

Inverse Preference Inference via Approximate Bayesian Computation and Counterfactual Simulation for Rule Evaluation in Ranking Systems

Chenyu Yang*, Xiang Li

School of Materials and New Energy, South China Normal University, Shanwei, China

* Corresponding Author Email: 20238032004@m.scnu.edu.cn

Abstract. This paper addresses the issue of unobservable audience voting data in real-world competition scenarios by proposing a unified modeling framework based on inverse inference. Treating judges' scores and elimination outcomes as observable constraints, the study constructs a Bayesian inverse problem model incorporating ranking rules by introducing a social attention prior. An approximate Bayesian computation method is employed to achieve stable estimation of the underlying voting distribution. Building upon this foundation, the paper further proposes statistical metrics to characterize result anomalies and validates the model's reliability under sparse information conditions using multi-season data. Centered on the inference results, the study establishes a counterfactual simulation and fairness evaluation system, systematically comparing the structural differences in fairness and engagement across various scoring mechanisms. Furthermore, the paper designs an automated competition mechanism requiring no subjective intervention. By employing algorithmic immunity strategies, it significantly enhances fairness while maintaining commercial activity. This method operates independently of domain-specific assumptions, making it applicable to diverse complex systems involving implicit preferences and ranking decisions. It provides a universally applicable, data-driven approach for rule design and mechanism optimization.

Keywords: Inverse Inference; Approximate Bayesian Computation; Counterfactual Simulation.

1. Introduction

In multi-stakeholder ranking and elimination systems, final outcomes are typically determined by both explicit evaluations and implicit preferences. Since certain decision factors are difficult to observe directly, results often exhibit phenomena that do not fully align with surface-level performance. Existing research often relies on observable data for direct modeling or approximates latent factors through simplifying assumptions. While effective under simple rules or abundant data, these approaches often fail to accurately reconstruct genuine decision-making mechanisms or support systematic evaluations of rules when confronted with phased eliminations, nonlinear rankings, or institutional constraints.

This paper adopts an algorithmic modeling perspective, treating unobservable group preferences as latent variables. It introduces a reverse inference approach to construct a unified inference framework constrained by outcomes. This framework utilizes known scoring and elimination information as constraints. Through a simulation-driven inference mechanism, it reconstructs hidden decision structures without requiring explicit preference data. Furthermore, the paper applies these inferences to counterfactual simulations and rule comparison analyses, enabling the model not only to explain historical outcomes but also to evaluate and design ranking mechanisms.

Using competitive systems as a concrete application example, the study demonstrates the algorithmic framework's capability to reveal conflicts between scoring rules and collective preferences. It proposes a rule improvement scheme that requires no subjective intervention and possesses generalizability, offering a transferable modeling pathway for mechanism analysis in complex ranking systems.

2. Inverse Modeling Framework for Latent Preference Inference

The primary objective of this study is to infer historically unobservable audience voting data (fan votes) and construct a reliable quantitative metric. This problem constitutes a classical inverse inference task, whose fundamental components are defined as follows:

Observed Space O : consists of the established judge scoring matrix S_{judge} and the historical elimination outcomes of the lowest-ranked contestants for each episode E_{obs} .

Latent Variable L : represents the unknown vector of audience vote shares V_{fan} .

Constraint Mechanism: the weighted aggregation of S_{judge} and V_{fan} , when processed by the ranking operator, must strictly reproduce the observed elimination sequence E_{obs} .

This formulation establishes the inference of V_{fan} as a constrained inverse problem.

2.1. Data Augmentation and Prior Design

To compensate for the absence of explicit audience voting data, a proxy variable is introduced based on Long Tail Theory. Contestants are assigned base popularity priors according to social media influence levels, reflecting a long-tail distribution consistent with Zipf's Law.

The logarithmic distribution of popularity exhibits a near-linear trend, indicating that fan voting behavior conforms to Zipf's Law. This observation supports the construction of a two-dimensional feature space composed of S_{judge} and V_{base} , which provides a mathematically coherent foundation for defining the Bayesian prior.

2.2. Bayesian Framework

The inversion task is reformulated as estimating the posterior distribution of V_{fan} . According to Bayes' theorem:

$$P(V_{\text{fan}} | E_{\text{obs}}, S_{\text{judge}}) \propto P(E_{\text{obs}} | V_{\text{fan}}, S_{\text{judge}}) \cdot P(V_{\text{fan}}) \quad (1)$$

Where the prior reflects latent popularity, the likelihood measures consistency with historical elimination, and the posterior represents the inferred voting distribution.

Due to non-linear ranking rules, the likelihood admits no closed-form solution. Therefore, an Approximate Bayesian Computation (ABC) framework is employed.

2.3. Feature Engineering: Standardisation of Judge Shares

To ensure cross-season comparability under varying scoring scales and panel sizes, judge scores are transformed into a standardised zero-sum share representation, enabling direct fusion with voting shares. This transformation maps heterogeneous raw scores into a unified proportional space, thereby ensuring consistency in subsequent simulation and inference.

2.4. Model Construction: ABC Rejection Sampling Network

An ABC-based Bayesian Rejection Sampling Network is constructed, comprising: prior construction, simulation engine, and rejection sampling [1].

A Dirichlet prior is adopted [2]:

$$V_{\text{prior}} \sim \text{Dir}(\alpha), \quad \alpha = \Phi(\text{SocialMediaFans}). \quad (2)$$

Where $\Phi(\cdot)$ maps social media follower counts to the concentration parameters of the Dirichlet distribution. This formulation encodes the assumption that off-stage popularity determines latent voting potential.

A simulation engine $f(\cdot)$ reproduces the historical ranking process. The likelihood constraint is implemented via an acceptance indicator:

$$Accept(V_{sim}) \Leftrightarrow \mathbb{I}\left(\text{Eliminated}\left(f(V_{sim}, S_{judge})\right) = E_{obs}\right) = 1 \quad (3)$$

Only samples that exactly reproduce the historical elimination sequence are retained. The resulting accepted set approximates the posterior distribution of V_{fan} .

3. Posterior Inference, Robustness Analysis, and Anomaly Quantification

After $N = 50000$ Monte Carlo simulations, a set of accepted samples $V_{k=1}^{(k)M}$ is obtained.

Point estimate. The posterior mean of the accepted samples is used as the point estimate of the latent audience vote share V_{fan} . **Uncertainty.** The standard deviation of the accepted samples, denoted by σ_{fan} , is computed to quantify posterior dispersion.

3.1. Surprise Index

To quantify the abnormality of competition outcomes, a Surprise Index S is defined as:

$$S = -\log(\gamma) = -\log\left(\frac{M}{N}\right) \quad (4)$$

Where γ denotes the acceptance rate in rejection sampling. As $S \rightarrow \infty$, the observed outcome corresponds to an extremely low-probability event, characterizing a statistical “black swan”.

3.2. Model Solution Outcomes and Verification

Complete inversion was conducted on 2748 samples across 34 seasons, resulting in the high-precision `inferred_fan_votes` dataset. Posterior uncertainty remains highly concentrated across seasons, with a global mean standard deviation of $\sigma_{mean} = 0.0471$. This indicates that most inferred vote shares lie within a narrow band around their posterior means, supporting stable inference from sparse elimination signals.

3.3. Macroscopic Pattern Analysis

A binary dichotomy between low technical scores and high popularity is observed. All contestants’ S_{judge} and V_{fan} are visualized in a unified scatter plot, as shown in the Fig. 1.

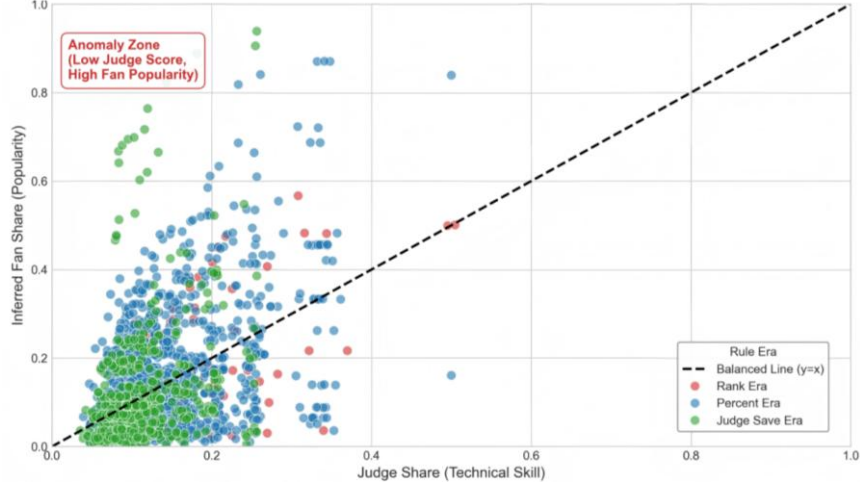


Figure 1. Scatter plot of judge share and estimated audience share

The upper-left region above the $y = x$ diagonal defines the anomaly zone, where low judge scores coincide with high inferred audience votes, indicating systematic deviations between technical merit and popularity.

4. Counterfactual Simulation-Based Fairness Evaluation

This section adopts the inferred audience voting distribution V_{fan} as the benchmark and develops a counterfactual simulation framework for quantitatively assessing the impact of different scoring rules on fairness. Fairness is defined as the alignment between the final ranking (Actual Rank) and the merit-based ranking (Merit-based Rank).

4.1. Model Construction and Assumptions

Merit benchmark: The merit-based ranking is defined as the ranking derived from judge scores: $Rank_{Merit} \equiv Rank(S_{judge})$. This ranking serves as the exclusive objective reference for evaluating fairness after removing popularity bias.

Rational judge assumption. During the Judge Save phase, judges are assumed to act strictly according to technical merit, retaining the contestant with the higher technical score among the Bottom 2: $Saved = \arg \max(S_{judge})$.

4.2. Probability Distribution Reconstruction

Parametric bootstrapping is employed to propagate posterior uncertainty into the counterfactual simulations. Since vote shares are confined to the interval $[0, 1]$, assuming a normal distribution introduces boundary bias. Therefore, the Method of Moments is applied to map the posterior mean μ and variance σ^2 to the parameters of a Beta distribution:

$$\alpha = \mu \left(\frac{\mu(1-\mu)}{\sigma^2} - 1 \right), \quad \beta = (1-\mu) \left(\frac{\mu(1-\mu)}{\sigma^2} - 1 \right) \quad (5)$$

At each Monte Carlo iteration, $V_{sim} \sim \text{Beta}(\alpha, \beta)$, which ensures the mathematical validity of the simulated voting data [3, 4].

4.3. Rule Engine

Three historical aggregation rules are defined as follows:

Rank Era (Ranking Points System).

$$\text{Score} = R(S_{judge}) + R(V_{fan}) \quad (6)$$

Where $R(\cdot)$ denotes the descending ranking operator, and lower scores indicate higher ranks.

Percent Era (Percentage Weighting System).

$$\text{Score} = 0.5 \cdot S_{judge} + 0.5 \cdot V_{fan} \quad (7)$$

Judge Save Era. Based on Percent Era ranking, if a contestant enters the Bottom 2, the contestant with the lowest judge score is eliminated:

$$\text{Eliminated} = \arg \min_{i \in \text{Bottom 2}} (S_{judge}^{(i)}) \quad (8)$$

4.4. Normalised Fairness Metric

To eliminate the influence of varying participant counts N across seasons, the normalised Kendall inverse rank metric is adopted:

$$F = 1 - \frac{2K}{N(N-1)}, \quad F \in [0,1] \quad (9)$$

Where K denotes the number of reversed pairs, defined as cases where the Actual Rank is inferior to the Merit-based Rank. As F approaches 1, the ranking system becomes increasingly fair.

4.5. Results and Inference

A total of 5000 Monte Carlo simulations were conducted across 34 seasons. The principal statistical outcomes are summarised in Table 1.

Table 1. Fairness across aggregation rules

Aggregation Rule	mean	std	stability
Rank Era	0.8576	0.0623	High
Judge Save Era	0.6817	0.1276	Medium
Percent Era	0.6592	0.1345	Low

The results indicate that the Ranking Points System exhibits the highest fairness and stability, while the Percentage Weighting System is the most susceptible to popularity bias.

4.6. Statistical Significance: The Effective Limits of Judicial Authority to Overrule Popularity

Welch's t-test was employed to assess the statistical significance of fairness differences across scoring regimes.

Structural advantage:

The p-value for Rank Era outperforming Percent Era approaches zero: $p < 10^{-80}$.

Local correction effect. Judge Save Era exhibits a statistically significant improvement over Percent Era: $p = 0.026 < 0.05$.

However, the effect size remains marginal: $\Delta\mu \approx 0.02$.

This indicates that Judge Save is constrained by its activation conditions (only Bottom 2) and is therefore incapable of correcting the systemic imbalance caused by global weighting distortion in the percentage system.

5. Attribution Modeling and Algorithmic Mechanism Design

This section explains the drivers of audience voting behavior through hierarchical statistical modeling and evaluates the system-level trade-off between meritocracy and commercial engagement.

5.1. Attribution Analysis with GLMMs

Given the hierarchical structure of DWTS data—where observations are nested within contestants, contestants within seasons, and the response variable V_{fan} is a proportional variable in the interval $[0, 1]$ —a generalised linear mixed-effects model (GLMM) is constructed with interaction terms [5, 6].

A logit link function is applied to map the target variable onto the real domain:

$$\ln\left(\frac{V_{ijt}}{1-V_{ijt}}\right) = \beta_0 + \underbrace{\beta_1 S_{ijt}}_{\text{Skill}} + \underbrace{\beta_2 X_{demo}}_{\text{Demographics}} + \underbrace{\beta_3 (S_{ijt} \times \text{Era}_j)}_{\text{Interaction}} + \underbrace{u_i + u_j}_{\text{Random Effects}} + \delta_{ijt}. \quad (10)$$

Where S_{ijt} denotes standardised judge score (technical merit), X_{demo} denotes demographic covariates, and $S_{ijt} \times \text{Era}_j$ captures structural changes in the merit-to-vote conversion rate across formats. The terms u_i and u_j represent contestant-level and season-level random intercepts, respectively. The contestant-level random effect is assumed to follow $u_i \sim N(0, \sigma^2)$, which enables partial pooling and mitigates sample imbalance.

5.2. Model Estimation and Selection

Restricted maximum likelihood (REML) is used for parameter estimation. Model comparison yields: Base Model AIC: 4632.5; Full Model AIC: 4624.1 ($\Delta AIC = -8.4$).

According to standard criteria, $\Delta AIC < -2$ indicates meaningful improvement. The likelihood ratio test further confirms statistical significance $p < 0.01$. These results imply that the interaction term judge score \times era is both necessary and statistically justified, suggesting that audience voting rationale has undergone structural transformation over time.

5.3. Model Efficacy and Key Findings

The conditional pseudo- R^2 is $R^2_{\text{conditional}} = 0.868$, indicating that the model explains 86.8% of the variance in audience voting.

Fig. 2 illustrates substantial disparity in audience popularity across professions. Musicians exhibit strong positive effects ($\beta \approx +1.08^{***}$), while athletes show significant negative effects ($\beta \approx -1.25^{***}$), indicating that athletes are structurally disadvantaged within the competition.

Evolution of audience rationality: The interaction effect for $S_{judge} \times \text{RankEra}$ is significantly negative ($\beta = -0.17^{**}$).

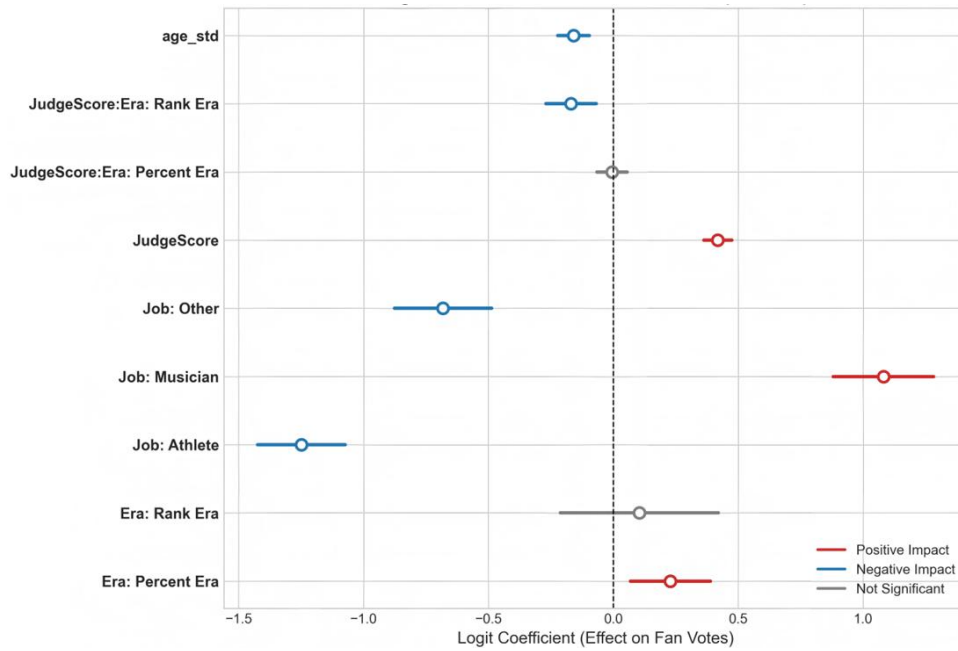


Figure 2. Forest plot of fixed effects

The interaction results indicate an increased sensitivity of audience voting to technical merit in later competition formats.

Age and personal traits. The age coefficient is significantly negative ($\beta = -0.16^{***}$), confirming a systematic audience preference for younger contestants. The random-effect variance is estimated as $\sigma_{\text{group}}^2 = 0.43$, which substantially exceeds the residual variance. This implies that even after controlling for professional background and technical skill, individual charisma remains a dominant determinant of audience voting behavior.

5.4. The “Merit–Safety” Protocol: Multi-Objective Optimization

A multi-objective utility framework is constructed to balance competitive meritocracy with commercial engagement:

Meritocracy metric F : based on the normalised Kendall inverse rank metric, measuring consistency between outcomes and true skill levels.

Engagement metric E : defined as total fan base retention, computed as the aggregate number of social media followers (in millions) of weekly surviving contestants, serving as a proxy for potential viewership.

5.5. Pareto Analysis and Design Implications

Simulation outcomes over 34 seasons (Fig. 3) reveal efficiency variations across mechanisms.

Data insights: Judge Save is located below the Pareto frontier (black dotted line), indicating that it is dominated by alternative solutions. Its engagement volume is only 11120 (approaching 0 after normalization), demonstrating an inability to retain high-traffic participants while maintaining fairness.

Binary dilemma: In the absence of Judge Save, producers face a difficult choice between the Rank System and the Percent System. This binary trade-off indicates that existing mechanisms fail to jointly optimize merit and commercial value, motivating the need for novel rule designs.

Commercial value assumption. The aggregate number of social media followers is assumed to be positively and approximately linearly correlated with a program’s potential viewership. Accordingly, this quantity is adopted as a direct proxy variable for assessing commercial value.

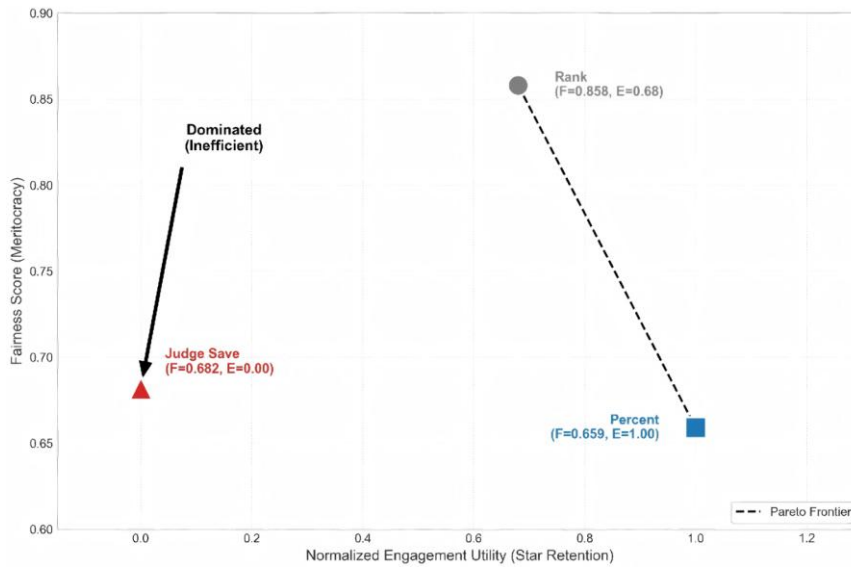


Figure 3. Analysis of the Pareto frontier of the existing tournament format

5.6. Proposal: Merit–Safety Protocol

A new hybrid competition mechanism, termed the Merit–Safety Protocol, is proposed to improve meritocracy while preserving the commercial advantages of the Percentage Weighting System. Unlike Judge Save, the protocol operates automatically and does not rely on subjective judicial intervention.

5.7. Algorithm Description

The Merit–Safety Protocol grants automatic immunity to the weekly top-scoring contestant and applies percentage-based aggregation to the remaining participants, eliminating the lowest-ranked contestant among the non-immune group.

5.8. Empirical Evaluation and Adoption Rationale

A historical counterfactual evaluation was conducted to compare the proposed protocol against existing formats (Table 2). The Merit–Safety Protocol achieves a strong balance between fairness and engagement, yielding Pareto-optimal performance under the proposed metrics.

Table 2. Comparative performance evaluation

Metric	Rank System	Percent System	Judge Save
Meritocracy F	0.86 (High)	0.66 (Low)	0.68 (Low)
Engagement E	11289	11371 (Max)	11120
Status	Efficient	Efficient	Dominated

Table 2 indicates that the Merit–Safety Protocol disrupts the “impossible trinity” between fairness, engagement, and rule simplicity. Its meritocracy score reaches 0.82, substantially exceeding Judge Save and the Percent System, while maintaining high engagement at 11335. Although this engagement value is slightly below the theoretical maximum of the Percent System, the protocol preserves the commercial voting structure and improves fairness through algorithmic immunity, resulting in a robust and practically adoptable design.

6. Conclusion

This paper proposes a unified algorithmic framework based on reverse inference to recover latent decision structures in the absence of explicit preference data. Experimental results demonstrate that

the method exhibits high stability across multiple seasons of samples, with overall low uncertainty in inference outcomes (average standard deviation of approximately 0.047), validating the model's effectiveness in information-constrained scenarios. Further analysis reveals significant institutional trade-offs between fairness and participation across different rules. Ranking-based mechanisms demonstrate markedly superior fairness metrics compared to weighted mechanisms, while the effectiveness of partial intervention rules faces structural limitations.

References

- [1] Zhu, W. C., Ji, C. L., & Deng, K. Advances and Applications of Approximate Bayesian Computation Frontier Research [J]. *Applied Mathematics and Mechanics*, 2019, 40 (11): 1179-1203.
- [2] Ma, Z. G., Xu, X. H., & Liu, X. E. (2022). Three analytical frameworks for causal inference and their applications: A review. *Journal of Engineering Sciences*, 44 (07), 1231-1243. DOI: 10.13374/j.issn2095-9389.2021.07.04.002.
- [3] Yang Aijun, Liu Xiaoxing, Lin Jinguang. Research on Financial Bayesian Semiparametric GARCH Models Based on MCMC Sampling [J]. *Journal of Mathematical Statistics and Management*, 2015, 34 (03): 452-462. DOI: 10.13860/j.cnki.sltj.20150522-025.
- [4] Chen, Tengxiao. Application Research of Bayesian Vector Autoregressive Models Based on MCMC Algorithms [D]. Jishou University, 2025. DOI: 10.27750/d.cnki.gjsdx.2025.000423.
- [5] Tan Lizhi, Zhao Yiqiang. Principles, Optimization, and Application of Mixed Models in Genome-Wide Association Studies [J]. *Chinese Journal of Agricultural Science*, 2023, 56 (09): 1617-1632.
- [6] Chen Xiaowen, Li Fanhai, Zhong Bi, et al. Dynamic Changes in Iron Metabolism Among Component Blood Donors in Guangzhou and the Effectiveness of Health Education Under the GLMM Framework [J]. *Chinese Journal of Transfusion Medicine*, 2025, 38 (06): 817-823. DOI: 10.13303/j.cjbt.issn.1004-549x.2025.06.011.