

REGRESSION ANALYSIS OF CLUSTERED CURRENT STATUS DATA UNDER THE ADDITIVE HAZARDS MODEL

LIU Yu-huan¹, WANG Cheng-yong²

(*1.School of Mathematics and Statistics, Wuhan University, Wuhan 430072, China*)

(*2.School of Mathematics and Computer Science, Hubei University of Arts and Science,
Xiangyang 441053, China*)

Abstract: In this paper, we discuss regression analysis of clustered current status data under the additive hazards model. Under the situation when the correlated failure times of interest may be related to cluster sizes, by proposing a within-cluster resampling (WCR) method, the limit distribution theory for the corresponding estimators are derived under some regularity conditions. Some simulation studies are conducted to assess the finite-sample behaviors of the estimators.

Keywords: additive hazards model; current status data; within-cluster resampling

2010 MR Subject Classification: 62G05; 62F10; 62F12

Document code: A **Article ID:** 0255-7797(2018)01-0067-08

1 Introduction

Case I interval-censored failure time data or current status data arise in many areas including demographical studies, economics, medical studies, reliability studies and social sciences, see e.g. [1–4]. By case I interval-censored data, we mean that the failure time of interest is not exactly observed but the observation on it is either left- or right-censored. A typical example of such data is given by a tumorigenicity study and in this case, the time to tumor onset is often of interest. However, it is usually not observable as the presence or absence of tumors in animals is usually known only at their death or sacrifice. In particular, clustered current status data are commonly encountered in biomedicine.

Many procedures were developed for regression analysis of interval-censored failure time data under various models. For example, Huang [3] developed the maximum likelihood approach for fitting the proportional hazards model to case I interval-censored data, Chen and Sun [5], Sun and Shen [6] discussed the same problem in the presence of clustering and competing risks, respectively. Hu and Xiang [7] considered the efficient estimation for semiparametric cuer models when one faces case II interval-censored data, Lin et al. [8], Chen

* **Received date:** 2016-07-01

Accepted date: 2016-11-02

Foundation item: Supported by National Natural Science Foundation of China (71371066).

Biography: Liu Yuhuan (1991–), female, born at Taian, Shandong, master degree candidate, major in statistics.

Corresponding author: Wang Chengyong.

and Sun [5] discussed the fitting of the additive hazards model to case I interval-censored data. However, these methods do not take the clustered data into account or assumes that the cluster size is completely random or noninformative, and it is well-known that this may not be true as the outcome of interest among individuals in a cluster may be associated with the size of the cluster. That is, we may have informative cluster sizes. In the following, we present one approach for the problem of the regression analysis of clustered current status data under the additive hazards model.

In the presence of informative cluster size, among others, Dunson et al. [9] proposed a Bayesian procedure that models the relationship between the failure times of interest and the cluster size through a latent variable. Williamson et al. [10] and Cong et al. [11] also considered the same problem and investigated a weighted score function (WSF) approach and a within-cluster resampling (WCR) procedure. However, it does not seem to exist an estimation procedure for regression analysis of clustered failure time data with informative cluster size under the additive hazards model framework and current status data.

The rest of the article is organized as follows. Section 2 proposes the model and some notations used in this paper. Section 3 gives the WCR method by using the inference procedure proposed by Lin et al. [8] under the additive hazards model for case I interval-censored failure time data, and Section 4 presents some extensive simulation studies to assess the performance of the proposed approach.

2 Notation and Model

Let $i = 1, \dots, n$ denote the independent clusters, and $j = 1, \dots, n_i$ denote the subjects within the i -th cluster. For subject j in the i -th cluster, for $i = 1, \dots, n$ and $j = 1, \dots, n_i$, let T_{ij} and C_{ij} denote the failure time of interest and the censoring or observation time, and let $Z_{ij}(t)$ be a p -dimensional vector of covariates that may depend on time t . It is assumed that the T_{ij} may be dependent for the subjects within the same cluster but are independent for subjects from different clusters. We assume that T_{ij} is conditionally independent of C_{ij} given $Z_{ij}(t)$.

We assume that the survival probabilities of individuals in a cluster depend on the size of that cluster. However, it just as noted in Cong et al. [11], the cause for cluster sizes being informative can be complicated and usually unknown, and some latent variables may implicitly affect the baseline hazard for each cluster and/or covariates. If cluster sizes are noninformative to survival, the usual marginal additive hazards model (see [12]) is

$$\lambda_{ij}(t | Z_{ij}) = \lambda_0(t) + \omega_i \beta_0' Z_{ij}(t), \quad (2.1)$$

where β_0 is the unknown vector of p -dimensional regression coefficient, ω_i is the cluster-specific random effect to account for within-cluster correlation in cluster i , and $\lambda_0(t)$ is the unknown baseline hazard function. If cluster sizes are ignorable (noninformative to survival), the usual marginal additive hazards model is applicable, given by

$$\lambda_{ij}(t | Z_{ij}) = \lambda_0(t) + \beta_0' Z_{ij}(t). \quad (2.2)$$

For each (i, j) , we define $N_{ij}(t) = I(C_{ij} \leq \min(t, T_{ij}))$, $\delta_{ij} = I(C_{ij} \leq T_{ij})$ and $Y_{ij}(t) = I(C_{ij} \geq t)$ and let $\lambda_c(t)$ denote the hazard function of the C_{ij} 's. Also define

$$\tilde{\lambda}_{ij}(t|Z_{ij}(s)) = \lambda_c(t) e^{-\Lambda_0(t)} e^{-\beta'_0 Z_{ij}^*(t)} := \lambda_0^c(t) e^{-\beta'_0 Z_{ij}^*(t)},$$

where $\Lambda_0(t) = \int_0^t \lambda_0(s) ds$ and $Z_{ij}^*(t) = \int_0^t Z_{ij}(s) ds$, and

$$M_{ij}(t) = N_{ij}(t) - \int_0^t Y_{ij}(u) \lambda_0^c(u) e^{-\beta'_0 Z_{ij}^*(u)} du.$$

Note that $M_{ij}(t)$ is a local square-integrable martingale with respect to the marginal filtration

$$\mathcal{F}_{ij}(t) = \sigma\{N_{ij}(u), Y_{ij}(u), Z_{ij}(u), 0 \leq u \leq t\}$$

(see Lin et al. [8]), and $\tilde{\lambda}_{ij}(t|Z_{ij}(s))$ satisfies the Cox proportional hazards model. However, due to the within-cluster dependence, $M_{ij}(t)$ is not a martingale with respect to the joint filtration generated by the history of all the failure, censoring and covariate information up to time t .

3 A Method Based on the Within-Cluster Resampling Technique

When cluster sizes are informative, the estimates and inference based on equation (2.2) may be incorrect. To account for informative cluster sizes, this section will propose a method based on the within-cluster resampling (WCR) technique. The basic idea behind the WCR-based procedure is that one observation is randomly sampled with replacement from each of the n clusters using the WCR approach (refer to Hoffman et al. [13]). For this, we randomly sample one subject with replacement from each of the n clusters, and suppose that the resampling process is repeated K times, where K is a large fixed number. Let τ denote a known time for the length of study period, the k -th resampled data set denoted by $\{C_{i,k}, \delta_{i,k}, Z_{i,k}(t); i = 1, \dots, n, 0 \leq t \leq \tau\}$, consists of n independent observations, which can be analyzed using model (2.2) for independent data set. Define $Y_{i,k}(t) = I(C_{i,k} \geq t)$ and $N_{i,k}(t) = \delta_{i,k} I(C_{i,k} \leq t)$, for the k -th resampled data, the partial likelihood function is

$$L_k(\beta) = \prod_{i=1}^n \left(\frac{\exp(-\beta' Z_{i,k}^*(C_{i,k}))}{\sum_{j=1}^n Y_{j,k}(C_{i,k}) \exp(-\beta' Z_{j,k}^*(C_{i,k}))} \right)^{\delta_{i,k}}, \quad (3.1)$$

and the partial likelihood score function and observed information matrix are

$$\begin{aligned} U_k(\beta) &= \sum_{i=1}^n \int_0^\tau \left(Z_{i,k}^*(t) - \frac{S_k^{(1)}(\beta, t)}{S_k^{(0)}(\beta, t)} \right) dN_{i,k}(t), \\ \Sigma_k(\beta) &= \sum_{i=1}^n \int_0^\tau \left(\frac{S_k^{(2)}(\beta, t)}{S_k^{(0)}(\beta, t)} - \left(\frac{S_k^{(1)}(\beta, t)}{S_k^{(0)}(\beta, t)} \right)^{\otimes 2} \right) dN_{i,k}(t), \end{aligned} \quad (3.2)$$

where

$$Z_{i,k}^*(t) = \int_0^t Z_{i,k}(s) ds,$$

$$S_k^{(b)}(\beta, t) = \frac{1}{n} \sum_{j=1}^n Y_{j,k}(t) (Z_{j,k}^*(t))^{\otimes b} e^{-\beta' Z_{j,k}^*(t)},$$

and $a^{\otimes b} = 1, a, aa'$ for $b = 0, 1$ and 2 . The maximum partial likelihood estimator (refer to [14]) $\hat{\beta}^k$ is the solution to $U_k(\beta) = 0$. Furthermore, Lin et al. [8] showed that $\sqrt{n}(\hat{\beta}^k - \beta_0)$ converges in distribution to a zero-mean normal random vector with covariance matrix can be consistently estimated by $n\Sigma_k^{-1}(\hat{\beta}^k)$, and so $\hat{\beta}^k$ is consistent.

As it is known to all that sample mean can reduce the system error, after repeating this procedure K times, the WCR estimator for β_0 can be constructed as the average of the K resample-based estimators, that is,

$$\hat{\beta}_{wcr} = \frac{1}{K} \sum_{k=1}^K \hat{\beta}^k.$$

Under some regularity conditions, we can show that $\sqrt{n}(\hat{\beta}_{wcr} - \beta_0)$ converges in distribution to a zero-mean normal random vector, and the covariance matrix can be consistently estimated by

$$\hat{\Sigma}_{wcr} = \frac{n}{K} \sum_{k=1}^K \Sigma_k^{-1}(\hat{\beta}_{wcr}) - \frac{n}{K} \sum_{k=1}^K (\hat{\beta}^k - \hat{\beta}_{wcr})(\hat{\beta}^k - \hat{\beta}_{wcr})'.$$

The proof of this result is sketched in Appendix. It does not need some special software to implement the proposed method. One can just input the data $\{C_{i,k}, \delta_{i,k}, Z_{i,k}^*(\cdot), i = 1, \dots, n\}$ into standard software for fitting the proportional hazards model with right-censored data.

4 Simulation Study

In this section, we conduct some simulations to assess the finite sample performance of the methods developed in the previous section. In the study, the failure times were generated from model (2.1) with $\lambda_0(\cdot) = 2$. The covariate process was assumed to be time independent for simplicity and generated from the Bernoulli distribution with success probability $p = 0.5$. The censoring times were generated from the exponential distribution with mean $1/\exp(\beta Z_i)$. The cluster sizes were randomly generated from uniform distribution $U\{2, 3, 4, 5, 6, 7\}$ regardless of the frailty values. Here we chose $\beta_0 = \pm 0.5, \pm 0.2$ and 0 . The censoring times were generated from the exponential distribution to achieve approximately 30%, 40%, 50% and 60%.

The results include the estimated bias (Bias) given by the average of the proposed estimates minus the true value, the sample standard deviation (SSE) of the proposed estimates,

the average of the proposed estimates of the standard errors (SEE), and the empirical 95% coverage probabilities (CP). All results listed in the following table are based on 500 replications with the number of clusters $n = 200, 300$ and $K = 500$. It can be seen from Table 1 that the proposed estimate seem to be unbiased, the proposed variance estimates also seem to be reasonable, and all estimates become better when the sample size increases.

Table 1: simulation results for estimates of β_0

cen%	β_0	$n = 200$				$n = 300$			
		BIAS	SSE	SEE	CP	BIAS	SSE	SEE	CP
30	-0.5	-0.0203	0.4416	0.4502	0.956	-0.0063	0.3683	0.3658	0.946
	-0.2	-0.0247	0.4265	0.4104	0.938	-0.0103	0.4165	0.4010	0.938
	0	-0.0051	0.4699	0.4776	0.948	0.0046	0.4312	0.4112	0.946
	0.2	0.0020	0.4718	0.4882	0.954	-0.0056	0.4423	0.4395	0.960
	0.5	-0.0011	0.5043	0.5159	0.950	0.0068	0.3865	0.3827	0.948
40	-0.5	-0.0112	0.4173	0.3923	0.928	-0.0074	0.3064	0.2950	0.944
	-0.2	-0.0164	0.3531	0.3573	0.952	-0.0064	0.3373	0.3349	0.946
	0	0.0026	0.4141	0.4240	0.950	0.0061	0.3456	0.3482	0.958
	0.2	-0.0148	0.4542	0.4358	0.928	0.0112	0.3609	0.3593	0.940
	0.5	-0.0022	0.4807	0.4602	0.952	-0.0043	0.3788	0.3827	0.954
50	-0.5	-0.0102	0.3302	0.3408	0.954	-0.0110	0.2861	0.2646	0.946
	-0.2	-0.0164	0.3531	0.3573	0.952	0.0047	0.2943	0.2898	0.944
	0	-0.0011	0.3658	0.3709	0.952	0.0025	0.3000	0.3057	0.952
	0.2	0.0114	0.3664	0.3853	0.956	0.0106	0.3186	0.3193	0.954
	0.5	0.0062	0.3989	0.4046	0.946	-0.0056	0.3233	0.3399	0.950
60	-0.5	-0.0031	0.3066	0.3063	0.944	-0.0058	0.2423	0.2465	0.950
	-0.2	0.0043	0.3195	0.3291	0.954	0.0036	0.2691	0.2658	0.952
	0	-0.0025	0.3335	0.3441	0.954	0.0015	0.2812	0.2805	0.944
	0.2	0.0085	0.3524	0.3615	0.953	0.0091	0.2882	0.2941	0.952
	0.5	0.0063	0.3999	0.3908	0.938	-0.0067	0.3210	0.3196	0.950

References

- [1] Andersen P K, Gill R D. Cox's regression model for counting processes: a large sample study[J]. Ann. Stat., 1982, 10: 1100–1120.
- [2] Jewell N P, van der Laan M. Generalizations of current status data with applications[J]. Lifetime Data Anal., 1995, 1: 101–110.
- [3] Huang J. Efficient estimation for the proportional hazards model with interval censoring[J]. Ann. Stat., 1996, 24: 540–568.
- [4] Rossini A J, Tsiatis A A. A semiparametric proportional odds regression model for the analysis of current status data[J]. J. Amer. Stat. Assoc., 1996, 91: 713–721.

- [5] Chen L, Sun J. A multiple imputation approach to the analysis of current status data with the additive hazards model[J]. *Comm. Stat. The. Meth.*, 2009, 38: 1009–1018.
- [6] Sun J, Shen J. Efficient estimation for the proportional hazards model with competing risks and current status data[J]. *Canad. J. Stat.*, 2009, 37: 592–606.
- [7] Hu T, Xiang L. Efficient estimation for semiparametric cure models with interval-censored data[J]. *J. Multi. Anal.*, 2013, 121: 139–151.
- [8] Lin D, Oakes D, Ying Z. Additive hazards regression with current status data[J]. *Biom.*, 1998, 85: 289–298.
- [9] Dunson D B, Chen Z, Harry J. Bayesian joint models of cluster size and subunitspecific outcomes[J]. *Biom.*, 2003, 63: 663–672.
- [10] Williamson J, Kim H Y, Manathuga A, Addiss D G. Modeling survival data with informative cluster size[J]. *Stat. Med.*, 2008, 27: 543–555.
- [11] Cong X, Yin G, Shen Y. Marginal analysis of correlated failure time data with informative cluster sizes[J]. *Biom.*, 2007, 63: 663–672.
- [12] Lin D, Ying Z. Semiparametric analysis of the additive risk model[J]. *Biom.*, 1994, 81: 61–71.
- [13] Hoffman E B, Sen P K, Weinberg C R. Within cluster resampling[J]. *Biom.*, 2001, 88: 1121–1134.
- [14] Xiao Z, Zhu Q. Moderate deviation of maximum likelihood estimators for truncated and censored data[J]. *J. Math.*, 2009, 29(3): 273–278.

Appendix: proofs of asymptotic normality of $\hat{\beta}_{wcr}$

We first assume that $1/n \sum_{i=1}^n 1/n_i \sum_{j=1}^{n_i} Y_{ij}(t) e^{-\beta'_0 Z_{ij}^*(t)} Z_{ij}^*(t)$, $1/n \sum_{i=1}^n 1/n_i \sum_{j=1}^{n_i} Y_{ij}(t) e^{-\beta'_0 Z_{ij}^*(t)}$, $1/n \sum_{i=1}^n Y_{i,k}(t) e^{-\beta'_0 Z_{i,k}^*(t)}$ and $1/n \sum_{i=1}^n Y_{i,k}(t) e^{-\beta'_0 Z_{i,k}^*(t)}$ uniformly converge to $\kappa(t)$, $\pi(t)$, $\tilde{\kappa}(t)$ and $\tilde{\pi}(t)$, respectively. For $i = 1, \dots, n; j = 1, \dots, n_i$ and some constant τ , we assume that $P\{Y_{ij}(t) = 1, 0 \leq t \leq \tau\} > 0$, $\int_0^\tau \lambda_c(t) dt < \infty$; $Z_{ij}(t)$ is bounded and the cluster sizes are finite.

Since $\hat{\beta}_k$ is the solution of the estimating equation $U_k(\beta) = 0$, and by the Taylor's expansion, we have

$$-U_k(\beta_0) = U_k(\hat{\beta}_k) - U_k(\beta_0) = \frac{\partial U_k(\beta_\xi)}{\partial \beta_\xi} (\hat{\beta}_k - \beta_0), \quad (5.1)$$

where β_ξ is on the line segment between $\hat{\beta}_k$ and β_0 . Rewriting (5.1) yields that

$$\sqrt{n}(\hat{\beta}_k - \beta_0) = \left(\frac{1}{n} \frac{\partial U_k(\beta_\xi)}{\partial \beta_\xi} \right)^{-1} \left(-\frac{1}{\sqrt{n}} U_k(\beta_0) \right).$$

Note that

$$\begin{aligned} \frac{1}{n} \frac{\partial U_k(\beta)}{\partial \beta} &= \frac{1}{n} \sum_{i=1}^n \int_0^\tau \frac{S_k^{(2)}(\beta, s) - \left(S_k^{(1)}(\beta, s) \right)^{\otimes 2}}{\left(S_k^{(0)}(\beta, s) \right)^2} dN_{i,k}(s) \\ &= \frac{1}{n} \sum_{i=1}^n \int_0^\tau \left(Z_{i,k}^*(s) - \bar{Z}^k(\beta, s) \right)^{\otimes 2} Y_{i,k}(s) e^{-\beta' Z_{i,k}^*(s)} \frac{d\bar{N}_k(s)}{S_k^{(0)}(\beta, s)} \\ &:= A_k(\beta), \end{aligned}$$

where $\bar{Z}^k(\beta, s) = S_k^{(1)}(\beta, s)/S_k^{(0)}(\beta, s)$ and $\bar{N}_k(s) = n^{-1} \sum_{i=1}^n N_{i,k}(s)$. Note that $A_k(\beta)$ is positive definite. Since the K resamples are identically distributed, it can be seen that $A_k(\beta_0)$ converges in probability to a deterministic and positive definite matrix denoted by \mathcal{A}_{wcr} .

Averaging over $k = 1, \dots, K$ resamples, it yields

$$\begin{aligned} \sqrt{n}(\hat{\beta}_{wcr} - \beta_0) &= \frac{1}{K} \sum_{k=1}^K \sqrt{n}(\hat{\beta}^k - \beta_0) = \frac{1}{K} \sum_{k=1}^K A_k(\beta_\xi)^{-1} \frac{-1}{\sqrt{n}} U_k(\beta_0) \\ &= -\mathcal{A}_{wcr}^{-1} \frac{1}{\sqrt{nK}} \sum_{q=1}^K U_k(\beta_0) + o_p(1). \end{aligned}$$

It is sufficient to show that $1/(K\sqrt{n}) \sum_{q=1}^K U_k(\beta_0)$ converges to a normal distribution as $n \rightarrow \infty$, changing the order of summation yields that

$$\begin{aligned} \frac{1}{\sqrt{nK}} \sum_{k=1}^K U_k(\beta_0) &= \frac{1}{\sqrt{n}} \sum_{i=1}^n \frac{1}{K} \sum_{k=1}^K \int_0^\tau (Z_{i,k}^*(t) - \bar{Z}^k(t)) dM_{i,k}(t) \\ &= \frac{1}{\sqrt{n}} \sum_{i=1}^n \frac{1}{K} \sum_{k=1}^K \int_0^\tau \left(Z_{i,k}^*(t) - \frac{\tilde{\kappa}(t)}{\tilde{\pi}(t)} \right) dM_{i,k}(t) + o_p(1) \\ &:= \frac{1}{\sqrt{n}} \sum_{i=1}^n \mathcal{U}_i(\beta_0) + o_p(1), \end{aligned}$$

where $\mathcal{U}_i(\beta_0)$, $i = 1, \dots, n$ are independent with zero mean and finite variance. By the multivariate central limit theorem, $n^{-1/2} K^{-1} \sum_{k=1}^K U_k(\beta_0)$ is asymptotically normal with zero mean and some positive definite covariance matrix. Combining with Slutsky's theorem, $\sqrt{n}(\hat{\beta}_{wcr} - \beta_0)$ converges in distribution to a normal random vector with zero mean and denote the consistent estimator of the covariance matrix by $\hat{\Sigma}_{wcr}$.

To obtain the consistent estimator of the covariance matrix, it is similar to Hoffman et al. [13], we first write

$$\text{var}(\hat{\beta}_k) = E\left(\text{var}(\hat{\beta}_k|\text{data})\right) + \text{var}\left(E(\hat{\beta}_k|\text{data})\right),$$

where the expectations on the right-hand side are over the resampling distribution for $\hat{\beta}_k$ given the data. By the fact of $E(\hat{\beta}_k|\text{data}) = \hat{\beta}_{wcr}$, it yields that

$$\text{var}(\hat{\beta}_{wcr}) = \text{var}(\hat{\beta}_k) - E(\text{var}(\hat{\beta}_k|\text{data})). \quad (5.2)$$

For each resampled data, $\text{var}(\hat{\beta}_k)$ can be consistently estimated by Σ_k . By averaging over the K resamples, the resulting estimator denoted by $K^{-1} \sum_{k=1}^K \hat{\Sigma}_k$ is also consistent. For the second term on the right-hand side of (5.2), since

$$E(\text{var}(\hat{\beta}_k|\text{data})) = E\left(\frac{1}{K} \sum_{k=1}^K (\hat{\beta}_k - \hat{\beta}_{wcr})(\hat{\beta}_k - \hat{\beta}_{wcr})'\right),$$

it can be estimated as the covariance matrix based on the K resamples estimators $\hat{\beta}_k$, that is

$$\Omega = \frac{1}{K} \sum_{k=1}^K (\hat{\beta}_k - \hat{\beta}_{wcr})(\hat{\beta}_k - \hat{\beta}_{wcr})'.$$

Thus the estimated variance-covariance matrix of $\hat{\beta}_{wcr}$ is

$$\tilde{\Sigma}_{wcr} = \frac{1}{K} \sum_{k=1}^K \hat{\Sigma}_k - \frac{1}{K} \sum_{k=1}^K (\hat{\beta}_k - \hat{\beta}_{wcr})(\hat{\beta}_k - \hat{\beta}_{wcr})'.$$

To show the consistency of $\tilde{\Sigma}_{wcr}$, it suffices to show that $\Omega - E(\Omega) \rightarrow 0$ in probability as $n \rightarrow \infty$. Actually, by applying the same arguments as those in the proof of Cong et al. [11], it can be shown that $\tilde{\Omega} - E(\tilde{\Omega}) \rightarrow 0$ in probability as $n \rightarrow \infty$. This completes the proof.

加法风险率模型下聚类的当前状态数据的回归分析

刘玉环¹, 王成勇²

(1. 武汉大学数学与统计学院, 湖北 武汉 430072)

(2. 湖北文理学院数学与计算机科学学院, 湖北 襄阳 441053)

摘要: 本文研究了加法风险率模型下聚类的当前状态数据 (I型区间删失数据) 的回归分析问题. 在相关的失效时间数据与簇类的规模有关的情形下, 本文提出了一个簇内再抽样方法, 并在一些正则条件下给出了相应估计量的极限分布理论. 最后通过模拟实验验证了估计量的有限样本行为.

关键词: 加法风险率模型; 当前状态数据; 簇内再抽样

MR(2010)主题分类号: 62G05; 62F10; 62F12 中图分类号: O212.1