

Article

Model Averaging for Accelerated Failure Time Models with Missing Censoring Indicators

Longbiao Liao and Jinghao Liu *

Department of Statistics and Data Science, School of Economics, Jinan University, Guangzhou 510632, China; lbl023@stu2022.jnu.edu.cn

* Correspondence: jinghao123@stu2021.jnu.edu.cn

Abstract: Model averaging has become a crucial statistical methodology, especially in situations where numerous models vie to elucidate a phenomenon. Over the past two decades, there has been substantial advancement in the theory of model averaging. However, a gap remains in the field regarding model averaging in the presence of missing censoring indicators. Therefore, in this paper, we present a new model-averaging method for accelerated failure time models with right censored data when censoring indicators are missing. The model-averaging weights are determined by minimizing the Mallows criterion. Under mild conditions, the calculated weights exhibit asymptotic optimality, leading to the model-averaging estimator achieving the lowest squared error asymptotically. Monte Carlo simulations demonstrate that the method proposed in this paper has lower mean squared errors compared to other model-selection and model-averaging methods. Finally, we conducted an empirical analysis using the real-world Acute Myeloid Leukemia (AML) dataset. The results of the empirical analysis demonstrate that the method proposed in this paper outperforms existing approaches in terms of predictive performance.

Keywords: model averaging; accelerated failure time model; censoring indicator

MSC: 62D10; 62N01; 62N02



Citation: Liao, L.; Liu, J. Model Averaging for Accelerated Failure Time Models with Missing Censoring Indicators. *Mathematics* **2024**, *12*, 641. <https://doi.org/10.3390/math12050641>

Academic Editor: Jiancang Zhuang

Received: 23 January 2024

Revised: 16 February 2024

Accepted: 20 February 2024

Published: 22 February 2024



Copyright: © 2024 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

1. Introduction

In some practical scenarios, we often need to select useful models from a candidate model set. A popular approach to address this issue is model selection. Methods such as the Akaike Information Criterion (AIC) [1], Mallows' Cp [2] and Bayesian Information Criterion (BIC) [3] are designed to identify the best model. However, in cases where a single model does not receive strong support from the data, these model-selection methods may overlook valuable information from other candidate models, leading to issues of model-selection uncertainty and bias [4].

To tackle these challenges and enhance prediction accuracy, several model-averaging techniques have been developed to leverage all information from the candidate models. Taking inspiration from AIC and BIC, Buckland et al. [5] proposed smoothed AIC (SAIC) and smoothed BIC (SBIC) methods based on AIC and BIC, respectively. Hansen [6] introduced the Mallows model-averaging (MMA) estimator, obtaining weights through the minimization of Mallows' Cp criterion. The MMA estimator asymptotically attains the minimum squared error among the model-averaging estimators in its class. Subsequently, Wan et al. [7] relaxed the constraints of Hansen [6], allowing for non-nested candidate models and continuous weights. In practical applications, many datasets exhibit heteroscedasticity. Therefore, it is essential to explore model-averaging methods tailored for heteroscedastic settings. Firstly, Hensen and Racine [8] proposed Jackknife model averaging (JMA), which determines weights by minimizing a cross-validation criterion. JMA significantly reduces Mean Squared Error (MSE) compared to MMA when errors are heteroscedastic. Secondly, Liu and Okui [9] modified the MMA method proposed by Hensen [6] to make it suitable for

heteroscedastic scenarios. Furthermore, Zhao et al. [10] extended [6]’s work by estimating the covariance matrix based on the weighted average of squared residuals corresponding to all candidate models. This approach improves the model average estimator under heteroskedasticity settings.

In survival analysis, the accelerated failure time (AFT) model provides a straightforward description of how covariates directly impact survival time and has consequently garnered widespread attention. There are several parameter-estimation methods for the Accelerated Failure Time (AFT) model, including Miller’s estimator [11], Buckley–James estimator [12], Koul–Susarla–Van Ryzin (KSV) estimator [13] and WLS estimator [14]. However, all these methods assume that the censoring indicator is observable. Therefore, Wang and Dinse [15] improved the KSV estimator to make it adaptable to situations where the censoring indicator is missing.

Under practical conditions, it is common to encounter situations where only the observed time is available and it is uncertain whether the event of interest has occurred. In such cases, data suffer from missingness in the censoring indicator. For example, in a clinical trial for lung cancer, a patient may die for unknown reasons and while the survival time is observed, it is uncertain whether the patient died specifically due to lung cancer. This situation leads to missingness in the censoring indicator. Previous studies have mainly addressed the issue of missingness in the censoring indicator under a specific model. Research on model averaging for right-censored data typically assumes that the censoring indicator is observable. Therefore, this paper adopts the inverse probability weighting method proposed by [15] to construct the response variable. Through appropriate weight-selection criteria, weights are chosen to build the model-averaged estimator for the accelerated failure time model. It significantly enhances the predictive performance of the model and mitigates the bias introduced by the selection of a single model. Compared to previous research, this paper makes two main contributions: First, it introduces a novel model-averaging method for the case of missingness in the censoring indicator. Second, the paper allows for heteroscedasticity and employs model-averaging techniques to estimate variance.

The remaining sections of this paper are organized as follows. In Section 2, we commence by introducing the notation and progressively delineate the methodology and associated theoretical properties of the proposed model-averaging approach. In Section 3, we report the Monte Carlo simulation results. In Section 4, we assess the predictive performance of the proposed model-averaging method against other approaches using the real-world Acute Myeloid Leukemia (AML) dataset. In Section 5, we provide a comprehensive summary of the entire paper and suggest future research directions in this area. All theorem proofs will be presented in Appendix A.

2. Methodology and Theoretical Property

We denote $Y = \log(T) = (Y_1, \dots, Y_n)'$, $C = \log(V) = (C_1, \dots, C_n)'$, where T represents the survival time and V denotes the censored time. $X = (X'_1, X'_2, \dots, X'_n)'$ denotes the covariate matrix for n independent observations, where $X_i = (x_{i1}, x_{i2}, \dots, x_{ip})$. The accelerated failure time model can be expressed as follows:

$$Y_i = \mu_i + e_i = \sum_{j=1}^p \beta_j x_{ij} + e_i, \quad (i = 1, \dots, n), \tag{1}$$

where e_i is the random error with $E(e_i|X_i) = 0$ and $E(e_i^2|X_i) = \sigma_i^2$.

We assume that there are M candidate models in the candidate model set. Where the m th candidate model contains p_m covariates. Following [7], these candidate model forms are non-nested. The m th candidate model is

$$Y_{mi} = \sum_{j=1}^{p_m} \beta_j x_{ij} + e_{mi}, \quad (i = 1, \dots, n), \tag{2}$$

for $m = 1, \dots, M$. The matrix form of (2) is

$$Y_m = X_m \beta_m + e_m, \tag{3}$$

where X_m is an $n \times p_m$ dimensional full column-rank matrix, $Y_m = (Y_{m1}, \dots, Y_{mn})'$, $\beta_m = (\beta_{1m}, \dots, \beta_{p_m m})'$, $e_m = (e_{m1}, \dots, e_{mn})'$.

In the case of right censored data, the response variable Y_i might be censored, making it unobservable. We only observe (Z_i, X_i, δ_i) , where $Z_i = \min(Y_i, C_i)$ and the censoring indicator $\delta_i = I(Y_i \leq C_i)$. Define a missingness indicator ξ_i which is 1 if δ_i is observed and is 0 otherwise. When the censoring indicators are missing, the observed data are $\{Z_i, X_i, \xi_i, \xi_i \delta_i\}$. For simplicity, we set $U_i = (Z_i, X_i)'$. In this paper, similar to [15], we assume the missing mechanism for δ to be:

$$P(\xi = 1|Z, X, \delta) = P(\xi = 1|Z).$$

This assumption is more stringent than the missing at random (MAR) condition yet less restrictive than the assumption of missing completely at random (MCAR).

Koul et al. [13] introduced a method that involves synthetic data for constructing linear regression models. Wang and Dinse [15] extended [13]'s method to address the situation where censoring indicators are missing. In our work, we follow the approach proposed by [15] to construct a response in the form of inverse probability weighting, specifically:

$$Y_{Wi} = \frac{\xi_i \delta_i}{\pi(Z_i)} + (1 - \frac{\xi_i}{\pi(Z_i)})m(U_i) \tag{4}$$

$$1 - G_n(Z_i) Z_i,$$

where $\pi(z) = E(\xi|Z = z)$, $m(u) = E(\delta|U = u)$. $G_n(\cdot)$ represents the cumulative distribution function of C . It is easy to observe that under the missing data mechanism in this paper:

$$E(Y_{Wi}|X_i) = \mu_i = X_i \beta.$$

Similar to Equation (2), we have:

$$Y_{Wi} = \sum_{j=1}^{p_m} \beta_j x_{ij} + e_{Wi}, \quad (i = 1, \dots, n), \tag{5}$$

where $E(e_{Wi}|X_i) = 0$, $\sigma_{Wi}^2 = \text{var}(e_{Wi}|X_i)$. This is expressed in matrix form as:

$$Y_W = X_m \beta_m + e_W, \tag{6}$$

where $Y_W = (Y_{W1}, \dots, Y_{Wn})'$, $e_W = (e_{W1}, \dots, e_{Wn})'$. And then the weighted least squares estimator of β_m :

$$\hat{\beta}_m = (X_m' D X_m)^{-1} X_m' D Y_W, \tag{7}$$

where $D = \text{diag}\{\frac{1}{\sigma_{W1}^2}, \dots, \frac{1}{\sigma_{Wn}^2}\}$.

Let $\mu_{mi} = E(Y_{Wi}|X_i)$; subsequently, the estimation for the m th candidate model $\mu_m = (\mu_{m1}, \dots, \mu_{mn})'$ is given by:

$$\hat{\mu}_m = X_m \hat{\beta}_m = X_m (X_m' D X_m)^{-1} X_m' D Y_W = P_m Y_W, \tag{8}$$

where $P_m = X_m (X_m' D X_m)^{-1} X_m' D$. Denote weight vector $w = (w_1, \dots, w_M)^T$, belonging to the set $\mathcal{H}_M = \{w \in [0, 1]^M : \sum_{m=1}^M w_m = 1\}$. The model-averaging estimator of μ is defined as follows:

$$\hat{\mu}_{G_n}(w) = \sum_{m=1}^M w_m X_m \hat{\beta}_m = \sum_{m=1}^M w_m X_m (X_m' D X_m)^{-1} X_m' D Y_W = P(w) Y_W, \tag{9}$$

for any $w \in \mathcal{H}_M$, where $P(w) = \sum_{m=1}^M w_m X_m (X'_m D X_m)^{-1} X'_m D$.

Define the square loss function $L_{G_n}(w) = \|\mu - \hat{\mu}(w)\|^2$, where $\|\cdot\|$ denotes the Euclidean norm. Then the risk function is defined as:

$$R_{G_n}(w) = E(L_{G_n}(w)) = \|P(w)\mu - \mu\|^2 + tr\{P(w)\Omega P'(w)\}, \tag{10}$$

where $\Omega = diag\{\sigma_{W1}^2, \dots, \sigma_{Wn}^2\}$. The derivation of (10) is as follows:

$$\begin{aligned} R_{G_n}(w) &= E[L_{G_n}(w)] \\ &= E[(\mu - \hat{\mu}(w))'(\mu - \hat{\mu}(w))] \\ &= E[u'\mu - 2u'P(w)Y_W + Y'_W P'(w)P(w)Y_W] \\ &= u'\mu - 2u'P(w)\mu + u'P'(w)P(w)\mu + tr(P'(w)\Omega P(w)) \\ &= (P(w)\mu - \mu)'(P(w)\mu - \mu) + tr(P(w)\Omega P'(w)). \end{aligned} \tag{11}$$

Regarding the choice of weights, a natural approach is to minimize the risk function to obtain the optimal weights. However, as shown in Equation (11), we recognize that the risk function includes the unknowns μ , which makes it infeasible to directly minimize the risk function to obtain the optimal weights. Therefore, we replace μ with Y_W and seek an unbiased estimator of the risk function as the criterion for weight selection.

Define the criterion for weight selection as

$$C_{G_n}(w) = \|Y_W - \hat{\mu}(w)\|^2 + 2tr\{P(w)\Omega\}. \tag{12}$$

It is not difficult to observe that $E(C_{G_n}(w)) = R_{G_n}(w) + \sum_{i=1}^n \sigma_{Wi}^2$. By disregarding a term that is independent of w , $C_{G_n}(w)$ serves as an unbiased estimator of the risk function.

In practice, $m(\cdot)$, $\pi(\cdot)$ and $G_n(\cdot)$ are usually unknown; therefore, we need to estimate them. Firstly regarding the estimation of $m(u)$, it is usually estimated by the Logit model. Suppose $m(u)$ is estimated by the parametric model $m_0(u; \theta)$, where $m_0(u; \theta) = \frac{e^{u\theta}}{1+e^{u\theta}}$. By the maximum likelihood estimation method, we can obtain the parameter estimate $\hat{\theta}_n$ for the parameter θ . $\pi(z)$ usually can be estimated nonparametrically by

$$\hat{\pi}_n(z) = \frac{\sum_{i=1}^n \xi_i W\left(\frac{z - Z_i}{b_n}\right)}{\sum_{i=1}^n W\left(\frac{z - Z_i}{b_n}\right)},$$

where $W(\cdot)$ is a kernel function and b_n is a bandwidth sequence. Next, we define $u(z) = E(\delta|Z = z)$, $u(z)$ estimated nonparametrically by

$$\hat{u}_n(z) = \frac{\sum_{i=1}^n \left(\delta_i \frac{\xi_i}{\hat{\pi}_n(Z_i)} K\left(\frac{z - Z_i}{h_n}\right)\right)}{\sum_{i=1}^n \left(\frac{\xi_i}{\hat{\pi}_n(Z_i)} K\left(\frac{z - Z_i}{h_n}\right)\right)},$$

where $K(\cdot)$ is a kernel function and h_n is a bandwidth sequence. We adopt the following estimator of $G_n(z)$:

$$\hat{G}_n(z) = 1 - \prod_{i: Z_i \leq z} \left(\frac{n - R_i}{n - R_i + 1}\right)^{1 - \hat{u}_n(Z_i)},$$

where R_i denotes the rank of Z_i .

Next, replacing $m(\cdot)$, $\pi(\cdot)$ and $G_n(\cdot)$ with $m_0(\cdot, \cdot)$, $\hat{\pi}_n(\cdot)$ and $\hat{G}_n(\cdot)$, we have:

$$\hat{Y}_{Wi} = \frac{\frac{\xi_i \delta_i}{\hat{\pi}_n(Z_i)} + \left(1 - \frac{\xi_i}{\hat{\pi}_n(Z_i)}\right) m_0(U_i, \hat{\theta}_n)}{1 - \hat{G}_n(Z_i)} Z_i.$$

And the corresponding weight selection criterion is as follows:

$$C_{\hat{G}_n}(w) = \|\hat{Y}_W - \hat{\mu}_{\hat{G}_n}(w)\|^2 + 2trace\{P(w)\Omega\}, \tag{13}$$

where $\hat{Y}_W = (\hat{Y}_{W1}, \dots, \hat{Y}_{Wn})$. The weights for minimizing $C_{\hat{G}_n}(w)$ are given by:

$$\tilde{w} = \arg \min_{w \in \mathcal{H}_M} C_{\hat{G}_n}(w). \tag{14}$$

Then, we enumerate the necessary regularity conditions for the asymptotic optimality.

- (C1) Let $S(t) = 1 - (1 - F(t)(1 - G_n(t)))$ and $\tau_H = \inf\{t : S(t) = 1\}$, where $F(t)$ is the cumulative distribution function of Y_i . Assume that $1 - G_n(\tau_H -) > 0$.
- (C2) There exists a positive constant k such that $\max_{1 \leq i \leq n} |\mu_i| \leq k$.
- (C3) Denote $\zeta_n = \inf_{w \in \mathcal{H}_M} R_{G_n}(w)$ and w_m^0 is an $M \times 1$ unit vector in which the m th element is 1 and the others are 0. For some integer $1 \leq J < \infty$ and some positive constant k such that $E(e_i^{4J}) \leq k < \infty$, assume

$$M \zeta_n^{-2J} \sum_{m=1}^M \left\{ R_{G_n}(w_m^0) \right\}^J \rightarrow 0.$$

- (C4) There exists $\epsilon > 0$ such that $\inf e_{Wi}^2 > \epsilon, i = 1, \dots, n$.
- (C5) $m(\cdot)$ and $\pi(\cdot)$ are bounded.
- (C6) $nh_n \rightarrow \infty$ and $nh_n^2 \rightarrow 0$.
- (C7) Let $\tilde{p} = \max_m p_m, \rho_{ii}^m$ denote the i th diagonal element of P_m . There exists a constant c such that $|\rho_{ii}^{(m)}| \leq cn^{-1} p_m$.

Condition (C1) is utilized in [16] and it ensures that $1 - G_n(t)$ is not equal to 0. Condition (C2) mandates that the conditional expectation of μ_i remains within bounded limits, in line with assumptions seen in prior research, including [7,17]. Condition (C3) is a requirement commonly found in model-averaging literature (e.g., [7,18]). Condition (C4) mandates the non-degeneracy of the covariance matrix Ω as $n \rightarrow \infty$. Similar assumptions can also be found in [9,10]. Similar to [15], Conditions (C5) and (C6) impose constraints on the bounds of $m(\cdot), \pi(\cdot)$ and bandwidth, respectively. Condition (C7) is frequently employed in the analysis of the asymptotic optimality of cross-validation methods, as seen in prior works like [8].

Theorem 1. Under Conditions (C1) to (C6),

$$\frac{L_{\hat{G}_n}(\tilde{w})}{\inf_{w \in \mathcal{H}_M} L_{\hat{G}_n}(w)} \xrightarrow{p} 1.$$

Theorem 1 establishes the asymptotic optimality of the model-averaging procedure employing weights \tilde{w} , as its squared loss converges to that of the infeasible best possible model average estimator.

In most cases, Ω is unknown and needs to be estimated. We estimate Ω using residuals derived from the model-averaging process: $\hat{e}(w) = \hat{Y}_W - \hat{\mu}(w) = \{\hat{e}_{W1}(w), \dots, \hat{e}_{Wn}(w)\}'$. Specifically, the estimator of Ω is

$$\hat{\Omega}(w) = \text{diag}\{\hat{\sigma}_{W1}^2(w), \dots, \hat{\sigma}_{Wn}^2(w)\}, \tag{15}$$

where $\hat{\sigma}_{Wi}^2 = \text{var}(\hat{e}_{Wi})$.

In the existing literature on model averaging, most estimates of variance are predominantly derived from the largest candidate model, as exemplified by works such as [6,16]. In contrast, our approach, following [10], leverages information from all candidate models for estimation rather than relying on a single model. Such an estimation method is more robust. Replacing Ω by $\hat{\Omega}(w)$ in (13), $C(w)$ becomes

$$\hat{C}_{\hat{G}_n}(w) = \|\hat{Y}_W - \hat{\mu}_{\hat{G}_n}(w)\|^2 + 2\text{trace}\{P(w)\hat{\Omega}(w)\}. \tag{16}$$

The weights that minimize $\widehat{C}_{\widehat{G}_n}(w)$ are as follows:

$$\widehat{w} = \arg \min_{w \in \mathcal{H}_M} \widehat{C}_{\widehat{G}_n}(w). \tag{17}$$

This weight selection criterion $\widehat{C}_{\widehat{G}_n}(w)$ is a cubic function of w .

Theorem 2. Under Conditions (C1) to (C7),

$$\frac{L_{\widehat{G}_n}(\widehat{w})}{\inf_{w \in \mathcal{H}_M} L_{\widehat{G}_n}(w)} \xrightarrow{p} 1.$$

3. Simulation

In the simulation study, we generate data from the accelerated failure time (AFT) model, $\log(T_i) = Y_i = \sum_{j=1}^{1000} \beta_j x_{ij} + e_i$, where $\beta_j = 1/j^2$; the observations of $X_i = (x_{i1}, x_{i2}, \dots, x_{i1000})$ are generated from a multivariate normal distribution with zero mean and covariance matrix $\Sigma = (\sigma_{ij})$ with $\sigma_{ij} = 0.5^{|i-j|}$. The errors e_i follow normal distribution $N(0, \gamma^2(x_{i2}^4 + 0.01))$. By varying the value of γ , we allow R^2 to range from 0.1 to 0.9. This variance specification closely resembles that of [8]. However, we introduce a small constant, 0.01, to ensure that the variances remain strictly positive. The censoring time C_i is generated from $N(C_0, 7)$. By varying the value of C_0 , we achieve censoring rates (CRs) of approximately 20%, 40%. We set sample sizes $n = 150, 300$. Here, our model configuration is set in a nested form, meaning the first m models include the first m regressors. The number of candidate models M was set to be $\lceil 3n^{1/3} \rceil$, where $\lceil x \rceil$ denote the smallest integer greater than x .

Based on the missing mechanism described in this paper, we assume that the probability of missing censoring indicators, denoted as $1 - \pi(z)$, is determined via a logistic model: $\log\{\frac{\pi(z)}{1-\pi(z)}\} = \theta_1 + \theta_2 z$. Following [15], we employed the uniform kernel function $W(x) = \frac{1}{2}$ for $|x| \leq 1$ and $W(x) = 0$ otherwise. Additionally, we used the bi-weight kernel function $K(x) = \frac{15}{16}(1 - 2x^2 + x^4)$ for $|x| \leq 1$ and $K(x) = 0$ otherwise. The bandwidths were $b_n = h_n = n^{-\frac{1}{3}} \max(Z)$. We estimated $m(u)$ under the logistic model: $\log\{\frac{m(u)}{1-m(u)}\} = \gamma_1 + \gamma_2 z + \gamma_3 x$. As highlighted by [19], when the data on δ are completely (or quasi-completely) separated, the maximum likelihood estimate of $\gamma = (\gamma_1, \gamma_2, \gamma_3)$ does not exist. In our simulation setup, the number of covariates significantly exceeds the sample size. Therefore, we employ the lasso method to estimate the parameters.

We compare the proposed Model-Averaging method for the Missing Censoring Indicators in the Heteroscedastic setting (HCIMA) with other classical model-selection and model-averaging methods in this article. Brief descriptions of these methods are provided below:

- The model-selection methods rely on AIC and BIC, where the AIC and BIC criterion for the m th model are defined as follows:

$$AIC(m) = \log(\widehat{\sigma}_{\widehat{G}_n m}^2) + 2n^{-1} \text{tr}(P_m),$$

and

$$BIC(m) = \log(\widehat{\sigma}_{\widehat{G}_n m}^2) + n^{-1} \text{tr}(P_m) \log(n),$$

where $\widehat{\sigma}_{\widehat{G}_n m}^2 = n^{-1} \|\widehat{Y}_W - \widehat{\mu}_{\widehat{G}_n m}\|^2$.

- Model methods based on SAIC and SBIC: The weights for the m th candidate model are given by:

$$w_{(m)}^{AIC} = \exp(-AIC_m/2) / \sum_{m=1}^M \exp(-AIC_m/2),$$

$$w_{(m)}^{BIC} = \exp(BIC_m/2) / \sum_{m=1}^M \exp(-BIC_m/2),$$

where $AIC_m = AIC(m) - \min(AIC)$, $BIC_m = BIC(m) - \min(BIC)$.

- Additionally, we compare our approach with the method that estimates the variance using the maximum candidate model (MCIMA). And the specifics of variance estimation and weight selection in their approach are as follows:

$$\hat{\sigma}_{\hat{G}_n} = (\hat{\sigma}_{\hat{G}_n1}, \dots, \hat{\sigma}_{\hat{G}_nn})^T = \sqrt{\frac{n}{n-M}}(I - P_M)\hat{Y}_W,$$

$$\hat{C}_n(w) = \|\hat{Y}_W - \hat{\mu}(w)\|^2 + trace\{P(w)\hat{\Omega}\},$$

where $\hat{\Omega} = diag\{\hat{\sigma}_{\hat{G}_n1}^2, \dots, \hat{\sigma}_{\hat{G}_nn}^2\}$.

In the simulation, we utilize the Mean Squared Error (MSE) to evaluate the performance of various methods, where the MSE is defined as $\frac{1}{n} \|\hat{\mu}_{\hat{G}_n} - \mu\|^2$. We present the mean of MSEs from 500 replications.

Figures 1 and 2, respectively, show the Mean Squared Error (MSE) values for various methods across 500 repetitions under different censored rates and sample sizes, with missing rates of 20% and 40%. In terms of Mean Squared Error (MSE), our proposed HCIMA method outperforms other approaches. Additionally, the MCIMA method performs better than existing methods in all cases except for when compared to HCIMA. Furthermore, it is evident that SAIC and SBIC outperform their respective AIC and BIC counterparts, further highlighting the advantages of model-averaging methods.

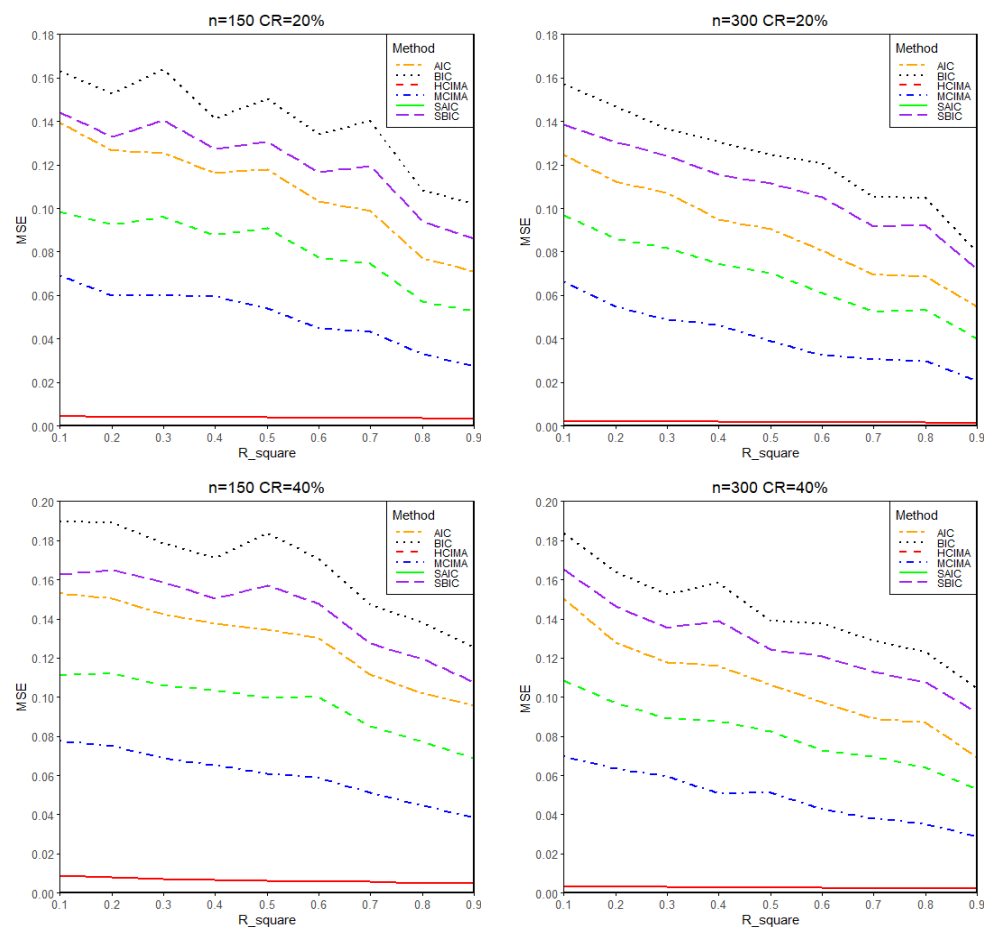


Figure 1. Mean Squared Errors (MSEs) of various methods under different sample sizes and censor rates at MR = 20%.

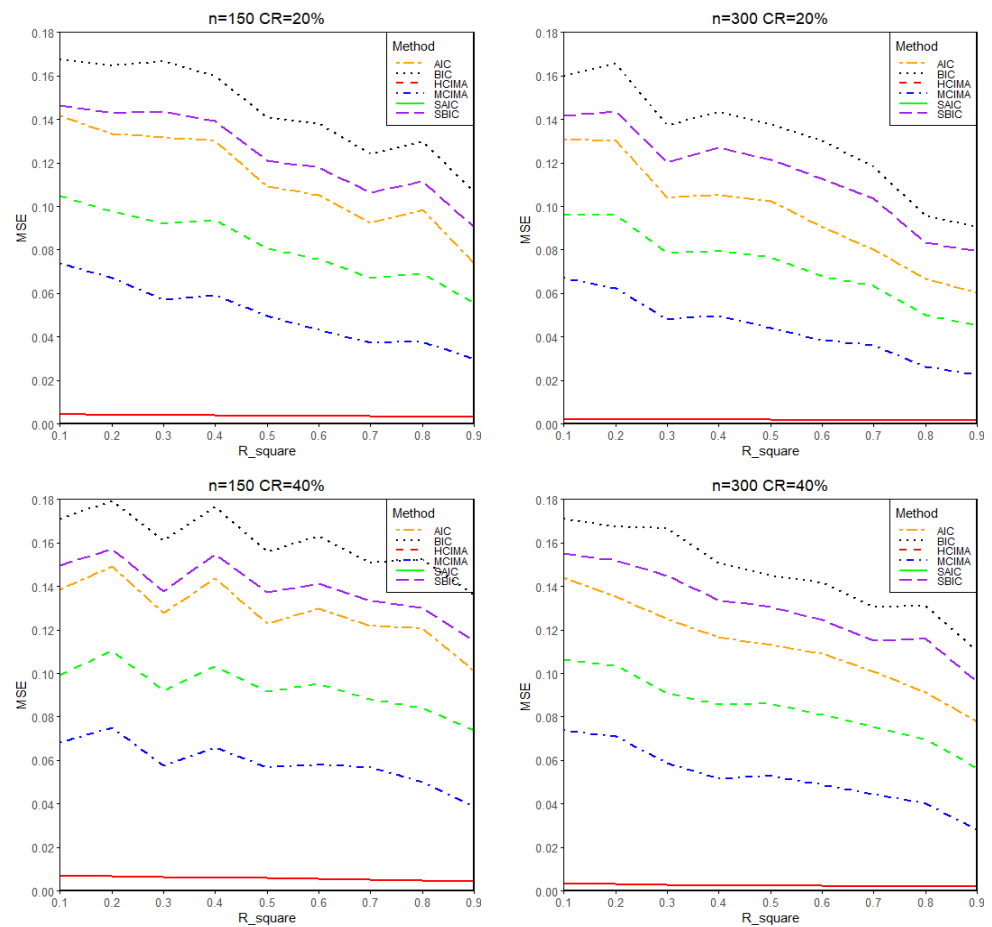


Figure 2. Mean Squared Errors (MSEs) of various methods under different sample sizes and censor rates at MR = 40%.

Comparing Figures 1 and 2, it is observed that the MSE at MR = 20% is slightly higher than at MR = 40%. The reason for this occurrence is that when $\zeta_i = 1, \delta_i = 0$, the signs of Y_{Wi} and Z_i are opposite. As MR increases, the occurrence of the $\zeta_i = 1, \delta_i = 0$ situation decreases. Although this result may seem counterintuitive, it does not affect the performance of the method proposed in this paper, which still keeps its advantages in this case.

4. Real Data Analysis

In this section, we assess the predictive performance of our proposed HCIMA method using the real Acute Myeloid Leukemia (AML) dataset. This dataset contains 672 samples, including 97 variables such as patient age, survival time, gender, race, mutation count, etc. For more specific information about this dataset, we refer the reader to https://www.cbioportal.org/study/clinicalData?id=aml_ohsu_2018 (accessed on 13 December 2023).

We selected ten variables for analysis: Cause Of Death, Age, Sex, Overall Survival Status, Overall Survival Months (Survival Time), Number of Cumulative Treatment Stages, Cumulative Treatment Regimen Count, Mutation Count, Platelet Count and WBC (White Bloodcell Count). After removing rows with missing values, we retained a total of 396 samples. We treat samples with unknown causes of death as missing censoring indicators. Among these 396 samples, 76 have unknown causes of death and 167 samples are still alive after the clinical trial ends. Therefore, the missing rate is approximately 19% and the censoring rate is 42%. We focus on the impact of seven variables, excluding “Cause Of Death” and “Overall Survival Status” on Survival Time. Therefore, we can construct $2^7 - 1 = 127$ non-nested candidate models.

We randomly select data from n_0 samples as the training dataset, while the remaining $n_1 = n - n_0$ samples are used as the testing dataset. We set the training dataset size to 50%, 60%, 70% and 80% of the total dataset size, respectively. Following [16,20], we employed the normalized mean squared prediction error (NMSPE) as the performance metric:

$$\text{NMSPE} = \frac{\sum_{i=n_0+1}^n (\hat{Y}_{Wi} - \hat{\mu}_i)^2}{\min_{m=1,2,\dots,M} \sum_{i=n_0+1}^n (\hat{Y}_{Wi} - \hat{\mu}_{mi})^2},$$

where $\hat{\mu}_i$ represents the predicted value and $\hat{\mu}_{mi}$ denotes the value of $\hat{\mu}$ for the m th model.

We calculate the mean, the standard deviation and the optimal rate of each method over these 1000 repetitions. Specifically, the optimal rate refers to the frequency at which the minimum value is achieved across these 1000 repetitions.

Table 1 displays the mean, optimal rates and standard deviations of NMSPE for each method over 1000 repetitions. Consistent with the simulation results, the HCIMA method exhibits the lowest average NMSPE and standard deviation and the highest optimal rate. The MCIMA method also performs well, ranking second after HCIMA. This indicates that the proposed model-averaging methods in this paper demonstrate superior predictive performance compared to other approaches.

Table 1. The mean, optimal rate and standard deviation of NMSPE.

	Method	AIC	SAIC	BIC	SBIC	MCIMA	HCIMA
50%	Mean	1.3628	1.3370	1.3517	1.3345	1.2765	1.2165
	Standard deviation	0.5283	0.5060	0.5123	0.4970	0.4039	0.3500
	Optimal rate	0.084	0.137	0.042	0.093	0.306	0.338
60%	Mean	1.3663	1.3388	1.3556	1.3404	1.2651	1.1800
	Standard deviation	0.5504	0.5166	0.5151	0.5068	0.4343	0.3066
	Optimal rate	0.094	0.119	0.049	0.091	0.288	0.359
70%	Mean	1.3347	1.3213	1.3361	1.3259	1.2451	1.1766
	Standard deviation	0.5324	0.5232	0.5288	0.5257	0.3433	0.2966
	Optimal rate	0.097	0.140	0.057	0.079	0.259	0.368
80%	Mean	1.2794	1.2619	1.2828	1.2628	1.2034	1.1504
	Standard deviation	0.4865	0.4714	0.4861	0.4777	0.2941	0.2030
	Optimal rate	0.083	0.165	0.063	0.129	0.240	0.320

5. Discussion

To address the uncertainty in model selection and enhance predictive accuracy, this paper proposes a novel model-averaging approach for the accelerated failure time model with missing indicators. Moreover, we establish asymptotic optimality under certain mild conditions. In Monte Carlo simulations, the method proposed in this paper exhibits lower mean squared errors compared to other model-selection and model-averaging methods. Empirical results demonstrate that the proposed method has a lower NMSPE compared to other approaches, indicating its superior predictive performance. This further underscores the applicability of the proposed method to real-life data scenarios with missing censoring indicators.

In this paper, we introduce the inverse probability weighted form of response variable proposed in [15]. The primary advantage of this form of response variable lies in its double robustness, making it less susceptible to the impact of model misspecification (if $\pi(\cdot)$ or $m(\cdot)$ is misspecified). However, as mentioned in [15], its drawback, compared to synthetic response [13], regression calibration and imputation [15], is a larger variance. Yet, in practical scenarios, the harm caused by model misspecification often outweighs the harm of higher variance. Therefore, in our work, we follow the recommendation of [15] to use

the inverse probability weighted form of the response variable. A future research direction is to further enhance this response variable for better applicability in the context of missing censoring indicators.

As far as we know, there is currently very limited research on model averaging for missing censoring indicators. Therefore, there are still many questions that deserve further investigation. There is potential for extending our approach to high-dimensional data in terms of data and in terms of models, exploration into partial linear models, generalized linear models and other extensions could be pursued.

Author Contributions: Conceptualization, L.L.; Methodology, L.L.; Writing—review and editing, L.L.; Software, J.L.; Data curation, J.L.; Writing—original draft, J.L. All authors have read and agreed to the published version of the manuscript.

Funding: This research received no external funding.

Data Availability Statement: Data are contained within the article.

Acknowledgments: We would like to thank the reviewers and editors for their careful reading and constructive comments.

Conflicts of Interest: The authors declare no conflict of interest.

Appendix A

In this appendix, we provide the proofs for Theorems 1 and 2. To facilitate the presentation, we begin with several lemmas.

Lemma A1. Under Conditions (C1) to (C3), there exists a positive constant c_1 such that

$$\max_{1 \leq i \leq n} E\left(e_{W,i}^{4J} \mid X_i\right) \leq c_1,$$

where J is given in Condition (C3).

Lemma A1 is consistent with Lemma 6.1 in [16] and under our specific conditions, the proof technique for Lemma 1 is the same as the proof technique for Lemma 6.1 in [16].

Lemma A2. Under Conditions (C1) to (C3),

$$\left| \frac{L_{G_n}(\mathbf{w})}{R_{G_n}(\mathbf{w})} - 1 \right| = o_p(1).$$

Under our specific conditions, we can prove this lemma using the same techniques as the proof of (A.3) in [7]. Therefore, we omit the proof here.

Lemma A3. Under Conditions (C1) to (C5), as $n \rightarrow \infty$, we have

$$\left\| \hat{Y}_W - Y_W \right\|^2 = o_p(1).$$

Proof of Lemma A3.

$$\begin{aligned} \left\| \hat{Y}_W - Y_W \right\|^2 &= \sum_{i=1}^n \left\{ \frac{\frac{\xi_i \delta_i}{\pi(Z_i)} + (1 - \frac{\xi_i}{\pi(Z_i)})m(U_i)}{1 - G_n(Z_i)} - \frac{\frac{\xi_i \delta_i}{\hat{\pi}_n(Z_i)} + (1 - \frac{\xi_i}{\hat{\pi}_n(Z_i)})m_0(U_i, \hat{\theta}_n)}{1 - \hat{G}_n(Z_i)} \right\}^2 Z_i \\ &\leq K \sum_{i=1}^n \left\{ \frac{1}{1 - G_n(Z_i)} - \frac{1}{1 - \hat{G}_n(Z_i)} \right\}^2 Z_i \\ &\leq C \left\{ n^{1/2} \max_{1 \leq i \leq n} \left| \frac{1}{1 - \hat{G}_n(Z_i)} - \frac{1}{1 - G_n(Z_i)} \right| \right\}^2 \left(\frac{1}{n} \mu^T \mu + \frac{1}{n} e_W^T e_W \right), \end{aligned}$$

where K is a constant. By Condition (C2), we have $\frac{1}{n}\mu^T\mu = O_p(1)$ and $\frac{1}{n}e_W^Te_W = O_p(1)$. According to [15], $\widehat{G}_n(Z_i) - G_n(Z_i) = o_p(1)$. Combined with Condition (C1), we have:

$$\left| \frac{\widehat{G}_n(z) - G_n(z)}{1 - G_n(z)} \right| = o_p(1),$$

and

$$\left| \frac{\widehat{G}_n(z) - G_n(z)}{1 - \widehat{G}_n(z)} \right| = o_p(1).$$

Similar to the proof of Lemma 6.2 in [16], we have:

$$n^{1/2} \max_{1 \leq i \leq n} \left| \frac{1}{1 - \widehat{G}_n(Z_i)} - \frac{1}{1 - G_n(Z_i)} \right| = o_p(1).$$

Furthermore, we can obtain

$$\|\widehat{Y}_W - Y_W\|^2 = o_p(1).$$

□

With the three lemmas mentioned above, we can now proceed to prove Theorem 1.

Proof of Theorem 1. First, we note that

$$\begin{aligned} C_{\widehat{G}_n}(\mathbf{w}) &= \|\widehat{Y}_W - \widehat{\mu}_{\widehat{G}_n}(\mathbf{w})\|^2 + 2\text{trace}\{P(\mathbf{w})\Omega\} \\ &= \|\widehat{Y}_W - \mu + \mu - \widehat{\mu}_{\widehat{G}_n}(\mathbf{w})\|^2 + 2\text{trace}\{P(\mathbf{w})\Omega\} \\ &= \|\widehat{Y}_W - \mu\|^2 + \|\mu - \widehat{\mu}_{\widehat{G}_n}(\mathbf{w})\|^2 + 2(\widehat{Y}_W - \mu)'(\mu - \widehat{\mu}_{\widehat{G}_n}(\mathbf{w})) + 2\text{trace}\{P(\mathbf{w})\Omega\} \\ &= \|e_W\|^2 + L_{\widehat{G}_n}(\mathbf{w}) + 2e'_W(\mu - P(\mathbf{w})\mu + P(\mathbf{w})\mu - \widehat{\mu}_{\widehat{G}_n}(\mathbf{w})) + 2\text{trace}\{P(\mathbf{w})\Omega\} \\ &= L_{\widehat{G}_n}(\mathbf{w}) + 2e'_W(I - P(\mathbf{w}))\mu + 2\text{trace}(P(\mathbf{w})\Omega) - 2e'_WP(\mathbf{w})e_W + \|e_W\|^2. \end{aligned}$$

Following [7], except for a term unrelated to w , to prove Theorem 1, we only need to verify

$$\sup_{\mathbf{w} \in \mathcal{H}_M} \frac{|e'_W(I - P(\mathbf{w}))\mu|}{R_{G_n}(\mathbf{w})} = o_p(1), \tag{A1}$$

$$\sup_{\mathbf{w} \in \mathcal{H}_M} \frac{|\text{trace}(P(\mathbf{w})\Omega) - 2e'_WP(\mathbf{w})e_W|}{R_{G_n}(\mathbf{w})} = o_p(1), \tag{A2}$$

$$\sup_{\mathbf{w} \in \mathcal{H}_M} \left| \frac{L_{\widehat{G}_n}(\mathbf{w})}{R_{G_n}(\mathbf{w})} - 1 \right| = o_p(1). \tag{A3}$$

We begin by proving Equation (A1). As per Equation (11), we can ascertain that:

$$R_{G_n}(\mathbf{w}_m^0) \geq \|P(\mathbf{w}_m^0)\mu - \mu\|^2, \tag{A4}$$

$$R_{G_n}(\mathbf{w}_m^0) \geq \text{trace}(P(\mathbf{w}_m^0)\Omega P^T(\mathbf{w}_m^0)). \tag{A5}$$

Furthermore, we denote the maximum eigenvalue of matrix A as $\lambda_{\max}(A)$; since P_m is an idempotent matrix, we have:

$$\lambda_{\max}(P_m) = 1, \tag{A6}$$

$$\lambda_{\max}\{P(\mathbf{w})\} \leq \sum_{m=1}^M w_m \lambda_{\max}\{P_m\} \leq 1. \tag{A7}$$

According to the proof of Theorem 1 in [21], we have:

$$\lim_{n \rightarrow \infty} \sup_{w \in \mathcal{H}_M} \lambda_{\max}(P(w)P(w)') < \infty. \tag{A8}$$

Applying the triangle inequality, Bonferroni’s inequality, Chebyshev’s inequality and Theorem 2 of [22], we can conclude, for any $\tau > 0$,

$$\begin{aligned} & P \left\{ \sup_{w \in \mathcal{H}_M} \frac{|e'_W(I - P(w))\mu|}{R_{G_n}(w)} > \tau \right\} \\ & \leq P \left\{ \sup_{w \in \mathcal{H}_M} \sum_{m=1}^M w_m |e'_W(I - P_m)\mu| > \tau \tilde{\zeta}_n \right\} \\ & = P \left\{ \max_{1 \leq m \leq M} |e'_W(I - P_m)\mu| > \tau \tilde{\zeta}_n \right\} \\ & = P \left\{ \left\{ |\langle e_W, A(w_1^0)\mu \rangle| > \tau \tilde{\zeta}_n \right\} \cup \dots \cup \left\{ |\langle e_W, A(w_M^0)\mu \rangle| > \tau \tilde{\zeta}_n \right\} \right\} \\ & \leq \sum_{m=1}^M P \left\{ |\langle e_W, A(w_m^0)\mu \rangle| > \tau \tilde{\zeta}_n \right\} \\ & \leq \sum_{m=1}^M E \left\{ \frac{\langle e_W, A(w_m^0)\mu \rangle^2}{\tau^2 \tilde{\zeta}_n^2} \right\} \\ & \leq C_1 \tau^{-2J} \tilde{\zeta}_n^{-2J} \sum_{m=1}^M \left\| \Omega(2J)^{1/2} A(w_m^0)\mu \right\|^{2J}, \end{aligned}$$

where $\langle \cdot, \cdot \rangle$ represents an inner product, $A(w) = I - P(w)$. C_1 is a constant, $\Omega(2J) = \text{diag}(\gamma_1^2(2J), \dots, \gamma_n^2(2J))$ and $\gamma_i^2(2J) = E(e_{Wi}^{2J} | X_i)^{1/2J}$. By Lemma A1, $\gamma_i^2(2J) < \infty$; thus, $\lambda_{\max}(\Omega(2J))^J = O(1)$. Hence, combining this with Equation (A4), we have:

$$\begin{aligned} & P \left\{ \sup_{w \in \mathcal{H}_M} |\langle e_W, A(w)\mu \rangle| / R_{G_n}(w) > \tau \right\} \\ & \leq C_1 \tau^{-2J} \tilde{\zeta}_n^{-2J} \lambda_{\max}(\Omega(2J))^J \sum_{m=1}^M \left\| A(w_m^0)\mu \right\|^{2J} \\ & \leq C'_1 \tau^{-2J} \tilde{\zeta}_n^{-2J} \sum_{m=1}^M \left\| A(w_m^0)\mu \right\|^{2J} \\ & \leq C'_1 \tau^{-2J} \tilde{\zeta}_n^{-2J} \sum_{m=1}^M \left(R_{G_n}(w_m^0) \right)^J. \end{aligned}$$

And together with condition (C3), we can prove Equation (A1). Next, we will prove (A2). Similar to the proof of Equation (A1), we have:

$$\begin{aligned} & P \left\{ \sup_{w \in \mathcal{H}_M} |\text{trace}[\Omega P(w)] - \langle e_W, P(w)e_W \rangle| / R_{G_n}(w) > \tau \right\} \\ & \leq \sum_{m=1}^M P \left\{ \left| \text{trace}[\Omega P(w_m^0)] - \langle e_W, P(w_m^0)e_W \rangle \right| > \tau \tilde{\zeta}_n \right\} \\ & \leq \sum_{m=1}^M E \left\{ \frac{[\text{trace}[\Omega P(w_m^0)] - \langle e_W, P(w_m^0)e_W \rangle]^2}{\tau^2 \tilde{\zeta}_n^2} \right\} \\ & \leq C_2 \tau^{-2J} \tilde{\zeta}_n^{-2J} \sum_{m=1}^M \left\{ \text{tr} \left[P(w_m^0)' \Omega(4J) P(w_m^0) \right] \right\}^J, \end{aligned}$$

where C_2 is a constant, $\Omega(4J) = \text{diag}(\gamma_1^2(4J), \dots, \gamma_n^2(4J))$ and $\gamma_i^2(4J) = E(e_{Wi}^{4J}|X_i)^{1/4J}$. By Lemma A1, $\gamma_i^2(4J) < \infty$; thus, $\lambda_{\max}(\Omega(4J))^J = O(1)$. Hence, combining Equation (A5) and condition (C3), we have:

$$\begin{aligned} & P \left\{ \sup_{\mathbf{w} \in \mathcal{H}_M} |\text{trace}[\Omega P(\mathbf{w})] - \langle e_W, P(\mathbf{w})e_W \rangle| / R_{G_n}(\mathbf{w}) > \tau \right\} \\ & \leq C_2 \tau^{-2N} \xi_n^{-2N} \lambda_{\max}[\Omega(4N)]^N \sum_{m=1}^M \text{tr} \left[P(\mathbf{w}_m^0)' P(\mathbf{w}_m^0) \right] \\ & \leq C_2 \tau^{-2J} \xi_n^{-2J} \left(\inf_i e_{Wi}^2 \right)^{-J} \sum_{m=1}^M \left\{ \inf_i e_{Wi}^2 \text{trace}(P^2(\mathbf{w}_m^0)) \right\}^J \\ & \leq C_3 \tau^{-2J} \xi_n^{-2J} \sum_{m=1}^M [R_{G_n}(\mathbf{w}_m^0)]^J = o(1). \end{aligned}$$

Next, we will prove Equation (A3). Note that

$$\begin{aligned} & \left| \frac{L_{\hat{G}_n}(\mathbf{w})}{R_{G_n}(\mathbf{w})} - 1 \right| \\ & = \left| \frac{L_{G_n}(\mathbf{w})}{R_{G_n}(\mathbf{w})} - 1 + \frac{L_{\hat{G}_n}(\mathbf{w}) - L_{G_n}(\mathbf{w})}{R_{G_n}(\mathbf{w})} \right| \\ & \leq \left| \frac{L_{G_n}(\mathbf{w})}{R_{G_n}(\mathbf{w})} - 1 \right| + \left| \frac{\|\mu - \hat{\mu}_{\hat{G}_n}(\mathbf{w})\|^2 - \|\mu - \hat{\mu}(\mathbf{w})\|^2}{R_{G_n}(\mathbf{w})} \right|. \end{aligned} \tag{A9}$$

From Lemma A2, we know that $\left| \frac{L_{G_n}(\mathbf{w})}{R_{G_n}(\mathbf{w})} - 1 \right| = o_p(1)$. Therefore, to prove (A3), it is sufficient to verify that the second part of Equation (A9) converges to 0 in probability.

$$\begin{aligned} & \left| \frac{\|\mu - \hat{\mu}_{\hat{G}_n}(\mathbf{w})\|^2 - \|\mu - \hat{\mu}(\mathbf{w})\|^2}{R_{G_n}(\mathbf{w})} \right| \\ & = \left| \frac{2(\mu - \hat{\mu}(\mathbf{w}))'(\hat{\mu}(\mathbf{w}) - \hat{\mu}_{\hat{G}_n}(\mathbf{w})) + \|\hat{\mu}(\mathbf{w}) - \hat{\mu}_{\hat{G}_n}(\mathbf{w})\|^2}{R_{G_n}(\mathbf{w})} \right| \\ & \leq \frac{2\{L_{G_n}(\mathbf{w})\}^{1/2} \|P(\mathbf{w})(\hat{Y}_W - Y_W)\|}{R_{G_n}(\mathbf{w})} + \frac{\|P(\mathbf{w})(\hat{Y}_W - Y_W)\|^2}{R_{G_n}(\mathbf{w})}. \end{aligned} \tag{A10}$$

According to Lemma A3, we have:

$$\|P(\mathbf{w})(\hat{Y}_W - Y_W)\|^2 \leq \lambda_{\max}(P(\mathbf{w})) \|\hat{Y}_W - Y_W\|^2 = O_p(1).$$

Combining this with Lemma A3, we can conclude that (A9) is of $o_p(1)$, which establishes the proof for (A3). \square

Proof of Theorem 2. It is evident from Equations (13) and (16) that:

$$\hat{C}_{\hat{G}_n}(\mathbf{w}) = C_{\hat{G}_n}(\mathbf{w}) + 2 \text{trace}\{P(\mathbf{w})\hat{\Omega}(\mathbf{w})\} - 2 \text{trace}\{P(\mathbf{w})\Omega\}.$$

In conjunction with Theorem 1, it is evident that to prove Theorem 2, we only need to establish:

$$\sup_{w \in \mathcal{H}_M} \left[\left| \text{trace}\{P(w)\widehat{\Omega}(w)\} - \text{trace}\{P(w)\Omega\} \right| / R_{G_n}(w) \right] = o_p(1). \tag{A11}$$

We denote $Q_m = \text{diag}(\rho_{11}^{(m)}, \dots, \rho_{nn}^{(m)})$ and $Q(w) = \sum_{m=1}^M w_s Q_m$. According to Lemma A1, we have:

$$\lambda_{max}(\Omega) = O(1). \tag{A12}$$

Considering the definition of $\widehat{\Omega}(w)$, and employing proof techniques similar to [10,23], we obtain:

$$\begin{aligned} & \sup_{w \in \mathcal{H}_M} \left[\left| \text{trace}\{P(w)\widehat{\Omega}(w)\} - \text{trace}\{P(w)\Omega\} \right| / R_{G_n}(w) \right] \\ &= \sup_{w \in \mathcal{H}_M} \left[\left| \{\widehat{Y}_W - P(w)\widehat{Y}_W\}' Q(w) \{\widehat{Y}_W - P(w)\widehat{Y}_W\} - \text{trace}\{Q(w)\Omega\} \right| / R_{G_n}(w) \right] \\ &= \sup_{w \in \mathcal{H}_M} \left[\left| \{e_W + \mu - P(w)\widehat{Y}_W\}' Q(w) \{e_W + \mu - P(w)\widehat{Y}_W\} - \text{trace}\{Q(w)\Omega\} \right| / R_{G_n}(w) \right] \\ &\leq \sup_{w \in \mathcal{H}_M} \left[\left| e_W' Q(w) e_W - \text{trace}\{Q(w)\Omega\} \right| / R_{G_n}(w) \right] \\ &\quad + 2 \sup_{w \in \mathcal{H}_M} \left[\left| e_W' Q(w) \{P(w)\widehat{Y}_W - \mu\} \right| / R_{G_n}(w) \right] \\ &\quad + \sup_{w \in \mathcal{H}_M} \left[\left| \{P(w)\widehat{Y}_W - \mu\}' Q(w) \{P(w)\widehat{Y}_W - \mu\} \right| / R_{G_n}(w) \right] \\ &\leq \sup_{w \in \mathcal{H}_M} \left[\left| e_W' Q(w) e_W - \text{trace}\{Q(w)\Omega\} \right| / R_{G_n}(w) \right] \\ &\quad + 2 \sup_{w \in \mathcal{H}_M} \left[\left| e_W' Q(w) \{P(w)\mu - \mu\} \right| / R_{G_n}(w) \right] \\ &\quad + 2 \sup_{w \in \mathcal{H}_M} \left[\left| e_W' Q(w) P(w) e_W - \text{trace}\{Q(w)P(w)\Omega\} \right| / R_{G_n}(w) \right] \\ &\quad + 2 \sup_{w \in \mathcal{H}_M} \left[\left| \text{trace}\{Q(w)P(w)\Omega\} \right| / R_{G_n}(w) \right] \\ &\quad + \sup_{w \in \mathcal{H}_M} \left[\left| \{P(w)\widehat{Y}_W - \mu\}' Q(w) \{P(w)\widehat{Y}_W - \mu\} \right| / R_{G_n}(w) \right] \\ &\equiv T_1 + T_2 + T_3 + T_4 + T_5. \end{aligned}$$

Let $\rho = \max_m \max_i \rho_{ii}^{(m)}$. According to condition (C7), we have:

$$\rho = O\left(n^{-1}\tilde{p}\right). \tag{A13}$$

Given the definition of $R_{G_n}(w)$ and condition (C4), the following equation holds:

$$\begin{aligned} R_{G_n}(w_m^0) &\geq \text{trace}\{P_m \Omega P_m^T\} \geq \epsilon \text{trace}(P_m) = \epsilon p_m, \\ \zeta_n &\rightarrow \infty \quad \text{and} \quad M \zeta_n^{-2J} = o(1). \end{aligned} \tag{A14}$$

From (A6), (A11), (A13), the Chebyshev inequality and Theorem 2 of [22], for any $\tau > 0$, there exist constants c_1 and c_2 such that:

$$\begin{aligned}
 P(T_1 > \tau) &\leq \sum_{m=1}^M P[|e'_W Q_m e_W - \text{trace}(Q_m \Omega)| > \tau \xi_n] \\
 &\leq \tau^{-2J} \xi_n^{-2J} \sum_{m=1}^M E \left\{ e'^T_W Q_m e_W - \text{trace}(Q_m \Omega) \right\}^{2J} \\
 &\leq c_1 \tau^{-2J} \xi_n^{-2J} \sum_{m=1}^M \{ \text{trace} \{ \Omega^{1/2}(4J) Q_m \Omega(4J) Q_m \Omega^{1/2}(4J) \} \}^J \\
 &\leq c_1 \tau^{-2J} \xi_n^{-2J} M \lambda_{\max}^{2J}(\Omega(4J)) \max_{1 \leq m \leq M} \{ \text{trace}(Q_m) \}^J \\
 &= \xi_n^{-2J} M \left\{ O(n^{-1} \tilde{p}^2) \right\}^J = o(1), \tag{A15}
 \end{aligned}$$

$$\begin{aligned}
 P(T_3/2 > \tau) &\leq \sum_{m=1}^M P \{ |e'_W Q_m P_m e_W - \text{trace}(Q_m P_m \Omega)| > \tau \xi_n \} \\
 &\leq \tau^{-2J} \xi_n^{-2J} \sum_{m=1}^M E [e'_W Q_m P_m e_W - \text{trace}(Q_m P_m \Omega)]^{2J} \\
 &\leq c_2 \tau^{-2J} \xi_n^{-2J} \sum_{m=1}^M \text{trace} \{ \Omega^{1/2}(4J) Q_m P_m \Omega(4J) P_m^T Q_m \Omega^{1/2}(4J) \}^J \\
 &\leq c_2 \tau^{-2J} \xi_n^{-2J} M \lambda_{\max}^{2J}(\Omega(4J)) \lambda_{\max}^J(P_m P_m^T) \max_{1 \leq m \leq M} \{ \text{trace}(Q_m^2) \}^J \\
 &= \xi_n^{-2J} M \left\{ O(n^{-1} \tilde{p}^2) \right\}^J = o(1), \tag{A16}
 \end{aligned}$$

$$\begin{aligned}
 T_2/2 &\leq \sup_{\mathbf{w} \in \mathbb{H}_M} \left\{ \|e_W\|^2 \rho^2 \|P(\mathbf{w})\mu - \mu\|^2 / R_{G_n}^2(\mathbf{w}) \right\}^{1/2} \\
 &\leq \|e_W\| \rho \xi_n^{-1/2} = \xi_n^{-1/2} O(n^{-1/2} \tilde{p}) = o(1), \tag{A17}
 \end{aligned}$$

$$\begin{aligned}
 T_4/2 &\leq \xi_n^{-1} \rho \lambda_{\max}(\Omega) \sup_{\mathbf{w} \in \mathcal{H}_M} [\text{trace}\{P(\mathbf{w})\}] \\
 &\leq \xi_n^{-1} \rho \lambda_{\max}(\Omega) \max_m \{ \text{trace}(P_m) \} \\
 &\leq \xi_n^{-1} \rho \lambda_{\max}(\Omega) \max_m \{ \lambda_{\max}(P_m) \} \max_m \{ \text{rank}(P_m) \} \\
 &= \xi_n^{-1} O(n^{-1} \tilde{p}^2) = \xi_n^{-1} O(n^{-1} \tilde{p}^2) = o(1), \tag{A18}
 \end{aligned}$$

$$\begin{aligned}
 T_5 &\leq \rho \sup_{\mathbf{w} \in \mathcal{H}_M} \left[\{P(\mathbf{w})\hat{Y}_W - \mu\}^T \{P(\mathbf{w})\hat{Y}_W - \mu\} / R_{G_n}(\mathbf{w}) \right] \\
 &= \rho \sup_{\mathbf{w} \in \mathcal{H}_M} [L_{G_n}(\mathbf{w}) / R_{G_n}(\mathbf{w})] = O(n^{-1} \tilde{p}). \tag{A19}
 \end{aligned}$$

Therefore, combining (A15)–(A19), along with Condition (C7), it is clear that Theorem 2 holds. □

References

1. Akaike, H. Information Theory and an Extension of the Maximum Likelihood Principle. In *Second International Symposium on Information Theory*; Petrov, B., Csáki, F., Eds.; Akadémiai Kiadó: Budapest, Hungary, 1973; pp. 267–281.
2. Mallows, C.L. Some Comments on Cp. *Technometrics* **1973**, *15*, 661–675.
3. Schwarz, G. Estimating the Dimension of a Model. *Ann. Stat.* **1978**, *6*, 15–18. [\[CrossRef\]](#)

4. Hjort, N.L.; Claeskens, G. Frequentist Model Average Estimators. *J. Am. Stat. Assoc.* **2003**, *98*, 879–899. [[CrossRef](#)]
5. Buckland, S.T.; Burnham, K.P.; Augustin, N.H. Model selection: An integral part of inference. *Biometrics* **1997**, *53*, 603–618. [[CrossRef](#)]
6. Hansen, B.E. Least squares model averaging. *Econometrica* **2007**, *75*, 1175–1189. [[CrossRef](#)]
7. Wan, A.T.; Zhang, X.; Zou, G. Least squares model averaging by Mallows criterion. *J. Econom.* **2010**, *156*, 277–283. [[CrossRef](#)]
8. Hansen, B.E.; Racine, J.S. Jackknife model averaging. *J. Econom.* **2012**, *167*, 38–46. [[CrossRef](#)]
9. Liu, Q.; Okui, R. Heteroscedasticity-robust Cp model averaging. *Econom. J.* **2013**, *16*, 463–472. [[CrossRef](#)]
10. Zhao, S.; Zhang, X.; Gao, Y. Model averaging with averaging covariance matrix. *Econom. Lett.* **2016**, *145*, 214–217. [[CrossRef](#)]
11. Miller, R. Least square regression with censored data. *Biometrika* **1976**, *63*, 449–464. [[CrossRef](#)]
12. Buckley, J.; James, I. Linear regression with censored data. *Biometrika* **1979**, *66*, 429–436. [[CrossRef](#)]
13. Koul, H.; Susarla, V.; Van Ryzin, J. Regression analysis with randomly right-censored data. *Ann. Stat.* **1981**, *9*, 1276–1288. [[CrossRef](#)]
14. He, S.; Huang, X. Central limit theorem of linear regression model under right censorship. *Sci. China Ser. A-Math.* **2003**, *46*, 600–610. [[CrossRef](#)]
15. Wang, Q.; Dinse, G.E. Linear regression analysis of survival data with missing censoring indicators. *Lifetime Data Anal.* **2011**, *17*, 256–279. [[CrossRef](#)]
16. Liang, Z.Q.; Chen, X.L.; Zhou, Y.Q. Mallows model averaging estimation for linear regression model with right censored data. *Acta Math. Appl. Sin. Engl. Ser.* **2022**, *38*, 5–23. [[CrossRef](#)]
17. Wei, Y.; Wang, Q.; Liu, W. Model averaging for linear models with responses missing at random. *Ann. Inst. Stat. Math.* **2020**, *73*, 535–553. [[CrossRef](#)]
18. Liu, Q.; Okui, R.; Yoshimura, A. Generalized least squares model averaging. *Econom. Rev.* **2016**, *35*, 1692–1752. [[CrossRef](#)]
19. Albert, A.; Anderson, J.A. On the existence of maximum likelihood estimates in logistic regression models. *Biometrika* **1984**, *71*, 1–10. [[CrossRef](#)]
20. Zhu, R.; Wan, A.T.; Zhang, X.; Zou, G. A Mallows-type model averaging estimator for the varyingcoefficient partially linear model. *J. Am. Stat. Assoc.* **2019**, *114*, 882–892. [[CrossRef](#)]
21. Dong, Q.K.; Liu, B.X.; Zhao, H. Weighted least squares model averaging for accelerated failure time models. *Comput. Stat. Data Anal.* **2023**, *184*, 107743. [[CrossRef](#)]
22. Whittle, P. Bounds for the moments of linear and quadratic forms in independent variables. *Theory Probab. Appl.* **1960**, *5*, 302–305. [[CrossRef](#)]
23. Qiu, Y.; Wang, W.; Xie, T.; Yu, J.; Zhang, X. Boosting Store Sales Through Machine Learning-Informed Promotional Decisions. 2023. Available online: http://www.mysmu.edu/faculty/yujun/Research/Maml_sales.pdf (accessed on 10 February 2024).

Disclaimer/Publisher’s Note: The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.