Are We Ranking Fairly? A Bias-Aware Bayesian Rank Aggregation Method for Peer- and Self-Evaluations

Introduction

- Rank aggregation (RA), which aims to combine multiple ordinal rankings into a single ranking, has broad applications in fields such as elections [1], genomic research [2], and educational evaluation [3].
- In peer-and self-assessments, subjective biases can lead to unreliable or unfair outcomes, and existing methods often overlook these unique data mechanisms.
- We propose a novel, bias-aware, Bayesian method **BayeSRank** (Bayesian Bias Detection in Peer-and Self-Ranking), validated through simulation studies and real-world data.



Figure: Illustration of Peer- and selfevaluation in a classroom setting. Generated by Adobe illustrator.

- We also invite you to consider these related interesting questions:
 - Should self-evaluation data be removed to mitigate self-bias?
 - □ Is self-bias indicative of a ranker's overall competence?
 - □ How will self-bias influence the aggregated results?

Data and Latent Variable Model

	Ranker 1	Ranker 2	Ranker 3	Ranker 4	Ranker 5	
Presenter 1	<i>r</i> ₁₁	<i>r</i> ₁₂	<i>r</i> ₁₃	<i>r</i> ₁₄	r ₁₅	•••
Presenter 2	<i>r</i> ₂₁	<i>r</i> ₂₂	<i>r</i> ₂₃	r ₂₄	r ₂₅	
Presenter 3	<i>r</i> ₃₁	r ₃₂	r ₃₃	r ₃₄	r ₃₅	•••
Presenter 4	<i>r</i> ₄₁	<i>r</i> ₄₂	<i>r</i> ₄₃	r_{44}	r ₄₅	•••
Presenter 5	r_{51}	r ₅₂	r ₅₃	r ₅₄	r_{55}	•••
:	:	:	:	:	:	•.

Table: A sample ranking matrix **R** of dimension $n \times n$, where the entry r_{ii} represents the rank given by Ranker *j* to Presenter *i*. The diagonal entries (r_{ii}) correspond to self-ranks.

- Assume a latent importance score, ω_{ij} , which is associated with the observed ranks r_{ii} of each presenter-ranker pair.
- The unknown true performance, μ_i , is estimated via a probabilistic model that establishes a linear relationship between ω_{ii} and μ_i :

$$\omega_{ij} = \mu_i + I(i = j)\beta_j + \epsilon_{ij}$$
, where

- $i \in \{1, ..., n\}$ index presenters being ranked;
- $j \in \{1, \dots, J\}$ index rankers (*note that J = n);
- ω_{ij} : Latent variable of each observed rank;
- μ_i : True academic performance of presenter *i*;
- β_i : Self-evaluation bias of ranker *j* with mean μ_B
- and variance σ_{B}^{2} ;
- ϵ_{ii} : Random error in the evaluation;
- $\sigma_{\epsilon,i}^2$: The overall variance of each presenter's ranking.

Haoyao Ruan¹, Kaiwen (Kevin) Wang², Xinlei (Sherry) Wang¹

Bayesian Hierarchical Model



Figure: Comparison of BayeSRank and BayeRank (model without selfevaluation term, β_i) performance across different data-simulation settings. From top to bottom, bar plots display the (1) Spearman correlation (in %), (2) Top-1 coverage rate (in %), and (3) Top-3 coverage rate (in %) for each method across varying values of **n** (number of presenters), **p** (correlation between latent variable and true performance, which indicates the ranking <u>quality of rankers), and μ_{β} (mean self-evaluation bias effect).</u> Each bar represents the mean estimate with 95% confidence intervals.





1. Department of Mathematics, University of Texas at Arlington, Arlington, Texas 2. Department of Statistical Science, Southern Methodist University, Dallas, Texas

• **Gibb Sampling [4]** is implemented to generate the posterior samples given the known derived **conditional distributions**:

$$\succ \mu_i \sim N\left(\frac{\sum_{j=1}^{J} \frac{\omega_{ij} - I(i=j)\beta_j}{\sigma_{\epsilon_{.j}}^2}}{1 + \sum_{j=1}^{J} \frac{1}{\sigma_{\epsilon_{.j}}^2}}, \frac{1}{1 + \sum_{j=1}^{J} \frac{1}{\sigma_{\epsilon_{.j}}^2}}\right), i = 1, \dots, n;$$

$$> \beta_{j} \sim N\left(\frac{\frac{\omega_{jj} - \mu_{j}}{\sigma_{\epsilon_{.j}}^{2}} + \frac{\mu_{\beta}}{\sigma_{\beta}^{2}}}{\frac{1}{\sigma_{\epsilon_{.j}}^{2}} + \frac{1}{\sigma_{\beta}^{2}}}, \frac{1}{\frac{1}{\sigma_{\epsilon_{.j}}^{2}} + \frac{1}{\sigma_{\beta}^{2}}}\right) \times I(-100, 100), \ j = 1, \cdots, J;$$

$$\succ \mu_{\beta} \sim N\left(\frac{\sum_{j=1}^{J} \beta_{j}}{J}, \frac{\sigma_{\beta}^{2}}{J}\right) \times I(-100, 100);$$

$$\succ \sigma_{\beta}^{2} \sim IG\left(\delta_{\beta} + \frac{J}{2}, \delta_{\beta} + \frac{1}{2}\sum_{j=1}^{J}(\beta_{j} - \mu_{\beta})^{2}\right);$$

$$\succeq \sigma_{\beta}^{2} \sim IC\left(\delta_{\beta} + \frac{J}{2}, \delta_{\beta} + \frac{1}{2}\sum_{j=1}^{n}(\alpha_{j} - \mu_{\beta})^{2}\right);$$

- $\succ \sigma_{\epsilon_j}^2 \sim IG\left(\delta_{\epsilon} + \frac{J}{2}, \, \delta_{\epsilon} + \frac{1}{2}\sum_{i=1}^n (\omega_{ij} \mu_i I(i=j)\beta_j)^2\right),$ $j = 1, \cdots, J;$
- $\succ \omega_{ij} \sim TN(\mu_i + I(i=j)\beta_j, \sigma_\beta^2; \omega_{(k-1,j)}, \omega_{(k+1,j)}), i = 1,$ \cdots , n; $j = 1, \cdots, J$.

Simulation Study Results



Figure: Spearman correlation results of BayeSRank vs. other methods [5-10] under varying (n, ρ and μ_{β}) simulation settings. * .selfrm indicates removal of self-valuation data. Within each column, color indicates the relative performance of each method (warmer = better, cooler = worse), while size reflects the magnitude of the average Spearman correlation between estimated and true rankings. BayeSRank consistently ranks among the top-performing methods.

Figure: Top 3 coverage rate (in %) with 95% confidence interval across varying (n, p and µ_B) simulation settings. Note that the x-axes are not linearly scaled for n and µ_B, which may visually exaggerate or flatten trends. BayeSRank (in black) consistently outperforms or matches other RA methods.



Conclusions

- Study results suggest that removal of selfevaluation data is not the optimal approach, as self-evaluations contribute to ranking accuracy when properly adjusted.
- Through simulations and real-world data examples, we demonstrate BayeSRank's superiority in generating interpretable, **unbiased** evaluations, especially in noisy data settings.
- Our work enhances fairness, transparency, and reliability in peer and self-evaluation systems, offering theoretical and practical implications for bias-aware ranking.

References

- Borda, JC de. "M'emoire sur les' elections au scrutin." Histoire de l'Acad'emie Royale des Sciences (1781). 2. Li, Xue, et al. "A Bayesian latent variable approach to aggregation of partial and top-ranked lists in genomic studies." Statistics in medicine 37.28 (2018): 4266-4278. 3. Gramzow, Richard H., et al. "Self-evaluation bias and academic performance: Some ways and some reasons
- why." Journal of Research in Personality 37.2 (2003): 41-61. 4. Geman, Stuart, and Donald Geman. "Stochastic relaxation, Gibbs distributions, and the Bayesian restoration
- of images." IEEE Transactions on pattern analysis and machine intelligence 6 (1984): 721-741. 5. Badgeley, Marcus A., Stuart C. Sealfon, and Maria D. Chikina. "Hybrid Bayesian-rank integration approach improves the predictive power of genomic dataset aggregation." Bioinformatics 31.2 (2015): 209-215.
- 6. Dwork, Cynthia, et al. "Rank aggregation methods for the web." Proceedings of the 10th international conference on World Wide Web. 2001.
- 7. Lin, Shili, and Jie Ding. "Integration of ranked lists via cross entropy Monte Carlo with applications to mRNA and microRNA studies." Biometrics 65.1 (2009): 9-18. 8. Kolde, Raivo, et al. "Robust rank aggregation for gene list integration and meta-analysis." Bioinformatics 28.4 (2012): 573-580.
- R packages . Schimek MG, Budinska E, Ding J, Kugler KG, Svendova V, Lin S, Pfeifer B (2022). _TopKLists: Inference, Aggregation and Visualization for Top-K Ranked Lists_. R package version 1.0.8, <https://CRAN.R-project.org/package=TopKLists>. 10. Kolde R (2022). _RobustRankAggreg: Methods for Robust Rank Aggregation_. R package version 1.2.1, < https://CRAN.Rproject.org/package=RobustRankAggreg>.