

Causes and effects in Dichotomous Comparative Judgments: an information-theoretical system with plausible mechanism

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

Dichotomous Comparative Judgment (DCJ, [Pollitt \(2012a\)](#), [Pollitt \(2012b\)](#)) requires judges to evaluate the relative manifestation of traits between pairs of stimuli, resulting in a dichotomous outcome indicating which stimulus exhibits the trait more strongly. Research has demonstrated DCJ's effectiveness and reliability in various domains ([Pollitt, 2012b](#); [Bartholomew et al., 2018](#); [van Daal et al., 2019](#); [Lesterhuis, 2018](#); [Bartholomew and Williams, 2020](#); [Boonen et al., 2020](#)). However, the literature lacks a clear and transparent depiction of the plausible mechanisms underlying DCJ data. Specifically, there is no detail explanation of how the different assessment factors can potentially influence the observed DCJ data. This study aims to fill this gap by applying the framework of causal analysis and Directed Acyclic Graphs (DAG; [Pearl \(2009\)](#)). Using this framework, the study will construct a scientific model to elucidate the causal assumptions and mechanisms inherent the system. This model will enable researchers to draw inferences about causal relationships from DCJ data. Subsequently, the study will translate this model into a probabilistic statistical model, aiming to derive statistical estimands for different targets of inference. The outcomes of this study will inform the planning of DCJ experiments and hold significance for researchers or analysts involved in education and assessment procedures who implement the DCJ methodology.

Keywords: comparative judgement, directed acyclic graph, causal analysis, probabilistic statistics

1. Introduction

In contemporary contexts, Thurstone's law of comparative judgment ([1927](#)) primarily refers to the method of *Dichotomous* Comparative Judgment (DCJ, [Pollitt, 2012a,b](#)). In DCJ, a judge assesses the relative manifestation of a *trait* within a pair of stimuli.

*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)

Preprint submitted to *Psychometrika*

July 17, 2024

This assessment results in a dichotomous value indicating which stimulus possesses a higher degree of the trait. After different judges perform multiple rounds of pairwise comparisons, an outcome vector is produced. This vector is modeled using the Bradley-Terry-Luce model (BTL, [Bradley and Terry, 1952](#); [Luce, 1959](#); [Pollitt, 2012a](#)), which produces a latent variable corresponding to the trait of interest. This latent variable is then used to rank the stimuli from lowest to highest or to evaluate the influence of certain variables on the stimuli's positions in the ranking.

DCJ has proven effective in assessing competencies and traits predominantly within the educational realm, as demonstrated by the works of [Pollitt \(2012b\)](#), [Bartholomew et al. \(2018\)](#), [van Daal et al. \(2019\)](#), [Lesterhuis \(2018\)](#), and [Bartholomew and Williams \(2020\)](#). However, its application transcends education, as exemplified by the work of [Boonen et al. \(2020\)](#). The methodology has also evolved to include multiple, as opposed to pairwise comparisons ([Luce, 1959](#); [Plackett, 1975](#)), and to accommodate comparisons with ordinal outcomes ([Tutz, 1986](#); [Agresti, 1992](#)). Overall, research suggests that DCJ offers an alternative and efficient approach to measurement and evaluation, characterized by its reliability and validity ([Lesterhuis, 2018](#); [van Daal, 2020](#); [Marshall et al., 2020](#)). Nevertheless, despite the method's widespread use, the literature lacks a transparent depiction of the DCJ system and the plausible mechanisms that give rise to DCJ data. Particularly, there is no detailed explanation of how different assessment factors can potentially influence the observed DCJ data.

According to [Verhavert et al. \(2019\)](#) and [van Daal \(2020\)](#), several assessment factors interact and contribute to the reliability of the DCJ method. These factors include the number and characteristics of the stimuli, their *proximity* in terms of the assessed trait, the number of comparison per stimulus, and the pairing algorithm used. Furthermore, since the method relies on judges' assessments, the number and characteristics of judges, their *discrimination* abilities, and the number of comparisons per judge also play pivotal roles. Moreover, when the stimuli represent sub-units of higher-levels units, factors such as the number and characteristics of these units, along with their *proximity* in terms of the assessed trait, can significantly influence the outcome. An example of this can be found in [van Daal et al. \(2019\)](#), where the authors assessed the *skills in academic writing* (trait) of Flemish university students, utilizing multiple argumentative essays (stimuli, sub-units) originating from various students (units).

2. Theoretical framework

2.1. Research questions and their estimands

2.2. A scientific model for the DCJ procedure

2.3. From the scientific to the Bradley-Terry-Luce model

3. Discussion

3.1. Limitations and further research

4. Conclusion

Declarations

Funding: The project was founded through the Research Fund of the University of Antwerp (BOF).

Conflict of interests: The authors declare no conflict of interest.

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://jriverspejo.github.io/paper2_manuscript/.

Authors' contributions: *Conceptualization:* S.G., S.D.M., T.vD., and J.M.R.E.; *Methodology:* S.D.M., T.vD., and J.M.R.E.; *Software:* J.M.R.E.; *Validation:* J.M.R.E.; *Formal Analysis:* J.M.R.E.; *Investigation:* J.M.R.E.; *Resources:* S.G., S.D.M., and T.vD.; *Data curation:* J.M.R.E.; *Writing - original draft:* J.M.R.E.; *Writing - review & editing:* S.G., S.D.M., and T.vD.; *Visualization:* J.M.R.E.; *Supervision:* S.G. and S.D.M.; *Project administration:* S.G. and S.D.M.; *Funding acquisition:* S.G. and S.D.M.

References

- Agresti, A., 1992. Analysis of ordinal paired comparison data. *Journal of the Royal Statistical Society* 41, 287–297. URL: <https://www.jstor.org/stable/2347562>, doi:10.2307/2347562.
- 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. 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.
- Bradley, R., Terry, M., 1952. Rank analysis of incomplete block designs: I. the method of paired comparisons. *Biometrika* 39, 324–345. URL: <http://www.jstor.com/stable/2334029>, doi:10.2307/2334029.
- Lesterhuis, M., 2018. The validity of comparative judgement for assessing text quality: An assessor’s perspective. Ph.D. thesis. University of Antwerp.
- Luce, R., 1959. On the possible psychophysical laws. *The Psychological 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.
- Pearl, J., 2009. *Causality: Models, Reasoning and Inference*. Cambridge University Press.
- Plackett, R., 1975. The analysis of permutations. *Journal of the Royal Statistical Society* 24, 193–202. URL: <https://www.jstor.org/stable/2346567>, doi:10.2307/2346567.
- 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.
- Thurstone, L., 1927. A law of comparative judgment. *Psychological Review* 34, 482–499. doi:10.1037/h0070288.
- Tutz, G., 1986. Bradley-terry-luce model with an ordered response. *Journal of Mathematical Psychology* 30, 306–316. doi:10.1016/0022-2496(86)90034-9.
- van Daal, T., 2020. Making a choice is not easy?!: Unravelling the task difficulty of comparative judgement to assess student work. Ph.D. thesis. University of Antwerp.
- van Daal, T., Lesterhuis, M., Coertjens, L., Donche, V., De Maeyer, S., 2019. 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.
- 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.