarily tests instrument validity, in SEM, **overidentification** allows the entire structural relationship hypothesized by the model to be tested against the empirical data. A well-fitting, overidentified model provides evidence that the complex web of hypothesized relationships (e.g., factor loadings, latent variable correlations, structural paths) is plausible. Conversely, a poor fit implies that the restrictions imposed by the overidentified structure are statistically incompatible with the observed data.

5. Tests for Overidentifying Restrictions: The Sargan and Hansen J Tests

To formally test the validity of the **overidentifying restrictions** (OIRs) in econometric models, researchers utilize specific statistical tools. The most prominent of these is the **Sargan test** (sometimes called the Sargan-Hansen test or simply the J-test). This test is designed to verify the joint null hypothesis that all instruments are valid--meaning they are uncorrelated with the error term--and that the structural equation is correctly specified.

The test statistic is asymptotically chi-squared distributed with degrees of freedom equal to the number of overidentifying restrictions (i.e., the number of instruments minus the number of endogenous regressors). A small J-statistic, corresponding to a high p-value, suggests that the null hypothesis should not be rejected, indicating that the instruments appear valid and the model is consistent with the data constraints. A large J-statistic, conversely, leads to the rejection of the null hypothesis, signalling potential misspecification, often due to invalid instruments (endogeneity of the instruments themselves).

It is important to differentiate between the Sargan test and the more general **Hansen J test**. The original Sargan test is applicable when the errors are homoskedastic. When using Generalized

Method of Moments (GMM) estimation, or when heteroskedasticity is suspected, the Hansen J test is preferred because it is robust to arbitrary forms of heteroskedasticity. In either case, the fundamental role of the test remains the same: it provides a quantitative assessment of the constraints imposed by the **overidentification** status of the model.

6. Consequences and Practical Challenges

The existence of **overidentification** is a double-edged sword. While it enables crucial specification tests, rejection of the overidentifying restrictions presents a significant practical challenge. When the J-test rejects the null hypothesis, the researcher faces an **identification problem** of a different sort: determining the source of the contradiction. The rejection only indicates that at least one of the restrictions is false, but it does not specify which instrument is invalid or if the assumed structural form is flawed.

Researchers must then engage in diagnostic detective work. This often involves running subsets of the instruments or using difference-in-Sargan tests (also known as C-tests) if applicable to isolate the potentially invalid instruments. Arbitrarily dropping instruments, however, carries the risk of data mining or making the model just-identified, thereby losing the ability to test instrument validity internally.

Furthermore, in cases of severe **overidentification**, if the instruments are weak (i.e., only poorly correlated with the endogenous variable), the resulting IV estimates can be severely biased towards the inconsistent OLS estimates, even if the instruments are technically valid. Consequently, the advantages of overidentification (testability) must be weighed against the potential statistical liabilities of including weak or marginally valid instruments, which can degrade the finite sample performance of the estimator.

7. Handling Overidentified Models

When an econometric model is successfully **overidentified** and the specification tests (like the J-test) do not reject the null hypothesis, the researcher benefits from the efficiency gains offered by methods like 2SLS or GMM, compared to just-identified estimators. The estimation procedure leverages all available information optimally.

However, if the overidentifying restrictions are rejected, several remedial steps are typically considered.

Revising Instrument Set: If external knowledge or theory suggests that a subset of instruments is more likely to be valid than others, the researcher may systematically remove suspect instruments and re-test the resulting, less-overidentified model.

Model Respecification: The rejection may indicate that the underlying structural equation is fundamentally misspecified. This requires theoretical adjustment of the model, such as including additional variables, changing functional forms, or altering the assumed linearity.

Alternative Estimation Techniques: In complex settings, switching from 2SLS to **Generalized Method of Moments** (GMM) often proves advantageous. GMM is a flexible estimation technique that naturally handles overidentified systems by minimizing a quadratic form of the sample moment conditions. When errors are not independent and identically distributed, GMM provides asymptotically efficient estimates, often preferred over 2SLS in highly overidentified contexts.

8. Key Characteristics

Redundancy of Information: The model structure imposes more constraints or provides more moments than the number of parameters requiring estimation.

Testable Restrictions: The surplus information creates overidentifying restrictions (OIRs) that allow for formal statistical testing of the model's internal consistency and assumptions (e.g., instrument exogeneity).

Efficiency Potential: Successfully estimated overidentified models (using efficient methods like 2SLS or GMM) typically yield more efficient parameter estimates than just-identified models.

Context Specificity: Crucial in simultaneous equation models (econometrics) and covariance structure analysis (SEM), where endogeneity or latent variables necessitate complex identification strategies.

Further Reading

[Sargan test \(Wikipedia\)](#)

[Instrumental variables estimation \(Wikipedia\)](#)

[Structural equation modeling \(Wikipedia\)](#)

[Generalized method of moments \(Wikipedia\)](#)