# Let's talk about Thurstone & Co.: An information-theoretical model for comparative judgments, and its statistical translation

Jose Manuel Rivera Espejo[a,*], Tine van Daal[a], Sven De Maeyer[a], Steven Gillis[b]

[a]*University of Antwerp, Training and education sciences,*

[b]*University of Antwerp, Linguistics,*

**Abstract**

This study revisits Thurstone's law of comparative judgments (CJ) by addressing two key limitations in traditional approaches. Firstly, it addresses the overreliance on the assumptions of Thurstone's Case V in the statistical analysis of CJ data. Secondly, it addresses the apparent disconnect between CJ's approach to trait measurement and hypothesis testing. We put forward a systematic approach based on causal analysis and Bayesian statistical methods, which results in a model that facilitates a more comprehensive understanding of the factors influencing CJ experiments while offering a robust statistical translation. The new model accommodates unequal dispersions and correlations between stimuli, enhancing the reliability and validity of CJ's trait estimation, thereby ensuring the accurate measurement and interpretation of comparative data. The paper highlights the relevance of this updated framework for modern empirical research, particularly in education and social sciences. This contribution advances current research methodologies, providing a robust foundation for future applications in diverse fields.

*Keywords:* Causal analysis, Directed Acyclic Graphs, Bayesian statistical methods, Thurstonian model, Comparative judgement, Probability, Statistical modeling

## 1. Introduction

In *comparative judgment* (CJ) studies, judges assess a specific trait or attribute across various stimuli by performing pairwise comparisons (Thurstone, 1927b,a). Each comparison produces a dichotomous outcome, indicating which stimulus is perceived to exhibit

---

*Corresponding author

*Email addresses:* `JoseManuel.RiveraEspejo@uantwerpen.be` (Jose Manuel Rivera Espejo), `tine.vandaal@uantwerpen.be` (Tine van Daal), `sven.demaeyer@uantwerpen.be` (Sven De Maeyer), `steven.gillis@uantwerpen.be` (Steven Gillis)

a higher trait level. For example, when assessing text quality, judges compare pairs of written texts (the stimuli) to determine the relative quality each text exhibit (the trait) (Laming, 2004; Pollitt, 2012b; Whitehouse, 2012; van Daal et al., 2016; Lesterhuis, 2018a; Coertjens et al., 2017; Goossens and De Maeyer, 2018; Bouwer et al., 2023).

Numerous studies have documented the effectiveness of CJ in assessing traits and competencies over the past decade. These studies have emphasized three aspects of the method's effectiveness: its reliability, validity, and practical applicability. Research on reliability indicates that CJ requires a relatively small number of pairwise comparisons (Verhavert et al., 2019; Crompvoets et al., 2022) to produce trait scores that are as precise and consistent as those generated by other assessment methods (Coertjens et al., 2017; Goossens and De Maeyer, 2018; Bouwer et al., 2023). Furthermore, evidence suggests that the reliability and time efficiency of CJ are comparable, if not superior, to those of other assessment methods when employing adaptive comparison algorithms (Pollitt, 2012b; Verhavert et al., 2022; Mikhailiuk et al., 2021). Meanwhile, research on validity suggests that scores generated by CJ can accurately represent the traits under measurement (Whitehouse, 2012; van Daal et al., 2016; Lesterhuis, 2018a; Bartholomew et al., 2018; Bouwer et al., 2023), while research on practical applicability highlights the method's versatility across both educational and non-educational contexts (Kimbell, 2012; Jones and Inglis, 2015; Bartholomew et al., 2018; Jones et al., 2019; Marshall et al., 2020; Bartholomew and Williams, 2020; Boonen et al., 2020).

Nevertheless, despite the increasing number of CJ studies, unsystematic and fragmented research approaches have left several critical issues unaddressed. The present study primarily focuses on two: the overreliance on the assumptions of Thurstone's Case V in

the statistical analysis of CJ data, and the apparent disconnect between CJ's approach to trait measurement and hypothesis testing. The following sections begin with a brief overview of Thurstone's theory and a detailed examination of these issues. Subsequently, the study introduces a theoretical model for CJ that builds upon Thurstone's theory, alongside its statistical translation, designed to address the two concerns simultaneously.

## 2. Thurstone's theory

In its most general form, Thurstone's theory addresses pairwise comparisons where a single judge evaluates multiple stimuli (Thurstone, 1927a, pp. 267). The theory posits that two key factors determine the dichotomous outcome of these comparisons: the discriminal process of each stimulus and their discriminal difference. The *discriminal process* captures the psychological impact each stimulus exerts on the judge or, more simply, his perception of the stimulus trait. The theory assumes that the discriminal process for any given stimulus forms a Normal distribution along the trait continuum (Thurstone, 1927a, pp. 266). The mode (mean) of this distribution, known as the *modal discriminal process*, indicates the stimulus position on this continuum, while its dispersion, referred to as the *discriminal dispersion*, reflects variability in the perceived trait of the stimulus.

Figure 1a illustrates hypothetical discriminal processes along a quality trait continuum for two written texts. The figure indicates that the modal discriminal process for Text B is positioned further along the continuum than that of Text A ($S_B > S_A$), suggesting that Text B exhibits higher quality. Additionally, the figure highlights that Text B has a broader distribution compared to Text A, which arises from its larger discriminal

3

dispersion $(\sigma_B > \sigma_A)$.



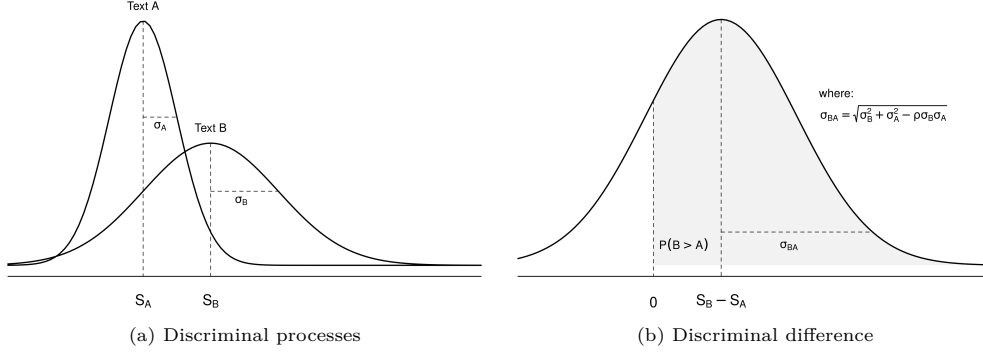| (a) Discriminal processes | (b) Discriminal difference |
|---|---|

Figure 1: Hypothetical discriminal processes and discriminant difference along a quality trait continuum for two written texts.

However, since the individual discriminal processes of the stimuli are not directly observable, the theory introduces *the law of comparative judgment*. This law posits that in pairwise comparisons, a judge perceives the stimulus with a discriminal process positioned further along the trait continuum as possessing more of the trait (Bramley, 2008, pp. 251). This suggests that the relative distance between stimuli, rather than their absolute positions on the continuum, likely defines the outcome of pairwise comparisons. Indeed, the theory assumes that the difference between the underlying discriminal processes of the stimuli, referred to as *the discriminal difference*, determines the observed dichotomous outcome. Moreover, the theory assumes that because the individual discriminal processes form a Normal distribution on the continuum, the discriminal difference will also conform to a Normal distribution (Andrich, 1978). In this distribution, the mode (mean) represents the relative separation between the stimuli, and its dispersion indicates the variability of that separation.

Figure 1b illustrates the distribution of the discriminal difference for the hypothetical

texts depicted in panel Figure 1a. The figure indicates that the judge perceives Text B as having significantly higher quality than Text A. This conclusion rests on two key observations: the positive difference between their modal discriminal processes ($S_B - S_A > 0$) and the probability area where the discriminal difference distinctly favors Text B over Text A, represented by the shaded gray area denoted as $P(B > A)$. As a result, the dichotomous outcome of this comparison is more likely to favor Text B over Text A.

## 3. The two critical issues in CJ literature

This section examines the two critical issues in the CJ literature that serve as the primary focus of this study. The first is the overreliance on Thurstone's Case V assumptions in the statistical analysis of CJ data. The second is the apparent disconnect between CJ's approach to trait measurement and hypothesis testing.

### 3.1. The Case V and the statistical analysis of CJ data

Thurstone observed that the general form of the theory, outlined in Section 2, created a trait scaling problem. The model required estimating more "unknown" parameters than the available pairwise comparisons (Thurstone, 1927a, pp. 267). To address this issue and facilitate the practical application of the theory, he developed five cases derived from this general form, each case progressively incorporating additional simplifying assumptions into the model.

In Case I, Thurstone assumed that pairs of stimuli maintained a constant correlation across all comparisons. In Case II, he allowed multiple judges to make comparisons instead of restricting evaluations to a single judge. In Case III, he introduced the assumption of zero correlation between stimuli. In Case IV, he assumed stimuli exhibited

Table 1: Thurstones cases and their asumptions

| Assumption | General form | Thurstone's | | | | | BTL model |
|---|---|---|---|---|---|---|---|
| | | Case I | Case II | Case III | Case IV | Case V | |
| Discriminal process (distribution) | Normal | Normal | Normal | Normal | Normal | Normal | Logistic |
| Discriminal dispersion (between stimuli) | Different | Different | Different | Different | Similar | Equal | Equal |
| Correlation (between stimuli) | One per pair | Constant | Constant | Zero | Zero | Zero | Zero |
| How many judges compare? | Single | Single | Multiple | Multiple | Multiple | Multiple | Multiple |

similar dispersions. Finally, in Case V, he replaced this assumption with the condition that stimuli had equal discriminal dispersions. Table 1 summarizes the assumptions of the general form and the five cases. For an in-depth discussion of these cases and their progression, refer to Thurstone (1927a) and Bramley (2008, pp. 248–253).

Despite relying on the most extensive set of simplifying assumptions (Bramley, 2008, pp. 253; Kelly et al., 2022, pp. 677), Case V remains the most widely used case in the CJ literature. This popularity stems mainly from its simplified statistical representation in the Bradley-Terry-Luce (BTL) model (Bradley and Terry, 1952; Luce, 1959). The BTL model mirrors the assumptions of Case V, with one key difference: while Case V assumes a Normal distribution for the stimuli's discriminal processes, the BTL model uses the more mathematically tractable Logistic distribution (Andrich, 1978; Bramley, 2008, pp. 254) (see Table 1). This substitution has little impact on the model's estimation or interpretation, as the Normal and Logistic distributions share similar statistical properties, differing only by a scaling factor of approximately 1.7 (van der Linden, 2017a, pp. 16).

However, Thurstone originally developed Case V to provide a "rather coarse scaling" of traits (Thurstone, 1927a, pp. 269), prioritizing statistical simplicity over precision

in trait measurement (Kelly et al., 2022, pp. 677). He explicitly warned against its untested application, stating that its use "should not be made without (an) experimental test" (Thurstone, 1927a, pp. 270), acknowledging that some assumptions could prove problematic when researchers assess complex traits or heterogeneous stimuli (Thurstone, 1927b, pp. 376). Consequently, given that modern CJ applications frequently involve such traits and stimuli, two main assumptions of Case V and, by extension, of the BTL model may not consistently hold in theory or practice: the assumption of equal dispersion and zero correlation between stimuli.

### 3.1.1. The assumption of equal dispersions between stimuli

According to the theory, discrepancies in the discriminal dispersions of stimuli shape the distribution of the discriminal difference, exerting a direct influence on the outcome of pairwise comparisons. Figure 2a illustrates this idea using a thought experiment in which a researcher can observe the discriminal processes for the texts shown in Figure 1a. Furthermore, the figure assumes that the discriminal dispersion for Text A remains constant and that the texts are uncorrelated ($\rho = 0$). The figure reveals that an increase in the uncertainty associated with the perception of Text B in comparison to Text A, ($\sigma_B - \sigma_A$), broadens the distribution of their discriminal difference. This broadening affects the probability area where the discriminal difference distinctly favors Text B over Text A, expressed as $P(B > A)$, ultimately influencing the comparison outcome. Additionally, the figure reveals that when the discriminal dispersions of the texts are equal ($\sigma_B - \sigma_A = 0$), the discriminal difference is more likely to favor Text B over Text A (shaded gray area), compared to situations where their dispersions differ.

In experimental practice, however, this process occurs in reverse. Researchers first

7

observe the comparison outcome and then use the BTL model to infer the discriminal difference between the stimuli and their respective discriminal processes (Thurstone, 1927b, pp. 373). For instance, when researchers observe a large sample of outcomes favoring Text B over Text A and correctly assume equal dispersions between the texts, the BTL model estimates a discriminal difference distribution that accurately represents the "true" discriminal difference of the texts. This scenario is illustrated in Figure 2a, where the model's discriminal difference distribution aligns with the "true" distribution, represented by the thick continuous line corresponding to $\sigma_B - \sigma_A = 0$. The accuracy of this estimated discriminal difference ensures reliable estimates of the texts' discriminal processes [(citation needed?){style="color:red;"}. Therefore, it is reasonable to assume that the outcome's ability to reflect the "true" differences between stimuli largely depends on the validity of the model's assumptions, particularly the assumption of equal dispersions.

However, Thurstone contended that the assumption of equal dispersions may not hold when researchers assess complex traits or heterogeneous stimuli (Thurstone, 1927b, pp. 376), as these traits and stimuli can introduce judgment discrepancies due to their unique features (van Daal et al., 2016; Lesterhuis, 2018b; Chambers and Cunningham, 2022). Indeed, evidence of this violation may already exist in the CJ literature as misfit statistics, which measure judgment discrepancies associated with specific stimuli (Pollitt, 2004, pp. 12; Goossens and De Maeyer, 2018, pp. 20). For example, labeling texts as "misfits" indicates that comparisons involving these texts result in more judgment discrepancies than those involving other texts (Pollitt, 2012a,b; van Daal et al., 2016; Goossens and De Maeyer, 2018). These discrepancies, in turn, suggest that

the discriminal differences for "misfit" texts have broader distributions, indicating that their discriminal processes may also exhibit more variation than that of other texts. A similar reasoning applies to "misfit" judges, whose evaluations deviate substantially from the shared consensus due to the unique characteristics of the stimuli or the judges themselves. Moreover, these "misfit" judges and their deviations can introduce additional statistical and measurement issues, which we discuss in Section 3.1.2.

Therefore, incorrectly assuming equal dispersions between stimuli is not without harm, as it can cause the BTL model to introduce significant statistical and measurement issues. For instance, the model may overestimate how accurately the outcome reflects the "true" discriminal differences between stimuli. This overestimation can lead researchers to draw spurious conclusions about these differences (McElreath, 2020, pp. 370) and, by extension, of the stimuli underlying discriminal processes. Figure 2a illustrates this issue when the model's discriminal difference distribution aligns with the thick continuous line for $\sigma_B - \sigma_A = 0$, while the "true" discriminal difference follows a discontinuous line where $\sigma_B - \sigma_A \neq 0$. Moreover, if researchers recognize that misfit statistics highlight critical differences in dispersions, the usual CJ practice of excluding stimuli based on these statistics (Pollitt, 2012a,b; van Daal et al., 2016; Goossens and De Maeyer, 2018) can unintentionally discard valuable information. This exclusion can introduce bias into trait estimates (Zimmerman, 1994; McElreath, 2020, chap. 12). The direction and magnitude of these biases are often unpredictable, as they depend on which stimuli are excluded from the analysis.
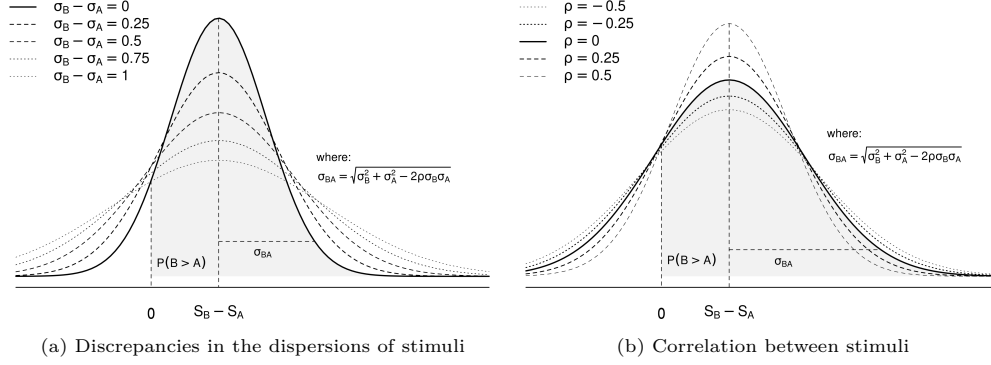
σ_B − σ_A = 0
σ_B − σ_A = 0.25
σ_B − σ_A = 0.5
σ_B − σ_A = 0.75
σ_B − σ_A = 1

where:
$\sigma_{BA} = \sqrt{\sigma_B^2 + \sigma_A^2 - 2\rho\sigma_B\sigma_A}$

P(B > A)   σ_BA

0   S_B − S_A

(a) Discrepancies in the dispersions of stimuli

ρ = −0.5
ρ = −0.25
ρ = 0
ρ = 0.25
ρ = 0.5

where:
$\sigma_{BA} = \sqrt{\sigma_B^2 + \sigma_A^2 - 2\rho\sigma_B\sigma_A}$

P(B > A)   σ_BA

0   S_B − S_A

(b) Correlation between stimuli

Figure 2: The effect of dispersion discrepancies and stimulus correlation on the distribution of the discriminal difference.

### 3.1.2. The assumption of zero correlation between stimuli

Denoted by $\rho$, the correlation measures the dependence of a judge's perception of the trait in one stimulus on his perception of the same trait in another. Like the discriminal dispersions, this correlation shapes the distribution of the discriminal difference and directly influences the outcomes of pairwise comparisons. Figure 2b illustrates this concept, again assuming the researcher can observe the discriminal processes for the texts shown in Figure 1a. The figure also considers that the discriminal dispersions for both texts remain constant.

Figure 2b reveals that as the correlation between the texts increases, the distribution of their discriminal difference becomes narrower. This narrowing affects the area under the curve where the discriminal difference distinctly favors Text B over Text A, denoted as $P(B > A)$, thus influencing the comparison outcome. Furthermore, the figure shows that when two texts are independent or uncorrelated $(\rho = 0)$, their discriminal difference is less likely to favor Text B over Text A (shaded gray area), compared to scenarios when the texts are highly correlated.

10

However, in experimental practice, researchers typically follow this process in reverse. For example, when they observe a large sample of outcomes favoring Text B over Text A and correctly assume zero correlation between the texts, the BTL model estimates a discriminal difference distribution that accurately represents the "true" discriminal difference of the texts. This scenario is illustrated with Figure 2b when the discriminal difference distribution of the model aligns with the "true" distribution, represented by the thick continuous line corresponding to $\rho = 0$. The accuracy of this discriminal difference estimation, in turn, ensures reliable estimates for the discriminal process of the texts [(citation needed?){style="color:red;"}.

Notably, Thurstone's Case V and the BTL model assume independent discriminal processes across comparisons. Thurstone attributed this independence to the cancellation of potential judges' biases, driven by two opposing and equally weighted effects occurring during the pairwise comparisons (Thurstone, 1927a, pp. 268). Andrich (1978) mathematically demonstrated this cancellation using the BTL model under the assumption of discriminal processes with additive biases. However, it is easy to imagine at least two scenarios where the zero correlation assumption almost certainly does not hold: when the pairwise comparison involves multidimensional, complex traits with heterogeneous stimuli and when an additional hierarchical structure is relevant to the stimuli.

In the first scenario, the intricate aspects of multidimensional, complex traits may introduce dependencies between the stimuli due to certain judges' biases that resist cancellation. Research on text quality suggests that when judges evaluate these traits, they often rely on various intricate characteristics of the stimuli to form their judgments (van Daal et al., 2016; Lesterhuis, 2018b; Chambers and Cunningham, 2022). These addi-

11

tional relevant characteristics, which are unlikely to be equally weighted or opposing, can unevenly influence judges' perceptions, creating biases in their judgments and, ultimately, introducing dependencies between stimuli (van der Linden, 2017b, pp. 346). For example, this could occur when a judge assessing the argumentative quality of a text places more weight on its grammatical accuracy than other judges, ultimately favoring texts with fewer errors but weaker arguments. While direct evidence for this specific scenario is lacking, studies such as Pollitt and Elliott (2003) demonstrate the presence of such biases, supporting the idea that the factors influencing pairwise comparisons may not always cancel out.

In the second scenario, the shared context or inherent connections created by additional hierarchical structures may further introduce dependencies between stimuli, a statistical phenomenon commonly known as clustering (Everitt and Skrondal, 2010). Although the CJ literature acknowledges the presence of such hierarchical structures, the statistical handling of this extra source of dependency between stimuli has been inadequate. For example, when CJ data includes multiple samples of stimuli from the same individuals, researchers often rely on (average) estimated BTL scores to conduct subsequent analyses and tests at the individual hierarchical level (Bramley and Vitello, 2019; Boonen et al., 2020; Bouwer et al., 2023; van Daal et al., 2017; Jones et al., 2019; Gijsen et al., 2021). However, this approach can introduce additional statistical and measurement issues, which we discuss in Section 3.2.

In any case, similar to Section 3.1.1, incorrectly assuming zero correlation between stimuli can lead the BTL model to introduce various statistical and measurement issues. For instance, the model could over- or underestimate the accuracy of the outcome in reflect-

ing the "true" discriminal differences between stimuli. This over- underestimation may result in spurious inferences about these differences and, by extension, about the stimuli's discriminal processes (Hoyle, 2023, pp. 341). Figure 2b also illustrates this scenario when the model's discriminal difference distribution aligns with the thick continuous line for $\rho = 0$, while the "true" discriminal difference follows any discontinuous line where $\rho \neq 0$. This missaligment can be due to the overlook of additional relevant traits, such as judges' biases, which cause dimensional mismatches in the BTL model, artificially inflating the reliability of the trait (Hoyle, 2023, pp. 341) or, even worse, introduce bias into the trait's estimates (Ackerman, 1989). Furthermore, researchers who exclude judges based on misfit statistics can risk discarding valuable information, further biasing the trait's estimates (Zimmerman, 1994; McElreath, 2020, chap. 12). Lastly, researchers who fail to account for hierarchical (grouping) structures can reduce the precision of model parameter estimates, which may amplify the overestimation of the trait's reliability (Hoyle, 2023, pp. 482).

## 3.2. The disconnect between trait measurement and hypothesis testing

Building on the previous section, it is clear that, despite its limitations, the BTL model is commonly used as the measurement model in CJ assessments. A measurement model specifies how manifest variables contribute to the estimation of latent variables (Everitt and Skrondal, 2010). For example, when evaluating text quality, researchers use the BTL model to process the dichotomous outcomes resulting from the pairwise comparisons (the manifest variables) to estimate scores that reflect the underlying quality level of the texts (the latent variable) (Laming, 2004; Pollitt, 2012b; Whitehouse, 2012; van Daal et al., 2016; Lesterhuis, 2018a; Coertjens et al., 2017; Goossens and De Maeyer, 2018; Bouwer

et al., 2023).

Researchers then typically use these estimated BTL scores, or their transformations, to conduct additional analyses or hypothesis tests. For example, these scores have been used to identify 'misfit' judges and stimuli (Pollitt, 2012b; van Daal et al., 2016; Goossens and De Maeyer, 2018), detect biases in judges' ratings (Pollitt and Elliott, 2003; Pollitt, 2012b), calculate correlations with other assessment methods (Goossens and De Maeyer, 2018; Bouwer et al., 2023), or test hypotheses related to the underlying trait of interest (Bramley and Vitello, 2019; Boonen et al., 2020; Bouwer et al., 2023; van Daal et al., 2017; Jones et al., 2019; Gijsen et al., 2021).

However, the statistical literature advises caution when using estimated scores for additional analyses and tests. A key consideration is that BTL scores are parameter estimates that inherently carry uncertainty. Ignoring this uncertainty can bias the analysis and reduce the precision of hypothesis tests. Notably, the direction and magnitude of such biases are often unpredictable. Results may be attenuated, exaggerated, or remain unaffected depending on the degree of uncertainty in the scores and the actual effects being tested (Kline, 2023, pp. 25; Hoyle, 2023, pp. 137). Finally, the reduced precision in hypothesis tests diminishes their statistical power, increasing the likelihood of committing type-I or type-II errors (McElreath, 2020).

In aggregate, researchers' inadequate handling of violations to the assumptions of equal dispersion and zero correlation between stimuli, along with the apparent disconnect between CJ's approach to trait measurement and hypothesis testing, can undermine the reliability of the trait estimates and ultimately compromise its validity (Perron and Gillespie, 2015, pp. 2). Consequently, adopting a more systematic and integrated approach to

14

examining what happens when judges compare two stimuli could offer several statistical and measurement benefits, including addressing these issues.

## 4. An updated theoretical and statistical model for CJ

This section presents a theoretical model for CJ that extends Thurstone's theory. The model systematically incorporates all factors involved when judges make pairwise comparisons. Additionally, the section develops the statistical translation of the theoretical model based on assumptions informed by the CJ theory.

*4.1. The theoretical model*

The (latent) discriminal difference of the stimuli directly determines the (manifest) outcome of the pairwise comparisons

$$D_{k_j a_{i1} b_{i2}} \quad O_{k_j a_{i1} b_{i2}}$$

Figure 3: Theoretical model A1$

The (latent) "perceived" discriminal processes for the stimuli directly determines their discriminal difference

$$T^*_{k_j a_{i1}}$$

$$D_{k_j a_{i1} b_{i2}} \quad O_{k_j a_{i1} b_{i2}}$$

$$T^*_{k_j b_{i2}}$$

Figure 4: Theoretical model A2$

15

The (latent) "true" discriminal processes for the stimuli and the judges' biases directly determines their (latent) "perceived" discriminal processes



Figure 5: Theoretical model A3$

without loosing generality, the (latent) "perceived" and "true" discriminal processes for the stimuli can be depicted in a vector for each judge, as in

*4.2. From theory to statistics*

## 5. Discussion

*5.1. Findings*

*5.2. Limitations and further research*

## 6. Conclusion

**Declarations**

**Financial interests:** The authors have no relevant financial interest to disclose.

**Non-financial interests:** The authors have no relevant non-financial interest to disclose.

**Ethics approval:** The University of Antwerp Research Ethics Committee has confirmed that no ethical approval is required.

**Consent to participate:** Not applicable

**Consent for publication:** All authors have read and agreed to the published version of the manuscript.

**Availability of data and materials:** No data was utilized in this study.

**Code availability:** All the code utilized in this research is available in the digital document located at: https://jriveraespejo.github.io/paper2_manuscript/.

**AI-assisted technologies in the writing process:** The authors used ChatGPT, an AI language model, during the preparation of this work. They occasionally employed the tool to refine phrasing and optimize wording, ensuring appropriate language use and enhancing the manuscript's clarity and coherence. The authors take full responsibility for the final content of the publication.

**CRediT authorship contribution statement:** *Conceptualization:* S.G., S.DM., T.vD., and J.M.R.E; *Methodology:* S.DM., T.vD., and J.M.R.E; *Software:* J.M.R.E.;

# References

Ackerman, T., 1989. Unidimensional irt calibration of compensatory and noncompensatory multidimensional items. Applied Psychological Measurement 13, 113–127. doi:10.1177/014662168901300201.

Andrich, D., 1978. Relationships between the thurstone and rasch approaches to item scaling. Applied Psychological Measurement 2, 451–462. doi:10.1177/014662167800200319.

Bartholomew, S., Nadelson, L., Goodridge, W., Reeve, E., 2018. Adaptive comparative judgment as a tool for assessing open-ended design problems and model eliciting activities. Educational Assessment 23, 85–101. doi:10.1080/10627197.2018.1444986.

Bartholomew, S., Williams, P., 2020. Stem skill assessment: An application of adaptive comparative judgment, in: Anderson, J., Li, Y. (Eds.), Integrated Approaches to STEM Education. Advances in STEM Education. Springer, pp. 331–349. doi:10.1007/978-3-030-52229-2_18.

Boonen, N., Kloots, H., Gillis, S., 2020. Rating the overall speech quality of hearing-impaired children by means of comparative judgements. Journal of Communication Disorders 83, 1675–1687. doi:10.1016/j.jcomdis.2019.105969.

Bouwer, R., Lesterhuis, M., De Smedt, F., Van Keer, H., De Maeyer, S., 2023. Comparative approaches to the assessment of writing: Reliability and validity of benchmark rating and comparative judgement. Journal of Writing Research 15, 497–518. doi:10.17239/jowr-2024.15.03.03.

Bradley, R., Terry, M., 1952. Rank analysis of incomplete block designs: I. the method of paired comparisons. Biometrika 39, 324–345. doi:10.2307/2334029.

Bramley, T., 2008. Paired comparison methods, in: Newton, P., Baird, J., Goldsteing, H., Patrick, H., Tymms, P. (Eds.), Techniques for monitoring the comparability of examination standards. GOV.UK., pp. 246—300. URL: https://assets.publishing.service.gov.uk/media/5a80d75940f0b62305b8d734/2007-comparability-exam-standards-i-chapter7.pdf.

Bramley, T., Vitello, S., 2019. The effect of adaptivity on the reliability coefficient in adaptive comparative judgement. Assessment in Education: Principles, Policy and Practice 71, 1–25. doi:10.1080/0969594X.2017.1418734.

Chambers, L., Cunningham, E., 2022. Exploring the validity of comparative judgement: Do judges attend to construct-irrelevant features? Frontiers in Education doi:10.3389/feduc.2022.802392.

Coertjens, L., Lesterhuis, M., Verhavert, S., Van Gasse, R., De Maeyer, S., 2017. Teksten beoordelen

met criterialijsten of via paarsgewijze vergelijking: een afweging van betrouwbaarheid en tijdsinvestering. Pedagogische Studien 94, 283–303. URL: https://repository.uantwerpen.be/docman/irua/e71ea9/147930.pdf.

Crompvoets, E., Béguin, A., Sijtsma, K., 2022. On the bias and stability of the results of comparative judgment. Frontiers in Education 6. doi:10.3389/feduc.2021.788202.

Everitt, B., Skrondal, A., 2010. The Cambridge Dictionary of Statistics. Cambridge University Press.

Gijsen, M., van Daal, T., Lesterhuis, M., Gijbels, D., De Maeyer, S., 2021. The complexity of comparative judgments in assessing argumentative writing: An eye tracking study. Frontiers in Education 5. doi:10.3389/feduc.2020.582800.

Goossens, M., De Maeyer, S., 2018. How to obtain efficient high reliabilities in assessing texts: Rubrics vs comparative judgement, in: Ras, E., Guerrero Roldán, A. (Eds.), Technology Enhanced Assessment, Springer International Publishing. pp. 13–25. doi:10.1007/978-3-319-97807-9_2.

Hoyle, R.e., 2023. Handbook of Structural Equation Modeling. Guilford Press.

Jones, I., Bisson, M., Gilmore, C., Inglis, M., 2019. Measuring conceptual understanding in randomised controlled trials: Can comparative judgement help? British Educational Research Journal 45, 662–680. doi:10.1002/berj.3519.

Jones, I., Inglis, M., 2015. The problem of assessing problem solving: can comparative judgement help? Educational Studies in Mathematics 89, 337–355. doi:10.1007/s10649-015-9607-1.

Kelly, K., Richardson, M., Isaacs, T., 2022. Critiquing the rationales for using comparative judgement: a call for clarity. Assessment in Education: Principles, Policy & Practice 29, 674–688. doi:10.1080/0969594X.2022.2147901.

Kimbell, R., 2012. Evolving project e-scape for national assessment. International Journal of Technology and Design Education 22, 135–155. doi:10.1007/s10798-011-9190-4.

Kline, R., 2023. Principles and Practice of Structural Equation Modeling. Methodology in the Social Sciences, Guilford Press.

Laming, D., 2004. Marking university examinations: Some lessons from psychophysics. Psychology Learning & Teaching 3, 89–96. doi:10.2304/plat.2003.3.2.89.

Lesterhuis, M., 2018a. The validity of comparative judgement for assessing text quality: An assessor's perspective. Ph.D. thesis. University of Antwerp. URL: https://hdl.handle.net/10067/1548280151162165141.

20

Lesterhuis, M., 2018b. When teachers compare argumentative texts: Decisions informed by multiple complex aspects of text quality. L1-Educational Studies in Language and Literature 18, 1–22. doi:`10.17239/L1ESLL-2018.18.01.02`.

Luce, R., 1959. On the possible psychophysical laws. The Psychologcal Review 66, 482–499. doi:`10.1037/h0043178`.

Marshall, N., Shaw, K., Hunter, J., Jones, I., 2020. Assessment by comparative judgement: An application to secondary statistics and english in new zealand. New Zealand Journal of Educational Studies 55, 49–71. doi:`10.1007/s40841-020-00163-3`.

McElreath, R., 2020. Statistical Rethinking: A Bayesian Course with Examples in R and STAN. Chapman and Hall/CRC.

Mikhailiuk, A., Wilmot, C., Perez-Ortiz, M., Yue, D., Mantiuk, R., 2021. Active sampling for pairwise comparisons via approximate message passing and information gain maximization, in: 2020 25th International Conference on Pattern Recognition (ICPR), pp. 2559–2566. doi:`10.1109/ICPR48806.2021.9412676`.

Perron, B., Gillespie, D., 2015. Reliability and Measurement Error, in: Key Concepts in Measurement. Oxford University Press. Pocket guides to social work research methods. chapter 4. doi:`10.1093/acprof:oso/9780199855483.003.0004`.

Pollitt, A., 2004. Let's stop marking exams, in: Proceedings of the IAEA Conference, University of Cambridge Local Examinations Syndicate, Philadelphia. URL: https://www.cambridgeassessment.org.uk/images/109719-let-s-stop-marking-exams.pdf.

Pollitt, A., 2012a. Comparative judgement for assessment. International Journal of Technology and Design Education 22, 157—170. doi:`10.1007/s10798-011-9189-x`.

Pollitt, A., 2012b. The method of adaptive comparative judgement. Assessment in Education: Principles, Policy and Practice 19, 281—300. doi:`10.1080/0969594X.2012.665354`.

Pollitt, A., Elliott, G., 2003. Finding a proper role for human judgement in the examination system. URL: https://www.cambridgeassessment.org.uk/Images/109707-monitoring-and-investigating-comparability-a-proper-role-for-human-judgement.pdf. research & Evaluation Division.

Thurstone, L., 1927a. A law of comparative judgment. Psychological Review 34, 482–499. doi:`10.1037/h0070288`.

Thurstone, L., 1927b. Psychophysical analysis. American Journal of Psychology , 368–89URL: https:

//brocku.ca/MeadProject/Thurstone/Thurstone_1927g.html. last accessed 20 december 2024.

van Daal, T., Lesterhuis, M., Coertjens, L., Donche, V., De Maeyer, S., 2016. Validity of comparative judgement to assess academic writing: examining implications of its holistic character and building on a shared consensus. Assessment in Education: Principles, Policy & Practice 26, 59–74. doi:`10.1080/0969594X.2016.1253542`.

van Daal, T., Lesterhuis, M., Coertjens, L., van de Kamp, M., Donche, V., De Maeyer, S., 2017. The complexity of assessing student work using comparative judgment: The moderating role of decision accuracy. Frontiers in Education 2. doi:`10.3389/feduc.2017.00044`.

van der Linden, W. (Ed.), 2017a. Handbook of Item Response Theory: Models. volume 1 of *Statistics in the Social and Behavioral Sciences Series*. CRC Press.

van der Linden, W. (Ed.), 2017b. Handbook of Item Response Theory: Statistical Tools. volume 2 of *Statistics in the Social and Behavioral Sciences Series*. CRC Press.

Verhavert, S., Bouwer, R., Donche, V., De Maeyer, S., 2019. A meta-analysis on the reliability of comparative judgement. Assessment in Education: Principles, Policy and Practice 26, 541–562. doi:`10.1080/0969594X.2019.1602027`.

Verhavert, S., Furlong, A., Bouwer, R., 2022. The accuracy and efficiency of a reference-based adaptive selection algorithm for comparative judgment. Frontiers in Education 6. doi:`10.3389/feduc.2021.785919`.

Whitehouse, C., 2012. Testing the validity of judgements about geography essays using the adaptive comparative judgement method. URL: https://filestore.aqa.org.uk/content/research/CERP_RP_CW_24102012_0.pdf?download=1. aQA Education.

Zimmerman, D., 1994. A note on the influence of outliers on parametric and nonparametric tests. The Journal of General Psychology 121, 391–401. doi:`10.1080/00221309.1994.9921213`.