

TITLE HERE

A Thesis Proposal Presented to
The Faculty of the School of Statistics
univ

In Partial Fulfillment
of the Requirements for the Degree of
Master of Science in Statistics
1st Semester A.Y. 2025-2026

by
Shaine Rosewel Matala

Contents

1	Introduction	3
1.1	Objective	3
1.2	Significance	3
1.3	Scope and Limitations	3
2	Related Literature	3
2.1	Joint confidence region for an overall ranking	3
2.1.1	Using Independence	4
2.1.2	Using Bonferroni Correction	4
2.2	T_1, T_2, T_3	4
3	Methodology	4
3.1	Joint confidence intervals for $\theta_1, \dots, \theta_K$ by using Parametric Bootstrap .	5
3.2	Joint confidence intervals for $\theta_1, \dots, \theta_K$ by using Nonrank-based method	5
3.3	Evaluation	6

1 Introduction

1.1 Objective

Rankings of government units derived from sample survey data are typically published without accompanying statistical statements that quantify uncertainty in estimated overall rankings (*add here uncertainty is just expressed for each element being ranked*). While the literature on quantifying overall uncertainty remains limited, existing methods overlook the potential correlation among ranks (*Literature that this is possible*). The objective of this study is to introduce a methodology that constructs joint confidence region for the true but unknown overall ranking while accounting for the correlation among them. In line with this, we also present ways to estimate correlation in a specific application—such as estimating the dependence structure among senatorial candidates’ rankings.

1.2 Significance

1.3 Scope and Limitations

2 Related Literature

2.1 Joint confidence region for an overall ranking

Klein et al. (2020) proposed an approach for quantifying overall rank uncertainty following the estimation of respondents’ average travel time to work in each K sampled geographical area.

Rank for the k th population as

$$r_k = \sum_{j=1}^K I(\theta_j \leq \theta_k) = 1 + \sum_{j:j \neq k} I(\theta_j \leq \theta_k), \text{ for } k = 1, \dots, K \quad (2.1)$$

The estimated overall ranking, computed on the basis of the estimates $\hat{\theta}_1, \dots, \hat{\theta}_K$, is denoted by $(\hat{r}_1, \dots, \hat{r}_K)$, where

$$\hat{r}_k = 1 + \sum_{j:j \neq k} I(\hat{\theta}_j \leq \hat{\theta}_k), \text{ for } k = 1, \dots, K \quad (2.2)$$

It follows that uncertainty in the estimators $\hat{\theta}_1, \dots, \hat{\theta}_K$ is carried over to the estimated ranking. As a result, a measure of this uncertainty should be reported alongside any resulting overall ranking.

True values, $\theta_1, \dots, \theta_K$ are unknown. It is assumed that for each $k \in \{1, 2, \dots, K\}$, there exists L_k and U_k such that

$$\theta_k \in (L_k, U_k) \quad (2.3)$$

For each $k \in \{1, 2, \dots, K\}$, define

$$\left. \begin{aligned} I_k &= \{1, 2, \dots, K\} - \{k\}, \\ \Lambda_{Lk} &= \{j \in I_k : U_j \leq L_k\}, \\ \Lambda_{Rk} &= \{j \in I_k : U_k \leq L_j\}, \\ \Lambda_{Ok} &= \{j \in I_k : U_j > L_k \text{ and } U_k > L_j\} = I_k - \{\Lambda_{Lk} \cup \Lambda_{Rk}\} \end{aligned} \right\} \quad (2.4)$$

For each $k \in \{1, 2, \dots, K\}$, and $j \in I_k$:

1. $j \in \Lambda_{Lk} \leftrightarrow (L_j, U_j) \cap (L_k, U_k) = \emptyset$ and (L_j, U_j) lies to the left of (L_k, U_k) ;
2. $j \in \Lambda_{Rk} \leftrightarrow (L_j, U_j) \cap (L_k, U_k) = \emptyset$ and (L_j, U_j) lies to the right of (L_k, U_k) ;
3. $j \in \Lambda_{Ok} \leftrightarrow (L_j, U_j) \cap (L_k, U_k) \neq \emptyset$

Λ_{Lk} , Λ_{Rk} , and Λ_{Ok} are mutually exclusive, and $\Lambda_{Lk} \cup \Lambda_{Rk} \cup \Lambda_{Ok} = I_k$

Equation 2.4 implies that for each $k \in \{1, 2, \dots, K\}$,

$$r_k \in \{|\Lambda_{Lk}| + 1, |\Lambda_{Lk}| + 2, |\Lambda_{Lk}| + 3, \dots, |\Lambda_{Lk}| + |\Lambda_{Ok}| + 1\} \quad (2.5)$$

Equation 2.5 demonstrates that smaller $|\Lambda_{Ok}|$ results in smaller difference between U_k and L_k . Collectively, these yield narrower confidence intervals for the overall ranks, which are desirable. Consequently, this results in a conservative confidence region whose joint coverage probability is at least as large as the nominal level, $1 - \alpha$ (See Equation 2.6).

$$P \left[\bigcap_{k=1}^K \{\theta_k \in (L_k, U_k)\} \right] \geq 1 - \alpha \quad (2.6)$$

2.1.1 Using Independence

2.1.2 Using Bonferroni Correction

2.2 T_1, T_2, T_3

<https://mgimond.github.io/Spatial/spatial-autocorrelation.html>

<https://cran.r-project.org/web/packages/simstudy/vignettes/corelationmat.html>

3 Methodology

This section introduces the proposed methodologies to obtain confidence regions for the unknown overall true ranking. The following cases are tackled: case when items ranked are assumed to have zero and nonzero correlation. Both approaches are based on parametric bootstrap. Sections 3.1 and 3.2 discuss the algorithms for the cases mentioned. Section

3.3 shows the algorithms used to assess the performance of the proposed approaches. This makes use of coverage and metrics to measure the tightness of the estimated confidence regions.

For sections 3.1 and 3.2, let $\theta_1, \theta_2, \dots, \theta_K$ be the true parameter values and $\hat{\theta}_1, \hat{\theta}_2, \dots, \hat{\theta}_K$ be the corresponding estimates.

3.1 Joint confidence intervals for $\theta_1, \dots, \theta_K$ by using Parametric Bootstrap

The rank-based parametric bootstrap approach assumes $\hat{\theta}_1, \hat{\theta}_2, \dots, \hat{\theta}_K$ to be independent but not identically distributed estimates, where $\hat{\theta}_k \sim N(\theta_k, \sigma_k^2)$, $k = 1, 2, \dots, K$. σ_k^2 is assumed known. Denote the corresponding ordered values by $\hat{\theta}_{(1)}, \hat{\theta}_{(2)}, \dots, \hat{\theta}_{(K)}$.

Algorithm 1 Computation of Joint Confidence Region using Parametric Bootstrap

1: **for** $b = 1, 2, \dots, B$ **do**

2: Generate $\hat{\theta}_{bk}^* \sim N(\hat{\theta}_k, \sigma_k^2)$, $k = 1, 2, \dots, K$ and let $\hat{\theta}_{b(1)}, \hat{\theta}_{b(2)}, \dots, \hat{\theta}_{b(K)}$ be the corresponding ordered values

	$k = 1$	$k = 2$	\dots	$k = K$
$b = 1$	$\hat{\theta}_{1(1)}^*$	$\hat{\theta}_{1(2)}^*$	\dots	$\hat{\theta}_{1(K)}^*$
$b = 2$	$\hat{\theta}_{2(1)}^*$	$\hat{\theta}_{2(2)}^*$	\dots	$\hat{\theta}_{2(K)}^*$
\vdots	\vdots	\vdots	\dots	\vdots
$b = B$	$\hat{\theta}_{B(1)}^*$	$\hat{\theta}_{B(2)}^*$	\dots	$\hat{\theta}_{B(K)}^*$

3: Compute

$$\hat{\sigma}_{b(k)}^* = \sqrt{\text{kth ordered value among } \{\hat{\theta}_{b1}^{*2} + \sigma_1^2, \hat{\theta}_{b2}^{*2} + \sigma_2^2, \dots, \hat{\theta}_{bK}^{*2} + \sigma_K^2\} - \hat{\theta}_{(k)}^{*2}}$$

4: Compute $t_b^* = \max_{1 \leq k \leq K} \left| \frac{\hat{\theta}_{b(k)}^* - \hat{\theta}_k^*}{\sigma_{b(k)}^*} \right|$

5: **end for**

6: Compute the $(1 - \alpha)$ -sample quantile of $t_1^*, t_2^*, \dots, t_B^*$, call this \hat{t} .

7: The joint confidence region of $\theta_{(1)}, \theta_{(2)}, \dots, \theta_{(K)}$ is given by

$$\mathfrak{R} = [\hat{\theta}_{(1)} \pm \hat{t} \times \hat{\sigma}_{(1)}] \times [\hat{\theta}_{(2)} \pm \hat{t} \times \hat{\sigma}_{(2)}] \times \dots \times [\hat{\theta}_{(K)} \pm \hat{t} \times \hat{\sigma}_{(K)}]$$

where $\hat{\sigma}_{(k)}$ is computed as

$$\hat{\sigma}_{(k)} = \sqrt{\text{kth ordered value among } \{\hat{\theta}_1^2 + \sigma_1^2, \hat{\theta}_2^2 + \sigma_2^2, \dots, \hat{\theta}_K^2 + \sigma_K^2\} - \hat{\theta}_{(k)}^2}$$

3.2 Joint confidence intervals for $\theta_1, \dots, \theta_K$ by using Nonrank-based method

The nonrank-based method assumes that $\hat{\boldsymbol{\theta}} = (\hat{\theta}_1, \hat{\theta}_2, \dots, \hat{\theta}_K) \sim N(\boldsymbol{\theta}, \boldsymbol{\Sigma})$. It accounts for potential correlation among items being ranked. For this case, an exchangeable

correlation, $\boldsymbol{\rho}$ (See Equation 3.1.), is assumed and used in the calculation of the variance covariance matrix (See Equation 3.2.).

$$\boldsymbol{\rho} = (1 - \rho) \mathbf{I}_K + \rho \mathbf{1}_K \mathbf{1}_K' \quad (3.1)$$

$$\boldsymbol{\Sigma} = \boldsymbol{\Delta}^{1/2} \boldsymbol{\rho} \boldsymbol{\Delta}^{1/2} \quad (3.2)$$

where $\boldsymbol{\Delta} = \text{diag} \{\sigma_1^2, \sigma_2^2, \dots, \sigma_K^2\}$, with known σ_k 's and ρ is studied for 0.1, 0.5, 0.9.

Algorithm 2 Computation of Joint Confidence Region using Nonrank-based Method

Let the data consist of $\hat{\theta}_1, \dots, \hat{\theta}_K$ and suppose $\boldsymbol{\Sigma}$ is known

- 1: **for** $b = 1, 2, \dots, B$ **do**
- 2: Generate $\hat{\boldsymbol{\theta}}_b^* \sim N_K(\hat{\boldsymbol{\theta}}, \boldsymbol{\Sigma})$ and write $\hat{\boldsymbol{\theta}}_b^* = (\hat{\theta}_{b1}^*, \hat{\theta}_{b2}^*, \dots, \hat{\theta}_{bK}^*)'$
- 3: Compute $t_b^* = \max_{1 \leq k \leq K} \left| \frac{\hat{\theta}_{bk}^* - \hat{\theta}_k}{\sigma_k} \right|$
- 4: **end for**
- 5: Compute the $(1 - \alpha)$ -sample quantile of $t_1^*, t_2^*, \dots, t_B^*$, call this \hat{t} .
- 6: The joint confidence region of $\theta_1, \theta_2, \dots, \theta_K$ is given by

$$\mathfrak{R} = [\hat{\theta}_1 \pm \hat{t} \times \sigma_1] \times [\hat{\theta}_2 \pm \hat{t} \times \sigma_2] \times \dots \times [\hat{\theta}_K \pm \hat{t} \times \sigma_K]$$

3.3 Evaluation

Algorithm 3 is used to calculate the coverage which is defined as the proportion of times that the true parameter values fall within the confidence interval for all K simultaneously. Ideally, this should be equal to 0.90 since $\alpha = 0.1$. It also calculates the average T_1, T_2 , and T_3 . Higher values of T_1 and T_2 indicate wider confidence intervals and are therefore less desirable, whereas higher values of T_3 are preferable.

Algorithm 3 Computation of Coverage Probability for Parametric Bootstrap

For given values of $\theta_1, \theta_2, \dots, \theta_K$ and thus $\theta_{(1)}, \theta_{(2)}, \dots, \theta_{(K)}$

1: **for** replications = 1, 2, ..., 5000 **do**

2: Generate $\hat{\theta}_k \sim N(\theta_k, \sigma_k^2)$, for $k = 1, 2, \dots, K$

3: Compute the rectangular confidence region \mathfrak{R} using Algorithm 1.

4: Check if $(\theta_{(1)}, \theta_{(2)}, \dots, \theta_{(K)}) \in \mathfrak{R}$ and compute

$$\begin{aligned} T_1 &= \frac{1}{K} \sum_{k=1}^K |\Lambda_{Ok}| \\ T_2 &= \prod_{k=1}^K |\Lambda_{Ok}| \\ T_3 &= 1 - \frac{K + \sum_{k=1}^K |\Lambda_{Ok}|}{K^2} \end{aligned}$$

5: **end for**

6: Compute the proportion of times that the condition in step 4 is satisfied and the average of T_1, T_2 , and T_3 .

Algorithm 4 is similar to Algorithm 3 but computes for the coverage and average T_1, T_2 , and T_3 for the nonrank-based method.

Algorithm 4 Computation of Coverage Probability for Nonrank-based Method

For given values of $\theta_1, \theta_2, \dots, \theta_K$ and Σ

1: **for** replications = 1, 2, ..., 5000 **do**

2: Generate $\hat{\theta} \sim N_K(\theta, \Sigma)$

3: Compute the rectangular confidence region \mathfrak{R} using Algorithm 2.

4: Check if $(\theta_1, \theta_2, \dots, \theta_K) \in \mathfrak{R}$ and compute T_1, T_2 , and T_3 .

5: **end for**

6: Compute the proportion of times that the condition in step 4 is satisfied and the average of T_1, T_2 , and T_3 .

Klein et al. (2020)

Bibliography

Klein, M., Wright, T., & Wieczorek, J. (2020). *A joint confidence region for an overall ranking of populations.*