

# ITEM RESPONSE THEORY (IRT)

Authored by  
**mohammad looti**

October 12, 2025

## RECOMMENDED CITATION

mohammad looti (2025). *ITEM RESPONSE THEORY (IRT)*. PSYCHOLOGICAL SCALES.  
Retrieved from <https://scales.arabpsychology.com/?p=42305>

## ITEM RESPONSE THEORY (IRT)

**Primary Disciplinary Field(s):** Psychometrics, Educational Measurement, Statistics

**Proponents:** Frederic M. Lord, Georg Rasch, Benjamin D. Wright

### 1. Core Principles

Item Response Theory (IRT) is a modern probabilistic framework of measurement used extensively in psychometrics and educational testing. It states that an individual's response to a test item is a function of an underlying, unobservable characteristic, known as a **latent trait** or ability (often denoted as  $\theta$ , or Theta). Unlike the older **Classical Test Theory (CTT)**, IRT focuses on the mathematical relationship between the probability of answering a specific item correctly and the level of the test taker's ability in that latent trait. This probabilistic relationship is defined by a non-linear regression function, typically logistic in nature, resulting in the **Item Characteristic Curve (ICC)**.

The central tenet of IRT is the separation of measurement parameters: item parameters (e.g., difficulty, discrimination) are considered independent of the specific sample of people taking the test, and person parameters (ability estimates) are considered independent of the specific set of items administered. This concept of **invariance** is highly advantageous, allowing IRT to overcome the limitations of CTT where item difficulty and person ability estimates are inextricably linked to the group and the specific test form used. The models utilized in Item Response Theory differ significantly in terms of the number of item parameters they include, ranging from highly restrictive one-parameter models to complex four-parameter and polytomous models designed for specific testing scenarios.

### 2. Theoretical Foundations and Advantages Over CTT

The development of IRT was largely driven by the inherent measurement inadequacies of Classical Test Theory (CTT). CTT relies on the observed score being a summation of a true score and random error, but its measures (such as reliability coefficients and standard errors) are heavily dependent on the specific population and the exact composition of the test. Consequently, CTT struggles to compare scores across different test forms or to select optimal items for a particular individual's ability level. IRT addresses these issues by modeling the relationship at the individual item level, rather than at the test level.

A key theoretical advantage of IRT lies in its ability to provide a **Standard Error of Measurement (SEM)** that is specific to the test taker's ability level, rather than a single, constant SEM for all scores, as is common in CTT. This means IRT can precisely quantify the measurement precision (or lack thereof) at different points along the latent trait continuum. For example, a test designed for

high-ability individuals will provide very precise measures (low SEM) for high-ability individuals but imprecise measures for low-ability individuals. This precision tailoring is critical for advanced applications like **Computerized Adaptive Testing (CAT)**, where items are selected dynamically to maximize measurement efficiency at the individual level.

### 3. Historical Development and Key Pioneers

Item Response Theory originated in the mid-20th century, emerging from the psychometric and mathematical statistics fields. Early conceptualization of latent trait models can be traced back to the work of scholars like Louis Guttman, but the formal statistical framework was solidified by Frederic M. Lord in the 1950s and 1960s, particularly through his landmark work on the relationship between item characteristics and test scores. Lord's contributions established the basic mathematical models for linking observed responses to latent ability.

A separate, yet parallel, line of development was championed by the Danish mathematician **Georg Rasch**. Rasch's work emphasized the requirement for "specific objectivity," leading to the development of the simplest IRT model, often referred to as the **Rasch Model** or the One-Parameter Logistic Model (1PL). The Rasch approach imposes stricter mathematical constraints than other IRT models, ensuring that the comparison of any two persons is independent of which items are used, and the comparison of any two items is independent of which persons are used. This emphasis on objective, fundamental measurement properties made the Rasch model particularly popular in educational and clinical settings prioritizing fairness and measurement linearity. The widespread adoption of various IRT models accelerated significantly in the late 1970s and 1980s, coinciding with improvements in computational power necessary to handle the iterative estimation algorithms (like Marginal Maximum Likelihood Estimation) required to calculate the complex parameters.

### 4. Fundamental Concepts: Item Characteristic Curves (ICC)

The **Item Characteristic Curve (ICC)** is the graphical representation of the IRT model, defining the relationship between the latent trait ( $\theta$ , on the x-axis) and the probability of a correct response ( $P(\theta)$ , on the y-axis). The shape and position of the ICC are entirely determined by the item parameters included in the chosen model, allowing test developers to visualize the functional properties of each item.

**Difficulty Parameter (b):** This parameter represents the location of the curve along the ability axis. Mathematically, it is the ability level ( $\theta$ ) at which a test taker has a 50% chance of answering the item correctly (assuming a two-parameter model where the guessing parameter is zero). Items with low 'b' values are easier, while those with high 'b' values are harder.

**Discrimination Parameter (a):** This parameter corresponds to the slope of the ICC at its point of

inflection (the difficulty level). A steeper slope indicates a higher discrimination parameter. Items with high discrimination are superior at distinguishing between individuals just above and just below the item's difficulty level. Items with low discrimination provide little information about the latent trait.

**Guessing Parameter (c):** Also known as the lower asymptote, this parameter estimates the probability that a person with extremely low ability will still answer the item correctly due to chance. This is crucial for multiple-choice tests, as it defines the minimum probability of success, even when ability is near zero.

The ICC is also used to derive the **Item Information Function (IIF)**, which plots the precision of measurement provided by a single item across the range of the latent trait. Items contribute the most information (i.e., are most useful for measurement) at the ability level corresponding to their difficulty parameter. By aggregating the IIFs for all items on a test, the **Test Information Function (TIF)** is created, which reveals where the entire test provides the greatest precision.

## 5. IRT Models: Parameter Variations

The nomenclature of IRT models is based on the number of item parameters allowed to vary during calibration, reflecting differing levels of complexity and assumptions about item behavior.

**One-Parameter Logistic Model (1PL):** This is the simplest model, exemplified by the Rasch Model. It incorporates only the **difficulty parameter (b)**. It operates under the strict assumption that all items have equal discrimination (the 'a' parameter is fixed, often to 1.0) and that the guessing parameter ('c') is zero. This model is often preferred when the goal is to achieve fundamental measurement properties and when the focus is on constructing scales with linear interval properties.

**Two-Parameter Logistic Model (2PL):** This model allows both the **difficulty parameter (b)** and the **discrimination parameter (a)** to vary. The 2PL model fits data better than the 1PL model when items vary substantially in their quality or effectiveness in differentiating between test takers. This flexibility, however, requires larger samples for stable estimation compared to the 1PL model.

**Three-Parameter Logistic Model (3PL):** The 3PL model includes the parameters for difficulty (b), discrimination (a), and the **pseudo-guessing parameter (c)**. This model is generally the standard choice for multiple-choice, high-stakes standardized tests, as it explicitly accounts for the effect of random guessing on the probability of a correct response. The inclusion of the guessing parameter is crucial for accurate ability estimates at the lower end of the ability spectrum.

Beyond these foundational models, specialized IRT models have been developed for specific response formats, such as **polytomous IRT models** (e.g., Graded Response Model, Partial Credit

Model) used for items that allow for multiple ordered categories of response, such as Likert scales or essay grading rubrics.

## 6. Model Fit and Assumptions

For IRT parameter estimates and resulting ability measures to be valid, the empirical data must adequately fit the chosen mathematical model. Assessing **model fit** is a critical step in any IRT application, relying on the verification of two primary assumptions.

The first crucial assumption is **Unidimensionality**, meaning that all items on the test measure a single, common underlying latent trait. If a test is intended to measure mathematical ability but some items also require high reading comprehension skills, the resulting measurement structure may be multidimensional. If a standard unidimensional IRT model (like 1PL, 2PL, or 3PL) is applied to multidimensional data, the item and person parameter estimates can be biased and misleading. While statistical methods, such as factor analysis (e.g., exploratory or confirmatory factor analysis) or dimensionality checks (e.g., Parallel Analysis), are often used to screen for unidimensionality, perfect unidimensionality is rare in practice; researchers must determine if the deviation from unidimensionality is substantively negligible.

The second major assumption is **Local Independence**. This assumption dictates that, once the effect of the underlying latent trait ( $\theta$ ) has been statistically accounted for, the responses to the items should be independent of one another. In simpler terms, the only reason an individual's response to Item 1 should correlate with their response to Item 2 is their shared level of the latent trait. Violations of local independence often occur when tests contain "testlets" (groups of items based on a common stimulus, like a reading passage) or when time constraints cause fatigue or speediness effects on later items. Statistical checks, often involving residual analysis or specialized dependency tests, are necessary to ensure that local independence holds, thereby validating the integrity of the item parameter estimates.

## 7. Criticisms and Methodological Limitations

Despite its sophisticated measurement capabilities, Item Response Theory faces several methodological and practical limitations that warrant consideration. One of the most significant drawbacks is the substantial **sample size requirement**. Robust and stable estimation of item parameters, particularly the discrimination ('a') and guessing ('c') parameters in the 2PL and 3PL models, necessitates large samples--often requiring several hundred to over a thousand test takers. This requirement can render IRT impractical or unreliable for smaller-scale research projects, clinical trials, or highly specialized assessments involving limited populations.

Furthermore, the mathematical complexity of IRT models requires a higher degree of statistical expertise compared to CTT. The process of model selection, parameter estimation, and rigorous

model-data fit testing often demands specialized software and trained psychometricians. If the fundamental assumptions, such as unidimensionality or local independence, are violated and this violation is not addressed through advanced modeling techniques (like Multidimensional IRT or specialized dependence models), the resulting parameter estimates may lack validity. Critics also point out that the 1PL (Rasch) model, while mathematically elegant and simple, is often too restrictive for real-world tests where item discrimination inevitably varies. Conversely, more complex models (2PL, 3PL) may introduce model identification issues or require unrealistic sample sizes to estimate all parameters reliably.

### Further Reading

[Item Response Theory - Wikipedia](#)

[Latent Variable Model - Wikipedia](#)

[Psychometrics - Wikipedia](#)

[Rasch Model - Wikipedia](#)

ARABPSYCHOLOGY.COM