

HERA: A Hierarchical-Compensatory, Effect-Size Driven and Non-parametric Ranking Algorithm using Data-Driven Thresholds and Bootstrap Validation

02 December 2025

Summary

In scientific disciplines ranging from clinical research to machine learning, researchers face the challenge of objectively comparing multiple algorithms, experimental conditions, or datasets across a variety of performance metrics. This process, often framed as Multi-Criteria Decision Making (MCDM), is critical for identifying state-of-the-art methods. However, traditional ranking approaches frequently suffer from limitations: they may rely on central tendencies that ignore data variability (Demšar 2006; Benavoli et al. 2016), depend solely on p-values which can be misleading in large samples (Wasserstein and Lazar 2016), or require subjective weighting of conflicting metrics (Taherdoost and Madanchian 2023).

HERA (Hierarchical-Compensatory, Effect-Size Driven Ranking Algorithm) is a MATLAB toolbox designed to automate this comparison process, bridging the gap between elementary statistical tests and complex decision-making frameworks. Unlike weighted-sum approaches that collapse multi-dimensional performance into a single scalar, HERA implements a **hierarchical-compensatory logic**. This logic integrates non-parametric significance testing (Wilcoxon signed-rank test), robust effect size estimation (Cliff's Delta, Relative Difference), and bootstrapping (e.g. Percentile and Cluster) to produce rankings that are both statistically robust and practically relevant. HERA is designed for researchers in biomedical imaging, machine learning, and applied statistics who need to compare method performance across multiple quality metrics in a statistically rigorous manner without requiring subjective parameter tuning.

Statement of Need

The scientific community increasingly recognizes the pitfalls of relying on simple summary statistics or p-values alone (Wasserstein and Lazar 2016). In benchmarking studies, specifically, several issues persist:

1. **Ignoring Variance:** Ranking based on mean scores fails to account for the stability of performance across different subjects or folds. A method might achieve a high average score due to exceptional performance on a few easy cases while failing catastrophically on others, yet still outrank a more consistent competitor.
2. **Statistical vs. Practical Significance:** A result can be statistically significant but practically irrelevant, especially in large datasets where even trivial differences yield $p < 0.05$. Standard tests do not inherently distinguish between these cases, potentially leading to the adoption of methods that offer no tangible benefit (Sullivan and Feinn 2012).
3. **Subjectivity in Aggregation:** Many MCDM methods require users to assign arbitrary weights to metrics (e.g., “Accuracy is 0.7, Speed is 0.3”). These weights are often chosen post-hoc or lack empirical justification, introducing researcher bias that can be manipulated to favor a specific outcome (Taherdoost and Madanchian 2023).
4. **Distributional Assumptions:** Parametric tests (e.g., t-test) assume normality, which is often violated in real-world benchmarks where performance metrics may be skewed, bounded, or ordinal (Romano et al. 2006).

HERA addresses these challenges by providing a standardized, data-driven framework. It ensures that a method is only ranked higher if it demonstrates a statistically significant and sufficiently large advantage, preventing “wins” based on negligible differences or noise. Unlike existing MCDM software packages such as the Python libraries `pyDecision` (Pereira et al. 2024) and `pymcdm` (Kizielewicz et al. 2023), or R’s `RMCD` (Najafi and Mirzaei 2025), which often implement classical methods like TOPSIS (Hwang and Yoon 1981), PROMETHEE (Brans and Vincke 1985), and ELECTRE (Roy 1968) that require user-defined weights or preference functions, HERA eliminates subjective parameterization by using data-driven thresholds derived from bootstrap resampling. Furthermore, HERA integrates statistical hypothesis testing directly into the ranking process, a feature absent in standard MCDM toolboxes. While the MATLAB ecosystem offers robust statistical functions, it currently lacks a dedicated, open-source toolbox that unifies this advanced MCDM method with bootstrap validation, forcing researchers to rely on ad-hoc scripts.

Methodological Framework

HERA operates on paired data matrices where rows represent subjects (or datasets) and columns represent the methods to be compared. The core innovation is its sequential logic, which allows for “compensation” between metrics based on strict statistical evidence.

Statistical Rigor and Effect Sizes

HERA quantifies differences using statistical significance and effect sizes to ensure practical relevance independent of sample size (Cohen 1988; Sullivan and Feinn 2012). A “win” always requires satisfying three conjunctive criteria, if not it is considered “neutral”:

- **Significance:** $p < \alpha_{\text{Holm}}$ (Holm-Bonferroni corrected). Pairwise comparisons use the Wilcoxon signed-rank test (Wilcoxon 1945), with p-values corrected using the step-down Holm-Bonferroni method (Holm 1979) to control the Family-Wise Error Rate (FWER).
- **Stochastic Dominance (Cliff’s Delta):** Cliff’s Delta ($d = P(X > Y) - P(Y > X)$) quantifies distribution overlap, is robust to outliers, and relates to common-language effect sizes (Cliff 1993; Vargha and Delaney 2000). The effect size d must exceed a bootstrapped threshold θ_d .
- **Magnitude (Relative Difference):** The Relative Difference (RelDiff) quantifies effect magnitude on the original metric scale, normalized by the mean absolute value. This normalization is formally identical to the Symmetric Mean Absolute Percentage Error (SMAPE) used in forecasting (Makridakis 1993) and conceptually related to the Response Ratio, which uses logarithmic ratios to compare effects across studies (Hedges et al. 1999). The metric enables scale-independent comparisons and facilitates the interpretation of percentage changes (Kampenes et al. 2007). RelDiff must exceed a threshold δ_{RelDiff} .

Dual Criteria & SEM Lower Bound HERA’s complementary logic requires both dominance and magnitude, preventing “wins” based on trivial consistent differences or noisy outliers (Lakens 2013). Thresholds are determined via Percentile Bootstrapping (lower $\alpha/2$ -quantile) (Rousset et al. 2021). To filter noise in low-variance datasets, the RelDiff threshold enforces a lower bound based on the Standard Error of the Mean (SEM), ensuring $\theta_r \geq \theta_{\text{SEM}}$. This approach is inspired by the concept of the “Smallest Worthwhile Change” (Hopkins 2004), but adapted for HERA to quantify the uncertainty of the group mean rather than individual measurement error.

Hierarchical-Compensatory Logic

The ranking process is structured as a multi-stage tournament. It does not use a global score but refines the rank order iteratively (see Fig. 1):

- **Stage 1 (Initial Sort):** Methods are initially ranked based on the win count of the primary metric M_1 . In case of a tie, Cliff’s Delta is used to break the tie.
- **Stage 2 (Compensatory Correction):** This stage addresses the trade-off between metrics. A lower-ranked method can “swap” places with a higher-ranked method if it shows a statistically significant and relevant superiority in a secondary metric M_2 . This effectively implements a lexicographic ordering with a compensatory component (Keeney and Raiffa

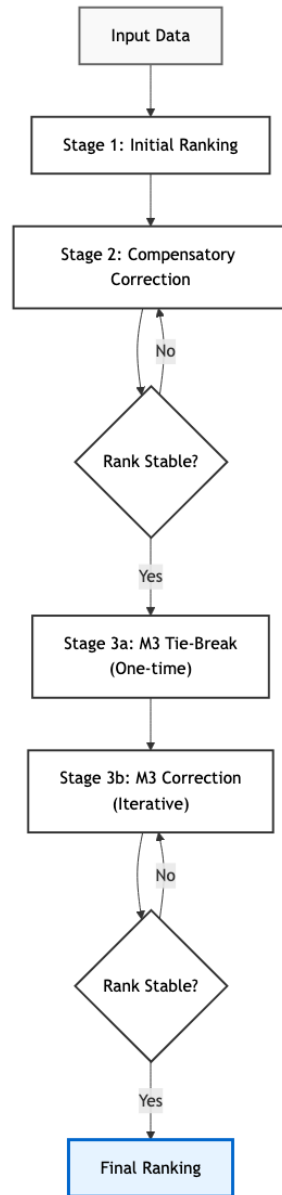


Figure 1: Hierarchical-Compensatory Ranking Logic

1976), allowing a method that is slightly worse in the primary metric but vastly superior in a secondary metric to improve its standing.

- **Stage 3 (Tie-Breaking):** This stage resolves “neutral” results using a tertiary metric M_3 . It applies two sub-logics to ensure a total ordering:
 - **Sublogic 3a:** A one-time correction if the previous metric is “neutral” based on the HERA criteria. This handles cases where two methods are indistinguishable in the second metric while still respecting the initial ranking.
 - **Sublogic 3b:** To resolve groups of remaining undecided methods, an iterative correction loop is applied if both M_1 and M_2 are “neutral”, iteratively using metric M_3 until a final stable ranking is found.

Validation and Uncertainty

HERA integrates advanced resampling methods to quantify uncertainty:

- **BCa Confidence Intervals:** Bias-Corrected and Accelerated (BCa) intervals are calculated for all effect sizes (DiCiccio and Efron 1996).
- **Cluster Bootstrap:** To assess the stability of the final ranking, HERA performs a cluster bootstrap resampling subjects with replacement (Field and Welsh 2007). This yields a 95% confidence interval for the rank of each method.
- **Power Analysis:** A post-hoc simulation with bootstrap estimates the probability of detecting a “win”, “loss” or “neutral” in all tested metrics given the data characteristics.
- **Sensitivity Analysis:** The algorithm permutes the metric hierarchy and aggregates the resulting rankings using a Borda Count (Young 1974) to evaluate the robustness of the decision against hierarchy changes.

Software Features

HERA offers a flexible configuration of up to three metrics (see Fig. 2). This allows users to adapt the ranking logic to different study designs and needs. It also provides a range of reporting options, data integration, and reproducibility features.

- **Automated Reporting:** Generates PDF reports, Win-Loss Matrices, Sankey Diagrams, and machine-readable JSON/CSV exports.
- **Reproducibility:** Supports fixed-seed execution and configuration file-based workflows. The full analysis state, including random seeds and parameter settings, is saved in a JSON file, allowing other researchers to exactly replicate the ranking results.
- **Convergence Analysis:** To avoid the common pitfall of using an arbitrary number of bootstrap iterations, HERA implements an adaptive algorithm. It automatically monitors the stability of the estimated confidence intervals and effect size thresholds, continuing the resampling process until the estimates converge within a specified tolerance, thus determining the

optimal number of iterations B dynamically (Pattengale et al. 2010). If the characteristics of the data for bootstrapping are known, the number of bootstrap iterations can be set manually.

- **Data Integration:** HERA supports seamless data import from standard formats (CSV, Excel) and MATLAB tables, facilitating integration into existing research pipelines. Example datasets and workflows demonstrating practical applications are included in the repository.
- **Accessibility:** HERA can be easily installed by cloning the GitHub repository and running a setup script, or deployed as a standalone application that requires no MATLAB license. An interactive command-line interface guides users through the analysis without requiring programming expertise, while an API and JSON Configuration allow for automated batch processing.

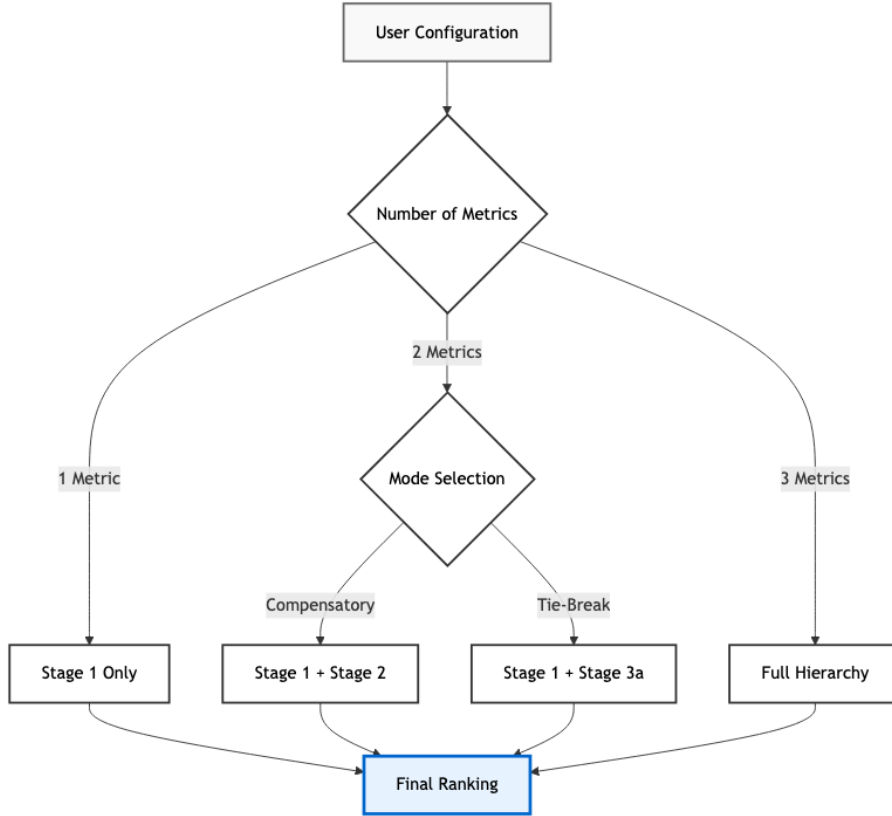


Figure 2: Flexible Configuration options for Ranking Logic

Acknowledgements

This software was developed at the Institute of Neuroradiology, Goethe University Frankfurt. I thank Prof. Dr. Dipl.-Phys. Ralf Deichmann (Cooperative Brain Imaging Center, Goethe University Frankfurt) for his support during the initial conceptualization of this project. I acknowledge Dr. med. Christophe Arendt (Institute of Neuroradiology, Goethe University Frankfurt) for his supervision and support throughout the project. I also thank Rejane Golbach PhD (Institute of Biostatistics and Mathematical Modeling, Goethe University Frankfurt) for her valuable feedback on the statistical methodology.

References

- Benavoli, A., G. Corani, and F. Mangili. 2016. “Should We Really Use Post-Hoc Tests Based on Mean-Ranks?” *Journal of Machine Learning Research* 17: 1–10. <https://jmlr.org/papers/v17/benavoli16a.html>.
- Brans, J. P., and P. Vincke. 1985. “A Preference Ranking Organization Method (the PROMETHEE Method for Multiple Criteria Decision-Making).” *Management Science* 31 (6): 647–56. <https://doi.org/10.1287/mnsc.31.6.647>.
- Cliff, N. 1993. “Dominance Statistics: Ordinal Analyses to Answer Ordinal Questions.” *Psychological Bulletin* 114 (3): 494–509. <https://doi.org/10.1037/0033-2909.114.3.494>.
- Cohen, J. 1988. *Statistical Power Analysis for the Behavioral Sciences*. 2nd ed. Lawrence Erlbaum Associates. <https://doi.org/10.4324/9780203771587>.
- Demšar, J. 2006. “Statistical Comparisons of Classifiers over Multiple Data Sets.” *Journal of Machine Learning Research* 7: 1–30. <https://jmlr.org/papers/v7/demsar06a.html>.
- DiCiccio, T. J., and B. Efron. 1996. “Bootstrap Confidence Intervals.” *Statistical Science* 11 (3): 189–228. <https://doi.org/10.1214/ss/1032280214>.
- Field, C. A., and A. H. Welsh. 2007. “Bootstrapping Clustered Data.” *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 69 (3): 369–90. <https://doi.org/10.1111/j.1467-9868.2007.00593.x>.
- Hedges, L. V., J. Gurevitch, and P. S. Curtis. 1999. “The Meta-Analysis of Response Ratios in Experimental Ecology.” *Ecology* 80 (4): 1150–56. [https://doi.org/10.1890/0012-9658\(1999\)080%5B1150:TMAORR%5D2.0.CO;2](https://doi.org/10.1890/0012-9658(1999)080%5B1150:TMAORR%5D2.0.CO;2).
- Holm, S. 1979. “A Simple Sequentially Rejective Multiple Test Procedure.” *Scandinavian Journal of Statistics* 6 (2): 65–70. <https://doi.org/10.2307/4615733>.

- Hopkins, W. G. 2004. "How to Interpret Changes in an Athletic Performance Test." *Sportscience* 8: 1–7. <https://www.sportsci.org/jour/04/wghtests.htm>.
- Hwang, C. L., and K. Yoon. 1981. *Multiple Attribute Decision Making: Methods and Applications*. Springer. <https://doi.org/10.1007/978-3-642-48318-9>.
- Kampenes, V. B., T. Dybå, J. E. Hannay, and D. I. K. Sjøberg. 2007. "A Systematic Review of Effect Size in Software Engineering Experiments." *Information and Software Technology* 49 (11–12): 1073–86. <https://doi.org/10.1016/j.infsof.2007.02.015>.
- Keeney, R. L., and H. Raiffa. 1976. *Decisions with Multiple Objectives: Preferences and Value Trade-Offs*. Wiley. <https://doi.org/10.1017/CBO9781139174084>.
- Kizielewicz, B., A. Shekhovtsov, and W. Salabun. 2023. "Pymcdm—the Universal Library for Solving Multi-Criteria Decision-Making Problems." *SoftwareX* 22: 101368. <https://doi.org/10.1016/j.softx.2023.101368>.
- Lakens, D. 2013. "Calculating and Reporting Effect Sizes to Facilitate Cumulative Science: A Practical Primer for t-Tests and ANOVAs." *Frontiers in Psychology* 4: 863. <https://doi.org/10.3389/fpsyg.2013.00863>.
- Makridakis, S. 1993. "Accuracy Measures: Theoretical and Practical Concerns." *International Journal of Forecasting* 9 (4): 527–29. [https://doi.org/10.1016/0169-2070\(93\)90079-3](https://doi.org/10.1016/0169-2070(93)90079-3).
- Najafi, A., and S. Mirzaei. 2025. "RMCDA: The Comprehensive r Library for Applying Multi-Criteria Decision Analysis Methods." *Software Impacts* 24: 100762. <https://doi.org/10.1016/j.simpa.2025.100762>.
- Pattengale, N. D., M. Alipour, O. R. P. Bininda-Emonds, B. M. E. Moret, and A. Stamatakis. 2010. "How Many Bootstrap Replicates Are Necessary?" *Journal of Computational Biology* 17 (3): 337–54. <https://doi.org/10.1089/cmb.2009.0179>.
- Pereira, V., M. P. Basilio, and C. H. T. Santos. 2024. "Enhancing Decision Analysis with a Large Language Model: pyDecision a Comprehensive Library of MCDA Methods in Python." *Journal of Modelling in Management*, ahead of print. <https://doi.org/10.1108/JM2-04-2024-0118>.
- Romano, J., J. D. Kromrey, J. Coraggio, and J. Skowronek. 2006. "Appropriate Statistics for Ordinal Level Data: Should We Really Be Using t-Test and Cohen's d for Evaluating Group Differences on the NSSE and Other Surveys?" *Proceedings of the Annual Meeting of the Florida Association of Institutional Research* (Cocoa Beach, FL). <https://www.researchgate.net/publication/237>

544991.

- Rousselet, G. A., C. R. Pernet, and R. R. Wilcox. 2021. “The Percentile Bootstrap: A Primer with Step-by-Step Instructions in R.” *Advances in Methods and Practices in Psychological Science* 4 (1): 1–10. <https://doi.org/10.1177/2515245920911881>.
- Roy, B. 1968. “Classement Et Choix En Présence de Points de Vue Multiples (La méthode ELECTRE).” *Revue Française d’informatique Et de Recherche Opérationnelle* 2 (V1): 57–75. <https://doi.org/10.1051/ro/196802V100571>.
- Sullivan, G. M., and R. Feinn. 2012. “Using Effect Size—or Why the p Value Is Not Enough.” *Journal of Graduate Medical Education* 4 (3): 279–82. <https://doi.org/10.4300/JGME-D-12-00156.1>.
- Taherdoost, H., and M. Madanchian. 2023. “Multi-Criteria Decision Making (MCDM) Methods and Concepts.” *Encyclopedia* 3 (1): 235–50. <https://doi.org/10.3390/encyclopedia3010006>.
- Vargha, A., and H. D. Delaney. 2000. “A Critique and Improvement of the CL Common Language Effect Size Statistics of McGraw and Wong.” *Journal of Educational and Behavioral Statistics* 25 (2): 101–32. <https://doi.org/10.3102/10769986025002101>.
- Wasserstein, R. L., and N. A. Lazar. 2016. “The ASA Statement on p-Values: Context, Process, and Purpose.” *The American Statistician* 70 (2): 129–33. <https://doi.org/10.1080/00031305.2016.1154108>.
- Wilcoxon, F. 1945. “Individual Comparisons by Ranking Methods.” *Biometrics Bulletin* 1 (6): 80–83. <https://doi.org/10.2307/3001968>.
- Young, H. P. 1974. “An Axiomatization of Borda’s Rule.” *Journal of Economic Theory* 9 (1): 43–52. [https://doi.org/10.1016/0022-0531\(74\)90073-8](https://doi.org/10.1016/0022-0531(74)90073-8).