

June 2022

Nonparametric Estimation of Transition Probabilities in Illness-Death Model based on Ranked Set Sampling

Ying Ma

University of South Florida

Follow this and additional works at: <https://digitalcommons.usf.edu/etd>

 Part of the [Biostatistics Commons](#), and the [Epidemiology Commons](#)

Scholar Commons Citation

Ma, Ying, "Nonparametric Estimation of Transition Probabilities in Illness-Death Model based on Ranked Set Sampling" (2022). *USF Tampa Graduate Theses and Dissertations*.
<https://digitalcommons.usf.edu/etd/9403>

This Dissertation is brought to you for free and open access by the USF Graduate Theses and Dissertations at Digital Commons @ University of South Florida. It has been accepted for inclusion in USF Tampa Graduate Theses and Dissertations by an authorized administrator of Digital Commons @ University of South Florida. For more information, please contact scholarcommons@usf.edu.

Nonparametric Estimation of Transition Probabilities in Illness-Death Model

Based on Ranked Set Sampling

by

Ying Ma

A dissertation submitted in partial fulfillment
of the requirements for the degree of
Doctor of Philosophy
with a concentration in Biostatistics
College of Public Health
University of South Florida

Major Professor: Henian Chen, M.D., Ph.D.
Jason Beckstead, Ph.D.
Getachew A. Dagne, Ph.D.
Wei Wang, Ph.D.
Ronee Wilson, Ph.D.

Date of Approval:
June 22, 2022

Keywords: Aalen-Johansen estimator, Markov process, censored data,
small for gestational age, colon cancer.

Copyright © 2022, Ying Ma

TABLE OF CONTENTS

List of Tables	iii
List of Figures	vi
Abstract	viii
1 Introduction.....	1
1.1 Simple Random Sampling	1
1.2 Ranked Set Sampling	2
1.2.1 Advantages of RSS	2
1.2.2 Application of RSS	3
1.2.3 RSS and Censored Data	4
1.3 Markov Process.....	6
1.4 Illness-Death Model.....	7
2 Chapter I: RSS Modified Aalen-Johansen Estimator for Transition Probabilities in Illness-Death Model: A Simulation Study	9
2.1 Ranked Set Sampling (RSS)	9
2.2 Illness-Death Model.....	9
2.3 Proposed Estimators.....	10
2.4 Simulation	12
2.5 Simulation Results	14
2.6 Discussion	18
3 Chapter II: Application of Aalen-Johansen Estimator in Healthy Start Dataset.....	50
3.1 Introduction.....	50
3.1.1 Small for Gestational Age.....	50
3.1.2 Healthy Start Project.....	51
3.1.3 Illness-Death Model Applied to Health Start Project for SGA Outcome	53
3.2 Methods	54
3.2.1 A Markov Model to Estimate Transition Probabilities for Healthy Start Study and Small for Gestational Age	54
3.2.2 Statistical Analysis.....	54
3.2.3 Population Examined in the Study.....	56
3.2.4 Setting	56
3.3 Results and Discussion	58
3.3.1 Results.....	58
3.3.2 Discussion.....	60
4 Chapter III: Application of Aalen-Johansen Estimator Based on RSS Design to Colon Cancer Dataset	78
4.1 Data Description	78
4.2 Methods	79

4.2.1	An RSS Modified Aalen Johansen Estimator for Colon Cancer Dataset	79
4.2.2	Statistical Analysis.....	80
4.3	Results.....	81
4.4	Discussion.....	82
5	Limitations	86
6	Contributions.....	87
	References.....	88

LIST OF TABLES

Table 1:	MSEs of Transition probabilities when sample size is 200 (set size: 2, cycle number: 100).....	21
Table 2:	MSEs of Transition probabilities when sample size is 200 (set size: 4, cycle number: 50).....	22
Table 3:	MSEs of Transition probabilities when sample size is 400 (set size: 2, cycle number: 200).....	23
Table 4:	MSEs of Transition probabilities when sample size is 400 (set size: 4, cycle number: 100).....	24
Table 5:	MSEs of Transition probabilities when sample size is 400 (set size: 8, cycle number: 500).....	25
Table 6:	MSEs of Transition probabilities when sample size is 800 (set size: 2, cycle number: 400).....	26
Table 7:	MSEs of Transition probabilities when sample size is 800 (set size: 4, cycle number: 200).....	27
Table 8:	MSEs of Transition probabilities when sample size is 800 (set size: 8, cycle number: 100).....	28
Table 9:	MSEs of Transition probabilities when sample size is 800 (set size: 16, cycle number: 50).....	29
Table 10:	MSEs of Transition probabilities when sample size is 1600 (set size: 4, cycle number: 400).....	30
Table 11:	MSEs of Transition probabilities when sample size is 1600 (set size: 8, cycle number: 200).....	31
Table 12:	MSEs of Transition probabilities when sample size is 1600 (set size: 16, cycle number: 100)	32
Table 13:	MSEs of Transition probabilities when sample size is 1600 (set size: 32, cycle number: 50)	33
Table 14:	Estimated MSEs of distribution function estimators at some percentiles for $n = 200$, the first censoring level	34
Table 15:	Estimated MSEs of distribution function estimators at some percentiles for $n = 400$,	

	the first censoring level	34
Table 16:	Estimated MSEs of distribution function estimators at some percentiles for $n = 800$, the first censoring level	35
Table 17:	Estimated MSEs of distribution function estimators at some percentiles for $n = 1600$, the first censoring level	35
Table 18:	Estimated MSEs of distribution function estimators at some percentiles for $n = 200$, the second censoring level.....	36
Table 19:	Estimated MSEs of distribution function estimators at some percentiles for $n = 400$, the second censoring level.....	36
Table 20:	Estimated MSEs of distribution function estimators at some percentiles for $n = 800$, the second censoring level.....	37
Table 21:	Estimated MSEs of distribution function estimators at some percentiles for $n = 1600$, the second censoring level.....	37
Table 22:	Estimated bias of distribution function estimators at some percentiles for $n = 200$, the first censoring level	38
Table 23:	Estimated bias of distribution function estimators at some percentiles for $n = 400$, the first censoring level	38
Table 24:	Estimated bias of distribution function estimators at some percentiles for $n = 800$, the first censoring level	39
Table 25:	Estimated bias of distribution function estimators at some percentiles for $n = 1600$, the first censoring level	39
Table 26:	Estimated bias of distribution function estimators at some percentiles for $n = 200$, the second censoring level.....	40
Table 27:	Estimated bias of distribution function estimators at some percentiles for $n = 400$, the second censoring level.....	40
Table 28:	Estimated bias of distribution function estimators at some percentiles for $n = 800$, the second censoring level.....	41
Table 29:	Estimated bias of distribution function estimators at some percentiles for $n = 1600$, the second censoring level.....	41
Table 30:	Estimated MSEs of transition probabilities P_{11}, P_{12} at some percentiles for $n = 200$, the second censoring level.....	42
Table 31:	Estimated bias of transition probabilities P_{11}, P_{12} at some percentiles for $n = 200$, the second censoring level.....	42
Table 32:	Estimated MSEs of transition probabilities P_{11}, P_{12} at some percentiles for $n = 400$,	

	the second censoring level.....	43
Table 33:	Estimated bias of transition probabilities P_{11} , P_{12} at some percentiles for $n = 400$, the second censoring level.....	43
Table 34:	Number of records for percentage of women choosing services and percentage of newborns with SGA in healthy start data set.....	61
Table 35:	Summary of demographic statistics in healthy start data set.....	62
Table 36:	Transition Probability at 2, 3, 4, 5, 6, 7, 8, 9, and last day from pregnancy in healthy start dataset.....	63
Table 37:	Transition probabilities by treatment groups (%) in healthy start dataset.....	64
Table 38:	Correlation coefficients and p values of linear regression between survival time and various variables in colonTP dataset.....	83
Table 39:	Squared errors of distribution function estimators at some percentiles for $n = 200$ in colonTP dataset.....	84
Table 40:	Percentage of RSS estimators of the real transition probabilities in colonTP dataset with a 200 sample size.....	85

LIST OF FIGURES

Figure 1:	An RSS procedure of obtaining a sample size of $k \times m$	44
Figure 2:	Illustration of progressive illness-death model.....	45
Figure 3:	The efficiency of the AJ estimators based on RSS with respect to SRS counterparts at different percentiles for $n = 200, 400$, the first censoring level.....	46
Figure 4:	The efficiency of the AJ estimators based on RSS with respect to SRS counterparts at different percentiles for $n = 800, 1600$, the first censoring level.....	47
Figure 5:	The efficiency of the AJ estimators based on RSS with respect to SRS counterparts at different percentiles for $n = 200, 400$, the second censoring level.....	48
Figure 6:	The efficiency of the AJ estimators based on RSS with respect to SRS counterparts at different percentiles for $n = 800, 1600$, the second censoring level.....	49
Figure 7:	Illustration of progressive illness-death model applied in Healthy Start project.....	65
Figure 8:	Transition probabilities of having an SGA infant by risk groups in healthy start data set	65
Figure 9:	Transition probabilities P00, P01 and P02 plots for teen mothers in healthy start data set	66
Figure 10:	Transition probabilities P11 and P12 plots for teen mothers in healthy start data set.....	67
Figure 11:	Transition probabilities P00, P01 and P02 plots for marriage groups in healthy start data set.....	68
Figure 12:	Transition probabilities P11 and P12 plots for marriage groups in healthy start data set	69
Figure 13:	Transition probabilities P00, P01 and P02 plots for race groups in healthy start data data set.....	70
Figure 14:	Transition probabilities P11 and P12 plots for race groups in healthy start data set.....	71
Figure 15:	Transition probabilities P00, P01 and P02 plots for obese status in healthy start data set	72
Figure 16:	Transition probabilities P11 and P12 plots for obese status in healthy start data set	73
Figure 17:	Transition probabilities P00, P01 and P02 plots for smoking groups in healthy start	

data set.....	74
Figure 18: Transition probabilities P11 and P12 plots for smoking groups in healthy start data set	75
Figure 19: Transition probabilities P00, P01 and P02 plots for education groups in healthy start data set.....	76
Figure 20: Transition probabilities P11 and P12 plots for education groups in healthy start data set	77

ABSTRACT

The ranked set sampling (RSS) design is applied widely in agriculture, environmental science, and medical research where the exact measurements of sampling units is costly, but sampling units can be ranked by a correlated concomitant variable. RSS is usually a cost-efficient alternate to simple random sampling (SRS) for selecting more representative samples. This study presents a novel methodology to investigate the nonparametric estimation of transition probabilities in illness-death model using the RSS design. We study the Aalen–Johansen estimator of transition probabilities in illness-death Markov model based on RSS design under random right censoring time and propose nonparametric estimators of the transition probabilities. We compare the performance of the suggested estimators with their SRS counterparts via simulation study, in which two censoring levels are considered. Our results show that the proposed estimator under RSS design outperforms its competitors in SRS design in many simulation scenarios. When sample size is big with the highest set number, the proposed estimator performs the best. Conventional and RSS modified Aalen Johansen estimators are applied to healthy start project and colon cancer dataset correspondingly for illustration. The Aalen-Johansen estimator under RSS design possesses higher efficiency as compared with its SRS competitor from simulation study and real research datasets.

1 Introduction

1.1 Simple Random Sampling

As we know that performing a study for an entire population is expensive and unnecessary. Sometimes it is considered impossible. Sampling technologies were developed to draw inferences about the target population from a sample with reduced costs.^{1,2} In order to produce reliable conclusions about the reference population, samples need to contain necessary sample size and sufficient representativeness.³ The most challenging aspect in reality is obtaining representative samples because it is essential in terms of generalizability.² Probability sampling is a sampling technique which gives every unit in the target population a known and nonzero chance of being selected. Compared with nonprobability sampling, probability sampling usually provides more representative samples, less selection bias and statistical inferences to the population.⁴ Simple random sampling (SRS) is the most fundamental and recognized procedure of selecting samples, in which each unit in the sampling pool has equal probability of being selected.⁵ If a population is denoted as Ω , in SRS, sample s of size n is selected unit by unit with or without replacement.⁵ There are two prerequisites for simple random sampling to perfectly carry out: the whole population is accessible to the researchers and the researchers have a list of all individuals from the population.⁶ Researchers can use computer program, lottery method or a table of random numbers to generate the random sample.⁷ SRS design has certain advantages, such as high internal and external validity, simple and straightforward method to analyze data, minimal knowledge needed for the population. All of the above make SRS design the most widely applied technique in observational studies.⁸ Other strengths of SRS include: every sampling unit has an independent probability of being selected; it is a facile method to understand and communicate; most statistical software have incorporated procedures to handle the inferential statistics.⁷

1.2 Ranked Set Sampling

Ranked Set Sampling (RSS) was first introduced by McIntyre in 1952 to improve the accuracy of estimating pasture.⁹ Agricultural and environmental monitoring data are usually observational.¹⁰ Since laboratory analysis of these units are expensive, it is better to obtain such data with representative samples of the population. The philosophy of RSS is that promising observational economy could be achieved if we could identify a large number of samples representing the population, yet only carefully select subsamples to be measured.¹⁰ In 1966, RSS was applied to estimate weights and forage yields, which showed considerably more efficiency than SRS.¹¹ However, when population pool is large and sample size is small, samples selected by SRS may not represent the population well and have reduced power for a given significant level.¹² In this situation, the uncertainty of estimating population mean is increased. To solve this problem, normally scientists increase sample size.^{10,12,13} However, the study economy is compromised after increasing sample size.

If we carry out the study with RSS of same sample size, more representative samples could be obtained via an underlying ranking process. When the outcome variable is expensive to measure, we could use a less costly covariate which is correlated with the outcome to rank a random sample. Then pick a unit according to its rank to perform full measurement and discard the rest. Repeat this procedure multiple times. The difference between this RSS process and previous SRS design is that though they both have same sample size, the number of subjects who participated the procedure is much higher for RSS. There is a silent stratification process during the RSS procedure, which provides underlying information and finally helps to obtain more representative samples to span the full range of values in the population.^{10,14}

1.2.1 Advantages of RSS

RSS is a sampling design technique to obtain more representative units from the population where measurement of the units is costly or time-consuming.¹⁵ This sampling technique is powerful, cost-effective and efficient.^{16,17} The improved efficiency is a consequence of additional information provided

by ranked but not measured units.¹⁷ What we randomly sampled are the subpopulations with relatively low, medium and high distribution. These subpopulations are formed without construction real strata but by ranking.¹⁰ Each subpopulation has its own distribution. These subpopulations' distributions comprise the parent distribution. Even though both RSS and SRS obtain a sample size of $k \times m$ (k is set number, m is cycle number.), the number of units really participated in an RSS procedure is $k^2 \times m$. However, only $k \times m$ units are measured. Through this process of randomly sampling and ranking at the same time, RSS obtained more regularly spaced samples than SRS.

Both McIntyre and Dell have indicated that RSS could provide unbiased estimator of population mean.^{9,18} Estimator of population mean provided by RSS is at least as precise as SRS.^{18,19} RSS not only provides unbiased estimation of population mean but also much smaller confidence interval.^{10,20,21} The sampling efficiency comparing RSS versus SRS is expressed as the relative precision (RP).

$$RP = \frac{\text{variance of sample average with SRS}}{\text{variance of sample average with RSS}}$$

It is proved that $1 \leq RP \leq \frac{k+1}{2}$.^{9,10} k is the number of ranks (set number). Since RP cannot be less than 1, RSS is always equal to or more efficient than SRS. Not only for the unbiased mean, theoretical investigation by Stokes showed that an estimator of variance provided by RSS is asymptotically unbiased regardless of ranking errors.²² Stokes also indicated that RSS could generate more precise correlation coefficient compared to SRS.²³ Stokes and Sager proved that empirical distribution function of a RSS is unbiased and has greater precision than SRS.²¹

1.2.2 Application of RSS

In 1966, RSS was firstly applied by L. K. Halls and T. R. Dell to estimate weights of browse and herbage in a pine forest in Texas and RSS was considered more efficient than SRS.¹¹ W. L. Martin et al., moved from applying RSS for estimating forage and pasture yields to shrub phytomass in a Appalachian oak forest in Virginia. The RSS provided both closer mean to the population and smaller variance than SRS.²⁴ J. M. Cobby et al., used RSS to estimate herbage mass clover contents in grazed swards.²⁵

Besides forestry and herbage, RSS is widely used in many other fields, such as agriculture and environmental monitoring.¹⁸ L. E. Nelson et al., used RSS technique to estimate annual maximum and minimum standing crop production of *Populus deltoides* plantations in the Mississippi River Valley.²⁶ N. A. Mode et al., applied RSS to measure stream habitat data in Pacific Northwest, which is correlated with salmon production.^{27,28} Mode mentioned that RSS is more efficient than SRS with same sample size even after accounting for cost of ranking. As set number increases, the precision of RSS increases.²⁷ Due to extensive sampling effort involved in studying the effect of spray deposit on the leaves of apple trees, R. A. Murray et al., applied RSS in the assessment of total deposit, which was found to be more efficient than SRS.²⁹ M. F. Al-Saleh et al., applied RSS to estimate average sheep milk yield and compared it with SRS. The relative saving of sampling units using RSS is between 32% and 44% compared with SRS to obtain the same precision of the estimator.³⁰ Omer Ozturk et al.'s study showed that RSS has a substantial improvement in estimating both mean and variance of seven month sheep weights at a research farm in Erzurum, Turkey.¹⁹ The authors pointed out that the RSS mean estimator is unbiased regardless of the accuracy of ranking and the population size. If the judgement ranking is imperfect, in the worst case, RSS estimator will be equivalent to SRS estimator.¹⁹ Therefore, there is nothing lost to use RSS where it is applicable.

In 2003, Paul H. Kvam applied RSS into stream water quality data from the National Stream Quality Accounting Network station on a river near Fredricksburg, VA. This is the first time RSS used in a binary outcome data.³¹ The author affirmed that if the success probability could be ranked from correlated covariates, RSS is able to be applied in binary outcome data. In his example, RSS showed superiority than SRS in both estimated precision and smaller confidence interval.³¹

1.2.3 RSS and Censored Data

Yu and Tam applied RSS to estimate the population mean and standard deviation for censored data with lognormal maximum likelihood and Kaplan-Meier methods.³² This is the first study to investigate RSS in censored data.³² The results were compared with the corresponding SRS estimators. They

considered four censoring levels (10th, 20th, 30th percentiles, no censoring) and two sample sizes (24, 120, RSS used different set size and cycle number combination to achieve same SRS sample size). Their results indicate that for all the censoring levels RSS estimators have smaller bias and mean squared error (MSE) than corresponding SRS estimators.³² For both RSS and SRS, as censoring levels increase, bias and MSE of estimators increase.³² Larger sample size indicates smaller estimator bias and MSE for both RSS and SRS sampling methods.³² When sample size is fixed, estimator bias and MSE decrease as set number increases for RSS sampling in all censoring levels.³²

Strzalkowska-Kominiak and Mahdizadeh modified Kaplan-Meier estimator based on RSS and compared it with traditional KM estimator under SRS.³³ Under all censoring levels (0.1, 0.2, 0.3) and sampling sizes (n=36, 60, 120, 240), RSS estimators showed superiority over SRS estimators in MSE. As sample size increases, MSE decreases for both RSS and SRS estimators. When sample size and censoring level are fixed, MSE is decreasing as set size increases for RSS estimators. However, they observed that as sample size increases the efficiency gain from RSS over SRS decreases.³³

Nematolahi et al., improved Kaplan-Meier estimator from Partially Rank-Ordered Set (PROS) Samples and they compared the estimator with RSS and SRS correspondingly.³⁴ In addition to RSS, PROS considers subsets. Instead of ranking elements within the set, it ranks subsets. An element is picked randomly from the lowest ranked subset to the highest one. In this way, it allows certain flexibility compared with RSS. PROS is equivalent to RSS when subset element number is 1. In their study, they considered two censoring levels (0.1 and 0.6). They indicate that with same sample size, KM estimators from PROS sampling method is more efficient than from RSS and SRS designs. The advantage of PROS over SRS sampling is much more obvious than over RSS method in both censoring levels, which could be as high as 3 times better than the SRS sampling. When sample size and censoring level are fix, efficiency of PROS estimator is improved with increased set size. Regardless of censoring level and ranking error, as sample size increases the efficiency of PROC estimator is enhanced.³⁴

The mean residual life (MRL) is defined as the expected additional lifetime given that a component has survived until time t . Estimations and properties of MRL were thoroughly

$$M(t) = E(X - t|X > t) = \frac{\int_t^{\infty} S(x)dx}{S(t)}$$

Studied based on simple random sampling. Zamanzade et al., estimated MRL based on RSS and compared it with corresponding RSS estimators.³⁵ Their simulation results showed that the MRL estimator in RSS setting is more efficient than in the SRS sampling for certain distributions. The confidence intervals based on RSS setting is narrower than its SRS counterparts.³⁵ Accelerated Failure Time (AFT) is a parametric model to provide a linear relationship between log of failure time and covariates. Samawi et al., proved from simulation study that using a Moving Extreme Ranked Set Sampling (MERSS) or an Extreme Ranked Set Sampling (ERSS) could improve power testing and hazard ratio testing during the procedure compared with SRS.³⁶ Discrete time survival analysis has a different data structure from continuous data model, in which every subject has multiple rows of data depending on the number of discrete times of following. Therefore, researchers turn to delicate sampling design to improve the extended time of analysis due to the enlarged dataset. Tutkun et al., performed a simulation study regarding discrete time survival analysis under RSS. Their results indicate that RSS has improved efficiency compared to SRS in three samples sizes (30, 60, 100) and two censoring rates (10%, 60%).³⁷

1.3 Markov Process

The Markov models are a class of stochastic models that assume a finite number of health states (clusters) and allows movement or transition from one state to the other.³⁸ The rate of movement from one state to the next are measured in terms of transition probabilities.³⁹ The transition probability from state a to state b ($0 \leq s < t$) is represented mathematically as

$$p_{ab}(s, t) = (P(Y(t) = b|Y(s) = a, H_{s-}))$$

H_{s-} represents all the historical information from the data along the interval $[0, s)$. The model in the above equation is assumed to be independent on H_{s-} . Any future evolution of the Markov process depends only on its current state and is independent of the previously visited states. This is the memoryless property which must be satisfied by all Markov models.⁴⁰

1.4 Illness-Death Model

The Markov model investigated in this dissertation is progressive illness-death model (also known as disability model).⁴¹ The model depicted in Figure 1 has initial state 0 “health”, intermediate state 1 “illness” and final absorption state 2 “death”. It assumes that individuals from an initial state may transit to an intermediate state and may finally enter a terminal state. Individuals also have the possibility to bypass the intermediate state and enter the final absorption state directly. Subjects could only stay in their current state or move forward. Illness-death model is one of the most popular multi-state Markov models.⁴²

This model has been applied in medical research to study the progression of human disease for a long time.^{43,44} Commenges et al., suggested that illness-death model is a better choice than survival model for studying the prevalence of dementia, since this model is able to provide the age specific incidence of dementia and mortality rate simultaneously.⁴⁵ Harezlak et al., used illness-death model to estimate the transition hazard parameter to multiple states in a longitudinal dementia study.⁴⁶ Frydman and Szarek applied illness-death model to infants born to HIV infected women for HIV positivity, HIV free survival and overall survival.⁴⁷ Moreover, the illness-death model has also been utilized to analyze liver cirrhosis,⁴⁸ bone marrow transplantation⁴⁹ and diabetes data.⁵⁰

Transition probabilities in illness-death model provide the probabilities of transition from one state to another. Odd Aalen proposed using counting process to estimate the cumulative transition intensities by Nelson-Aalen estimator in her doctoral dissertation.^{51,52} However cumulative transition intensities is difficult to interpret in medical or observational research, transition probabilities is needed to be developed. Aalen and Johansen made a breakthrough by generalizing the Kaplan-Meier estimator to Markov process to estimate the transition probabilities with a nonparametric method.^{53,54} Based on the inverse probability of censoring weighting (IPCW) principle, Datta and Satten proposed an estimator to extend the Aalen-Johansen estimators to data with dependent censoring.⁵⁵ Meira-Machado et al., proposed non-Markovian estimator based on the landmark methodology to substitute Aalen-Johansen estimator in 2006.^{56,57} Later they modified this non-Markovian estimator based on presmoothing

method.⁵⁸ However, this paper focus on conventional Aalen-Johansen estimator to estimator the transition probabilities in illness-death model.

Though there are extensive studies regarding RSS and survival analysis, all those studies have discussed only one event and censoring level. There is no reported investigation about RSS and data having more than one events. In this study, we present a novel methodology to investigate the nonparametric estimation of transition probabilities in illness-death Markov model using the RSS design. The Aalen–Johansen estimator of transition probabilities in illness-death model based on RSS under random right censoring is proposed. The performance of suggested estimator was compared with its counterparts under SRS design via simulation study. The conventional and RSS modified Aalen-Johansen estimators are applied to the healthy start data set and a real-word colon cancer data set respectively for illustration.

2 Chapter I: RSS modified Aalen–Johansen estimator of transition probabilities in illness-death Model: A simulation study

2.1 Ranked Set Sampling (RSS)

A preliminary condition for RSS is that a set of units drawn from population could be ranked by an uncovariate measured covariate without fully measurement of the outcome variable.⁵⁹ This covariate needs to be correlated with the interested outcome. An RSS sampling procedure with set number k and cycle number m is illustrated in figure 1. First, a set with sample size k is randomly drawn from the population. These k units are ranked with respect to a covariate, X . The unit with the lowest rank is taken for full measurement and the remaining units in the sample are discarded. Next, another set of sample size k is randomly drawn. These k units are ranked with respect to X . The unit with the second lowest rank is taken for full measurement and the remaining units in the sample are discarded. This process continues until the k^{th} set with sample size k is randomly drawn and ranked with respect to X . The unit with the highest rank is taken for full measurement. Until now, k units are selected. This process is referred to as a cycle. In figure 1 selected units are marked by dark blue color. The cycle can repeat multiple times, refer to as m . An RSS sample size is defined as $N = k \times m$.

2.2 Illness-Death Model

For two states a, b and two time points $s < t$, there are transition probabilities

$$p_{ab}(s, t) = P(Y(t) = b | Y(s) = a)$$

As indicated in figure 2, illness-death model has three states: State 0, the disease-free state; State 1, the diseased state; and State 2, the absorbing state or dead state. There are five transition probabilities in the model: $p_{00}(s, t)$, $p_{01}(s, t)$, $p_{02}(s, t)$, $p_{11}(s, t)$, and $p_{12}(s, t)$.

In the above stochastic process, there is a random variable T_{ab} ($0 \leq a \leq b \leq 2$), which represents the transition time from state a to state b . There are two conditions in this model: subjects visit intermediate state 1 or not. For subjects who do not visit intermediate state, there is a random variable T_{02} (time from state 0 to state 2). For those who visit state 1, there are two more random variables: T_{01} (time from state 0 to state 1), T_{12} (time from state 1 to state 2) with $T_{02} = T_{01} + T_{12}$. $Z = \min(T_{01}, T_{02})$ is the sojourn time in state 0. $\sigma = I(T_{01} < T_{02})$ is the indicator of visiting state 1. Therefore, $T = Z + \sigma T_{12}$ is the total survival time. For those who went through intermediate state 1, $\sigma = I(T_{01} < T_{02}) = 1$ and $T = Z + T_{12}$. For subjects who bypass state 1, $\sigma = I(T_{01} < T_{02}) = 0$ and $T = Z$. $Z < T$ indicates subjects visit the intermediate state.

The transition probabilities can be expressed as depending on the joint distribution of (Z, T) as following:

$$\begin{aligned}
p_{00}(s, t) &= P(Z > t | Z > s) \\
p_{01}(s, t) &= P(Z \leq t, T > t | Z > s) \\
p_{02}(s, t) &= P(T \leq t | Z > s) \\
p_{11}(s, t) &= P(Z \leq t, T > t | Z \leq s, T > s) \\
p_{12}(s, t) &= P(T \leq t | Z \leq s, T > s)
\end{aligned}$$

2.3 Proposed Estimators

Let C be a censoring variable, which is independent of Z and T . We introduce $\delta_0 = I(Z \leq C)$ to be the censoring indicator of Z and $\delta_1 = I(T \leq C)$ to be the censoring indicator of T . For a censored variable Z , we define $\tilde{Z} = \min(Z, C)$. For a censored variable T , we define $\tilde{T} = \min(T, C)$. For a censored version of visiting intermediate state 1 or not, we introduce indicator $\beta = I(\tilde{Z} < \tilde{T})$. For an observed subject, it he/she visits state 1, $\beta = I(\tilde{Z} < \tilde{T}) = 1$, otherwise $\beta = 0$. After considering censoring in practice, \tilde{Z} is the observed sojourn time in state 0. $\tilde{T}_{12} = \tilde{T} - \tilde{Z}$ is the observed sojourn time in state 1. The available data of illness death model via RSS design with set number k and cycle number m ($RSS(k, m)$) are $(\tilde{Z}_l, \tilde{T}_l, \delta_{0l}, \delta_{1l}, \beta_l), 1 \leq l \leq m, i. i. d.$

The proposed Aalen–Johansen (AJ) estimator for transition probabilities based on the $RSS(k, m)$ sampling in illness death model $(\tilde{Z}, \tilde{T}, \delta_0, \delta_1, \beta)$ is defined as the following.

For transition probability from state 0 to state 0, we have

$$\hat{P}_{00RSS}^{AJ}(s, t) = \frac{1}{k} \sum_{r=1}^k \hat{P}_{00[r]}(s, t),$$

$$\hat{P}_{00[r]}(s, t) = \prod_{s < \tilde{Z}_l \leq t, l=1}^m \left(1 - \frac{\delta_{0l}}{R_0(\tilde{Z}_l)}\right)^{1_{\{s < Y_{[r]l}^* \leq t\}}} \quad (r = 1, \dots, k),$$

where $Y_{[r]1}^*, \dots, Y_{[r]m}^*$ are ordered values of the units of the r^{th} rank and $R_0(t) = \sum_{l=1}^m I(\tilde{Z}_l \geq t)$.

For transition probability from state 1 to state 1, we have

$$\hat{P}_{11RSS}^{AJ}(s, t) = \frac{1}{k} \sum_{r=1}^k \hat{P}_{11[r]}(s, t),$$

$$\hat{P}_{11[r]}(s, t) = \prod_{s < \tilde{T}_l \leq t, \beta_l=1, l=1}^m \left(1 - \frac{\delta_{1l}}{R_1(\tilde{T}_l)}\right)^{1_{\{s < Y_{[r]l}^* \leq t\}}} \quad (r = 1, \dots, k),$$

where $R_1(t) = \sum_{l=1}^m I(\tilde{Z}_l < t \leq \tilde{T}_l)$. Then the modified transition probability from state 0 to state 1 is proposed as

$$\hat{P}_{01RSS}^{AJ}(s, t) = \frac{1}{k} \sum_{r=1}^k \hat{P}_{01[r]}(s, t),$$

$$\hat{P}_{01[r]}(s, t) = \sum_{l=1}^m \hat{P}_{00[r]}(s, \tilde{Z}_l^-) \hat{P}_{11[r]}(\tilde{Z}_l, t) I(s < \tilde{Z}_l \leq t) \frac{\beta_l}{R_0(\tilde{Z}_l)} \quad (r = 1, \dots, k),$$

Finally, it is obvious to propose the following transition probabilities from state 0 to state 2, and from state 1 to state 2. Since in Aalen–Johansen (AJ) transition probabilities, $P_{00} + P_{01} + P_{02} = 1$ and $P_{11} + P_{12} = 1$.

$$\hat{P}_{02RSS}^{AJ}(s, t) = \frac{1}{k} \sum_{r=1}^k \hat{P}_{02[r]}(s, t),$$

$$\hat{P}_{02[r]}(s, t) = 1 - \hat{P}_{00[r]}(s, t) - \hat{P}_{01[r]}(s, t) \quad (r = 1, \dots, k),$$

$$\hat{P}_{12RSS}^{AJ}(s, t) = \frac{1}{k} \sum_{r=1}^k \hat{P}_{12[r]}(s, t),$$

$$\hat{P}_{12[r]}(s, t) = 1 - \hat{P}_{11[r]}(s, t) \quad (r = 1, \dots, k)$$

2.4 Simulation

In this section, a simulations study was performed to investigate properties of the proposed estimators. More specifically, the estimators for RSS modified transition probabilities (P_{00}, P_{01}, P_{02}) introduced in section 2.3 were considered.

First, a population of 20,000 subjects were simulated. All subjects in this population went through illness-death model process. Subjects were considered separately as two groups. The first group went through intermediate state 1 ($\sigma = I(T_{01} < T_{02}) = 1$), while the second group bypassed the intermediate state and went directly to the absorbing state 2 ($\sigma = I(T_{01} < T_{02}) = 0$).

For the first group of individuals ($\sigma = 1$), they have Z , the sojourn time in the initial state 0, and $T_{12} = T - Z$, the time in intermediate state 1. These two successive gap times were simulated according to the following bivariate distribution:⁶⁰

$$F_{1,2}(x, y) = F_1(x)F_2(y)[1 + \theta\{1 - F_1(x)\}\{1 - F_2(y)\}]$$

Where the marginal distribution functions F_1 and F_2 are exponential distribution with rate parameter 1. This is a multivariate distribution corresponding to Farlie-Gumbel-Morgenstern copula. The parameter θ controls the dependence level between successive gap times Z and $T - Z$. Since illness-death model is a Markov process, the two successive gap times Z and $T - Z$ are independent. The parameter θ is set to be 0 in the study.

For the second group of individuals ($\sigma = 0$), they only have Z , the sojourn time in the initial state 0. Z is simulated as an exponential distribution with rate parameter 1. σ , a parameter used to control the number of individuals going through intermediate state 1, is simulated independently according to a Bernoulli distribution with parameter $p = 0.4$ or 0.8 to adjust the censoring level.

The simulation procedure is as following:

- 1) Independently generate $V_1 \sim Uniform(0,1)$
- 2) Independently generate $V_2 \sim Uniform(0,1)$
- 3) $U_1 = V_1, A = -1, B = 1$

$$4) U_2 = V_2$$

$$5) Z = \log\left(\frac{1}{1-U_1}\right)$$

$$6) \sigma \sim \text{Bernoulli}(p), p = 0.4 \text{ or } 0.7$$

$$7) T = Z + \sigma * \log\left(\frac{1}{1-U_2}\right)$$

An independent censoring time variable C is generated according to uniform distribution to control the censoring levels in successive gap times Z and $T - Z$. C is simulated to be $Uniform(0, 2)$ and $Uniform(0, 1.5)$ to adjust the censoring level. When σ follows a Bernoulli distribution with parameter $p = 0.4$ and C is simulated as $Uniform(0, 2)$, the final transition probabilities of the simulated population are: $P_{00} = 0.1367, P_{01} = 0.1057, P_{02} = 0.7577$, which is censoring level 1. When σ follows a Bernoulli distribution with parameter $p = 0.8$ and C is simulated according to $Uniform(0, 1.5)$, the final transition probabilities of the simulated population are: $P_{00} = 0.1957, P_{01} = 0.2918, P_{02} = 0.5125$, which is censoring level 2. Therefore, from perspective of P_{02} , censoring level 1 is 76% and censoring level 2 is 51%.

The following sample sizes are considered under the above censoring levels: $n = 200, 400, 800, 1600$. For each of the above sample size, 500 samples are generated under both RSS and SRS sampling design. Under RSS design, for each of the above sample size various set numbers are considered and corresponding cycle numbers are generated to reach the desired sample size. When sample size is 200, set number k is 2 or 4 with cycle number m equal to 100 or 50. When sample size is 400, k is 2, 4 or 8 with cycle number m equal to 200, 100 or 50 respectively. When sample size is 800, k is 2, 4, 8 or 16 with m equal to 400, 200, 100 or 50. When sample size is 1600, k is 4, 8, 16 or 32 with m equal to 400, 200, 100 or 50 correspondingly. R package ‘‘RSSampling’’ is applied to select RSS samples.⁶¹ Estimators are computed under both sampling designs. R package ‘‘TP.idm’’ is utilized to calculate transition probabilities.⁴³ Their mean square errors (MSEs) and bias compared with population are

computed under 25, 50, 75 and 95 percentiles of the distribution function.³⁴ A relative MSE with SRS MSE over RSS MSE is computed to compare the performance of two estimators.

$$Relative\ MSE = \frac{SRS\ MSE}{RSS\ MSE}$$

When relative MSE is larger than 1, it means RSS estimators have superiority over their SRS counterparts. Otherwise, SRS estimators perform better.

2.5 Simulation Results

Tables 1-13 display the MSEs of Aalen–Johansen (AJ) estimators based on RSS or SRS with sample size from 200 to 1600 under two censoring levels. In Table 1, when sample size is 200 and set number is 2, for both censoring levels the superiority of the RSS estimators compared to the SRS estimators is more obvious for transition probabilities P_{02} . In Table 2, for the same sample size as set number increased to 4, the efficiency ascendancy of the RSS estimators over SRS estimators becomes more obvious for both transition probabilities P_{00} and P_{02} . Table 3 reports that as sample sized increases to 400, MSE decreases for all transition probabilities at all time points. In Table 4, for the same sample size 400, when set size increases to 4 the improved efficiency of RSS estimators over SRS estimators becomes more obvious for transition probabilities P_{00} and P_{02} . In Table 5, as set number increases to 8, the superiority of RSS estimators over its SRS counterparts turns out to be more obvious with a relative MSE equal to 2.29 for P_{02} at 50 percentiles (censoring level 1) and 2.03 for P_{02} at 75 percentiles (censoring level 2). In Table 6, as sample size increases to 800, MSEs become even smaller, since power is increased. In Table 7, as set number increases to 4, the highest efficiency of RSS estimator over SRS estimator is 1.7789 which happens for transition probability P_{02} at 75% percentiles when censoring level is 2. In Table 8, for both censoring levels, 11 out 12 RSS estimators perform better than their SRS competitors with exceptions occurring for transition probability P_{01} at 25 percentiles. Table 9 shows that when set number increases to 16, the efficiency of RSS estimator for P_{02} at 25 percentiles is 2.7398 times as good as the corresponding SRS estimator. In Table 10, as sample size increases to 1600 the MSEs of both RSS and SRS designs reduce with a minimum value of 0.0682×10^{-3} and a maximum value of 0.6958×10^{-3} . Table 11 indicates

that for the same sample size as set number increased to 8, the efficiency superiority of the RSS estimators over SRS estimators becomes significant for transition probabilities P_{00} and P_{02} at both censoring levels. In Table 12, as set number increases to 16, the superiority of RSS estimators over SRS estimators escalates to another level with a relative MSE equal to 3.03 for P_{02} at 25 percentiles when censoring level is 1. In Table 13, when set number rises to 16, RSS estimator for transition probability P_{02} at 25 percentiles (censoring level 1) is 4.09 times as efficient as the corresponding SRS estimator.

Tables 14 to 17 report the estimated MSEs of SRS and RSS distribution function estimators at 25%, 50%, 75% and 95% percentiles for various sample sizes when censoring level is 1. Under a fixed sample size, the set number k was changed to investigate the effect of set number on MSE in RSS design. When set number k is equal to 1, the columns indicate results of SRS sampling. When RSS estimator is better than its SRS counterpart, the number is set bold. From Table 14, it is shown that when sample size is 200, censoring level is 1 with a set number of 2, the MSEs of all RSS estimators are smaller than their corresponding SRS counterparts except for P_{01} at 50 percentiles. When set number is increased to 4, the MSEs from RSS design become even smaller except for P_{01} at 25 and 50 percentiles. In Table 15 when $n = 400$ with same censoring level, among 12 estimators of P_{00} based on RSS, 11 estimators are better than their SRS counterparts, only 1 estimator is worse (1.442×10^{-3} , for P_{00} at 95% when $k = 8$). In the same table, when the transition probability is P_{02} , all 12 RSS estimators perform better than SRS counterparts. However, RSS design does not always dominate over SRS. For the same sample size and censoring level, when transition probability is P_{01} , 5 out of 12 RSS estimators are better than their SRS counterparts. When sample size increases, MSEs decrease for both RSS and SRS designs with a minimum MSE of 0.2812×10^{-3} and a maximum MSE of 2.1278×10^{-3} . In Table 16, when sample size escalates to 800, MSEs continue reducing for both RSS and SRS designs with a minimum value of 0.984×10^{-3} and a maximum value of 0.9446×10^{-3} . As set number escalates from 2 to 16, RSS estimators keep improving at 7 out of 12 percentiles (P_{00} at 50, 75, 95 percentiles, P_{01} at 50, 95 percentiles, P_{02} at 25, 75 percentiles). Among the total of 48 RSS estimators, 39 perform better than their SRS counterparts. In Table 17, when sample size increases to 1600, the decreasing pattern of MSE with escalated set number continues for

RSS modified transition probabilities P_{00} and P_{02} except for P_{00} at 95 percentiles and P_{02} at 95 percentiles. From example, when set number increases from 4 to 32, the MSE reduced from 0.1304×10^{-3} to 0.0946×10^{-3} for P_{00} at 25 percentiles. Tables 18 to 21 show the estimated MSEs of SRS and RSS distribution function estimators at 25%, 50%, 75% and 95% percentiles for various sample sizes but with a censoring level of 2. When censoring level of P_{02} increases from 0.24 to 0.49 and sojourn time in state 0, P_{00} , rises from 0.14 to 0.2, the overall efficiency advantage of RSS over SRS design does not change. In Table 18, when sample size is 200, among all 24 RSS estimators 22 show improved efficiency over their SRS competitors. Table 19 reports that change of censoring level does not affect the overall decrease pattern of MSE as sample size increases. When sample size climbs from 200 to 400, the MSEs of both RSS and SRS designs reduce with a minimum value of 0.1586×10^{-3} and a maximum value of 3.2158×10^{-3} . In table 21, the decreasing trend of MSEs as set number escalates is similar for censoring level 2 as for censoring level 1. When set number increases from 4 to 32, 9 out of 12 percentiles have reduced MSE except for P_{00} at 95 percentiles, P_{01} at 25 percentiles and P_{01} at 95 percentiles.

Tables 22 to 29 are estimated bias of distribution function estimators based on SRS or RSS designs at some percentiles for different sample sizes and censoring levels. In Table 22, overall, when sample size is small, both RSS and SRS estimators have large bias with the minimum value of 0.0195 and a maximum value of 0.0522. 20 out 24 biases of RSS estimators are smaller than the corresponding SRS estimators. In table 23, as sample size jumps from 200 to 400, the overall bias is reduced for estimators from both RSS and SRS designs with the minimum value of 0.0127 to a maximum value of 0.037. In Table 24, the improvement of bias by RSS design is obvious for transition probabilities P_{00} and P_{02} . For transition probability P_{00} , except for the estimator when set number is 8 at 95 percentiles (0.0208) is slightly greater than the corresponding SRS estimator (0.0205), all other RSS estimators are more competitive than their relative counterparts. For transition probability P_{02} , 10 out of 12 RSS estimators have improved bias. In Table 25, when sample size becomes the largest, 1600, the improvement of bias for estimators from RSS design becomes even more apparent. In total, 43 out 48 RSS estimators have smaller bias than their corresponding estimators under SRS design. Tables 26 to 29 present estimated

bias of distribution function estimators based on SRS or RSS designs at some percentiles for various sample sizes, but the censoring level is 2. In Table 26, for the same sample size, when censoring level is 2 bias is slightly higher for estimators from both RSS and SRS designs than censoring level 1. The minimum value in this table is 0.0176, while the maximum value is 0.0668. Though the censoring level is changed, the overall improvement of bias for estimators from RSS design over SRS design is still clear for transition probabilities P_{00} and P_{02} . In Table 27, all RSS estimators have smaller bias than their relative SRS counterparts for transition probability P_{02} . In Table 28, as sample size increases to 800, 13 out of 16 RSS estimators have improved bias than corresponding estimators from SRS design for transition probability P_{00} and 15 out of 16 RSS estimators show decreased bias values compared to estimators based on SRS design for P_{02} . Table 29 indicates that when sample size climbs to 1600, the overall bias is the smallest for both RSS and SRS designs with a minimum value of 0.0038 and a maximum value of 0.0209. The general superiority of estimators from RSS design over SRS design becomes significant with 43 out of 48 RSS estimators having smaller bias than their SRS competitors. The exceptions are bias for P_{00} at 95 percentiles when k is equal to 4 (0.0162) and k is equal to 32 (0.0166) compared to their corresponding SRS estimator (0.0159), bias for P_{01} at 25 percentiles when $k = 4$ (0.0088) and $k = 16$ (0.0088) compared with the relative SRS estimator (0.0087).

Transition probabilities P_{11} and P_{12} could also be improved by RSS design by ranking T_{12} . Table 30 presents the estimated MSEs of transition probabilities P_{11} and P_{12} for sample size 200, censoring level 2. We could notice that P_{11} and P_{12} have equal MSEs, which is due to $P_{11} = 1 - P_{12}$. When set number is 2, 6 out of 8 RSS estimators have improved MSE compared to their SRS counterparts. For the same sample size, when set number increases to 4, all 8 RSS estimators perform better than SRS competitors. Table 31 shows the bias of transition probabilities P_{11} and P_{12} for sample size 200, censoring level 2. This result corresponds with the MSEs results. Table 32 displays the estimated MSEs of transition probabilities P_{11} and P_{12} when sample size is 400 and censoring level is 2. We can see that when sample size increases the overall MSEs decrease with a minimum value of 0.0014 and a maximum value of 0.0085. This is due to increased power. When set number is 2, only 4 out of 8 estimators having lower MSEs than their SRS

counterparts. When set number increases to 4, among 8 RSS estimators 6 perform better than corresponding SRS estimators. The dominant advantage of RSS estimator continues when set number is increased to 8. Table 33 shows the bias of transition probabilities P_{11} and P_{12} for sample size 400, censoring level 2. The results correspond well with Table 32.

Figure 3 to 6 display the efficiency of the Aalen–Johansen estimator based on RSS with respect to SRS counterpart at different percentiles when sample size and censoring level vary. As shown in figure 3, AJ transition probabilities P_{00} and P_{02} based on RSS design is more efficient than SRS design in most cases since the majority of relative MSEs are larger than 1 with the exception happen for P_{00} when $N = 400$, $k = 8$. For transition probability P_{01} , both sampling methods indicate comparable efficiency with most of their relative MSEs around 1. In Figure 4, as sample size increases to 1600, the efficiency advantage of RSS estimators over their corresponding SRS estimators becomes more obvious for P_{00} and P_{02} with many relative MSEs above 1.5. For transition probability P_{02} , when sample size is fixed, relative MSE escalates as set number increases. The relative MSE is above 4 when set size is 32 at 25 percentiles, which means the RSS estimator is 4 times as efficient as its SRS competitor in this scenario. When censoring level is increased, the overall superiority of RSS modified AJ estimator over SRS design is not changed. In Figure 5, for transition probabilities P_{00} and P_{02} , most of the relative MSEs are larger than one, which indicates estimators under RSS design provide closer values to the population. This trend becomes more obvious when sample size increases to 800 and 1600. In Figure 6, for transition probability P_{02} , almost all relative MSEs are above 1. When sample size is fixed, as set number escalates the advantages of RSS design over SRS design for estimator efficiency is also enlarged.

2.6 Discussion

In medical studies, when measurement of desired variables is expensive or the disease scientists are interested is rare, RSS could obtain representative samples by ranking a less costly concomitant variable of the sampling units. Aalen-Johansen estimator is a widely applied technique in medical research fields to estimate the transition probabilities between health, illness, and death states in progressive illness-death

Markov model. However, this estimator was not investigated under an RSS environment. In this study, we proposed a modified Aalen-Johansen estimator in association with a randomly right-censored RSS sampling design. The performance of this new estimator was compared with its counterpart in SRS design through simulation study to show the efficiency of the new estimator. It would be appealing to apply the novel estimator to illness death Markov model for calculating transition probabilities when the measurement of interested variable is expensive but there is an economic concomitant variable.

The preceding sections established that the efficiency of Aalen-Johansen estimator under RSS design is improved in majority simulation scenarios compared to its SRS counterpart with the same sample size. Our simulation results suggest that for transition probabilities P_{00} and P_{02} , when sample size is fixed the superiority of the new estimator over conventional Aalen-Johansen estimator under SRS design escalates as set number increases. The same phenomenon was observed in RSS and Partially Rank-Ordered Set (PROS) modified Kaplan-Meier estimators previously.^{33,34} This is mainly due to when set number increases, the effect of ranking takes place which contributes to the selection of more representative samples. As we discussed before, for RSS, when sample size is $N = mk$, the number of subjects who truly participate in the process is mk^2 . As set number k increases, there are more subjects involved in the underlying ranking process. Though sample size stays same as SRS, more information is provided by RSS process. The results of simulation study also indicate that as sample size increases, the MSE of Aalen-Johansen estimator under both RSS and SRS designs decreases, since the power of study increases. This corresponds well with previous studies regarding RSS and PROS modified Kaplan-Meier estimators.^{33,34}

The improvement in the efficiency of the estimator based on the RSS design was only observed for transition probabilities P_{00}, P_{02} , but not for transition probability P_{01} , which corresponds well with the Markov property of illness death model. In Markov process any future evolution of the Markov process depends only on its current state and is independent of the previously visited states. Since progressive illness death model is a stochastic process, random variables T_{01} (time from state 0 to state 1), T_{12} (time from state 1 to state 2) and T_{02} (time from state 0 to state 2) are independent.⁶⁰ When ranking variable is

T_{02} , the transition probabilities that are influenced are P_{00} and P_{02} , which justifies the Markov property of illness death model simultaneously. However, though the ranking process put no effect on transition probability P_{01} , the efficiency of the P_{01} estimator based on the RSS design did not perform worse than the efficiency of the corresponding estimator based on the SRS method with the same sample size, which corresponds with previous studies suggesting that there is no harm if we do RSS design.^{9,10}

Our simulation results show that change in the censoring level does not affect the superiority of RSS modified estimator over its SRS competitor. Under both censoring levels, the efficiency of the estimator based on the RSS sampling design is better than the efficiency of the estimator based on the SRS method with the same sample sizes. For both censoring levels, by increasing the set size in RSS while the sample size stays the same, the RSS Aalen-Johansen estimator has improved MSE and bias than its SRS counterpart. The similar trend of censoring level on estimator efficiency was observed in previous studies involving RSS design in survival analysis and censored data.^{32-34,37} We believe it would be appealing to apply the introduced methodology to medical studies when measurement of interested variable is costly or the disease is rare.

Table 1. MSEs of transition probabilities when sample size is 200 (set size: 2, cycle number: 100).

Censoring level	Time point	P ₀₀				P ₀₁				P ₀₂			
		0.25 (494)	0.50 (987)	0.75 (1481)	0.95 (1875)	0.25 (494)	0.50 (987)	0.75 (1481)	0.95 (1875)	0.25 (494)	0.50 (987)	0.75 (1481)	0.95 (1875)
1	RSS MSE	1.203	1.416	1.8094	2.4926	0.5734	0.9442	1.3414	2.3626	0.9606	1.4108	2.169	3.6052
	SRS MSE	1.3656	1.7638	2.0354	3.1796	0.625	0.909	1.4436	2.7198	1.1448	1.6222	2.4028	4.3172
	$\frac{SRS\ MSE}{RSS\ MSE}$	1.1352	1.2456	1.1249	1.2756	1.0923	0.9627	1.0762	1.1512	1.1918	1.1498	1.1078	1.1975
2	Time point	0.25 (374)	0.50 (748)	0.75 (1122)	0.95 (1421)	0.25 (374)	0.50 (748)	0.75 (1122)	0.95 (1421)	0.25 (374)	0.50 (748)	0.75 (1122)	0.95 (1421)
	RSS MSE	1.1386	1.64	2.1342	3.5782	0.989	1.5298	2.2954	5.7448	0.4936	1.0338	1.7634	4.716
	SRS MSE	1.2288	1.7086	2.155	4.2988	1.0092	1.6052	2.609	7.0796	0.5586	1.4554	2.4034	6.1004
	$\frac{SRS\ MSE}{RSS\ MSE}$	1.0792	1.0418	1.0097	1.2014	1.0204	1.0493	1.1366	1.2323	1.1317	1.4078	1.3629	1.2936

*MSEs are values listed in the table time 10⁻³.

Table 2. MSEs of transition probabilities when sample size is 200 (set size: 4, cycle number: 50).

Censoring level	Time point	P ₀₀				P ₀₁				P ₀₂			
		0.25 (494)	0.50 (987)	0.75 (1481)	0.95 (1875)	0.25 (494)	0.50 (987)	0.75 (1481)	0.95 (1875)	0.25 (494)	0.50 (987)	0.75 (1481)	0.95 (1875)
1	RSS MSE	1.1194	1.2624	1.6582	2.46	0.6688	1.0306	1.2398	1.9136	0.713	1.1904	1.6038	3.2912
	SRS MSE	1.3656	1.7638	2.0354	3.1796	0.625	0.909	1.4436	2.7198	1.1448	1.6222	2.4028	4.3172
	$\frac{SRS\ MSE}{RSS\ MSE}$	1.2199	1.3972	1.2275	1.2925	0.9345	0.8820	1.1644	1.4213	1.6056	1.3627	1.4982	1.3117
2	Time point	0.25 (374)	0.50 (748)	0.75 (1122)	0.95 (1421)	0.25 (374)	0.50 (748)	0.75 (1122)	0.95 (1421)	0.25 (374)	0.50 (748)	0.75 (1122)	0.95 (1421)
	RSS MSE	1.0694	1.4876	1.8656	3.051	1.0136	1.7982	2.5184	4.8754	0.4792	0.903	1.763	3.5986
	SRS MSE	1.2288	1.7086	2.155	4.2988	1.0092	1.6052	2.609	7.0796	0.5586	1.4554	2.4034	6.1004
	$\frac{SRS\ MSE}{RSS\ MSE}$	1.1491	1.1486	1.1551	1.409	0.9957	0.8927	1.036	1.4521	1.1657	1.6117	1.3632	1.6952

*MSEs are values listed in the table time 10⁻³.

Table 3. MSEs of transition probabilities when sample size is 400 (set size: 2, cycle number: 200)

Censoring level	Time point	P ₀₀				P ₀₁				P ₀₂			
		0.25 (494)	0.50 (987)	0.75 (1481)	0.95 (1875)	0.25 (494)	0.50 (987)	0.75 (1481)	0.95 (1875)	0.25 (494)	0.50 (987)	0.75 (1481)	0.95 (1875)
1	RSS MSE	0.6074	0.703	0.8576	1.2196	0.2986	0.4882	0.6682	1.202	0.4808	0.7678	1.0432	1.739
	SRS MSE	0.658	0.9102	0.9584	1.4228	0.285	0.454	0.6364	1.4032	0.5572	0.8574	1.1706	2.1278
	$\frac{SRS\ MSE}{RSS\ MSE}$	1.0833	1.2950	1.1175	1.1666	0.9545	0.9299	0.9524	1.1674	1.1589	1.1170	1.1221	1.2236
2	Time point	0.25 (374)	0.50 (748)	0.75 (1122)	0.95 (1421)	0.25 (374)	0.50 (748)	0.75 (1122)	0.95 (1421)	0.25 (374)	0.50 (748)	0.75 (1122)	0.95 (1421)
	RSS MSE	0.5844	0.8768	1.0922	2.2054	0.4774	0.822	1.354	3.0468	0.264	0.5378	0.8876	2.394
	SRS MSE	0.6214	0.8646	1.062	1.8	0.4518	0.769	1.293	3.2158	0.2676	0.6798	1.2176	2.5728
	$\frac{SRS\ MSE}{RSS\ MSE}$	1.0633	0.9861	0.9723	0.8162	0.9464	0.9355	0.9549	1.0555	1.0136	1.2644	1.3718	1.0747

*MSEs are values listed in the table time 10⁻³.

Table 4. MSEs of transition probabilities when sample size is 400 (set size: 4, cycle number: 100).

Censoring level	Time point	P ₀₀				P ₀₁				P ₀₂			
		0.25 (494)	0.50 (987)	0.75 (1481)	0.95 (1875)	0.25 (494)	0.50 (987)	0.75 (1481)	0.95 (1875)	0.25 (494)	0.50 (987)	0.75 (1481)	0.95 (1875)
1	RSS MSE	0.5676	0.6406	0.831	1.1972	0.3238	0.4808	0.6346	1.118	0.3538	0.604	0.8662	1.5456
	SRS MSE	0.658	0.9104	0.9584	1.4228	0.285	0.454	0.6364	1.4032	0.5572	0.8574	1.1706	2.1278
	$\frac{SRS\ MSE}{RSS\ MSE}$	1.1593	1.4212	1.1533	1.1884	0.8802	0.9443	1.0028	1.2551	1.5749	1.4195	1.3514	1.3767
2	Time point	0.25 (374)	0.50 (748)	0.75 (1122)	0.95 (1421)	0.25 (374)	0.50 (748)	0.75 (1122)	0.95 (1421)	0.25 (374)	0.50 (748)	0.75 (1122)	0.95 (1421)
	RSS MSE	0.5568	0.7088	0.9802	1.5862	0.4982	0.787	1.2174	2.5236	0.2176	0.4766	0.8024	2.0142
	SRS MSE	0.6214	0.8646	1.062	1.8	0.4518	0.769	1.293	3.2158	0.2676	0.6798	1.2176	2.5728
	$\frac{SRS\ MSE}{RSS\ MSE}$	1.1160	1.2198	1.0835	1.1348	0.9069	0.9771	1.0621	1.2743	1.2298	1.4264	1.5174	1.2773

*MSEs are values listed in the table time 10⁻³.

Table 5. MSEs of transition probabilities when sample size is 400 (set size: 8, cycle number: 50).

Censoring level	Time point	P ₀₀				P ₀₁				P ₀₂			
		0.25 (494)	0.50 (987)	0.75 (1481)	0.95 (1875)	0.25 (494)	0.50 (987)	0.75 (1481)	0.95 (1875)	0.25 (494)	0.50 (987)	0.75 (1481)	0.95 (1875)
1	RSS MSE	0.4714	0.5148	0.7994	1.442	0.323	0.453	0.6784	1.0656	0.2812	0.3752	0.7912	1.8046
	SRS MSE	0.658	0.9104	0.9584	1.4228	0.285	0.454	0.6364	1.4032	0.5572	0.8574	1.1706	2.1278
	$\frac{SRS\ MSE}{RSS\ MSE}$	1.3958	1.7685	1.1989	0.9867	0.8824	1.0022	0.9381	1.3168	1.9815	2.2852	1.4795	1.1791
2	Time point	0.25 (374)	0.50 (748)	0.75 (1122)	0.95 (1421)	0.25 (374)	0.50 (748)	0.75 (1122)	0.95 (1421)	0.25 (374)	0.50 (748)	0.75 (1122)	0.95 (1421)
	RSS MSE	0.4834	0.6744	0.8824	1.6454	0.473	0.7204	1.2074	2.4508	0.1586	0.3866	0.5998	1.6136
	SRS MSE	0.6214	0.8646	1.062	1.8	0.4518	0.769	1.293	3.2158	0.2676	0.6798	1.2176	2.5728
	$\frac{SRS\ MSE}{RSS\ MSE}$	1.2855	1.2820	1.2035	1.0940	0.9552	1.0675	1.0709	1.3121	1.6873	1.7584	2.03	1.5944

*MSEs are values listed in the table time 10⁻³.

Table 6. MSEs of transition probabilities when sample size is 800 (set size: 2, cycle number: 400).

Censoring level	Time point	P ₀₀				P ₀₁				P ₀₂			
		0.25 (494)	0.50 (987)	0.75 (1481)	0.95 (1875)	0.25 (494)	0.50 (987)	0.75 (1481)	0.95 (1875)	0.25 (494)	0.50 (987)	0.75 (1481)	0.95 (1875)
1	RSS MSE	0.3322	0.3732	0.4326	0.6322	0.1514	0.236	0.3486	0.5518	0.2414	0.3944	0.4814	0.9222
	SRS MSE	0.3256	0.409	0.4772	0.706	0.1494	0.238	0.3308	0.6372	0.2696	0.3848	0.5956	0.9118
	$\frac{SRS\ MSE}{RSS\ MSE}$	0.9801	1.0959	1.1031	1.1167	0.9868	1.0085	0.9489	1.1537	1.1168	0.9757	1.2372	0.9887
2	Time point	0.25 (374)	0.50 (748)	0.75 (1122)	0.95 (1421)	0.25 (374)	0.50 (748)	0.75 (1122)	0.95 (1421)	0.25 (374)	0.50 (748)	0.75 (1122)	0.95 (1421)
	RSS MSE	0.2866	0.4286	0.614	1.0384	0.2412	0.4454	0.6816	1.6006	0.124	0.2904	0.4634	1.2606
	SRS MSE	0.31	0.4346	0.5014	0.8528	0.2286	0.3896	0.7086	1.5468	0.14	0.3308	0.6468	1.2396
	$\frac{SRS\ MSE}{RSS\ MSE}$	1.0816	1.014	0.8166	0.8213	0.9478	0.8747	1.0396	0.9664	1.129	1.1391	1.3958	0.9833

*MSEs are values listed in the table time 10⁻³.

Table 7. MSEs of transition probabilities when sample size is 800 (set size: 4, cycle number: 200).

Censoring level	Time point	P ₀₀				P ₀₁				P ₀₂			
		0.25 (494)	0.50 (987)	0.75 (1481)	0.95 (1875)	0.25 (494)	0.50 (987)	0.75 (1481)	0.95 (1875)	0.25 (494)	0.50 (987)	0.75 (1481)	0.95 (1875)
1	RSS MSE	0.2734	0.3416	0.3698	0.6246	0.132	0.2354	0.3404	0.6134	0.1748	0.3046	0.4262	0.797
	SRS MSE	0.3256	0.409	0.4772	0.706	0.1494	0.238	0.3308	0.6366	0.2696	0.3848	0.5956	0.9118
	$\frac{SRS\ MSE}{RSS\ MSE}$	1.1909	1.1973	1.2904	1.1303	1.1318	1.0110	0.9718	1.0378	1.5423	1.2633	1.3975	1.1440
2	Time point	0.25 (374)	0.50 (748)	0.75 (1122)	0.95 (1421)	0.25 (374)	0.50 (748)	0.75 (1122)	0.95 (1421)	0.25 (374)	0.50 (748)	0.75 (1122)	0.95 (1421)
	RSS MSE	0.2486	0.3544	0.4704	0.8282	0.2388	0.3582	0.5548	1.4196	0.1138	0.2142	0.3636	1.0166
	SRS MSE	0.31	0.4346	0.5014	0.8528	0.2286	0.3896	0.7086	1.5468	0.14	0.3308	0.6468	1.2396
	$\frac{SRS\ MSE}{RSS\ MSE}$	1.247	1.2263	1.0659	1.0297	0.9573	1.0877	1.2772	1.0896	1.2302	1.5444	1.7789	1.2194

*MSEs are values listed in the table time 10⁻³.

Table 8. MSEs of transition probabilities when sample size is 800 (set size: 8, cycle number: 100).

Censoring level	Time point	P ₀₀				P ₀₁				P ₀₂			
		0.25 (494)	0.50 (987)	0.75 (1481)	0.95 (1875)	0.25 (494)	0.50 (987)	0.75 (1481)	0.95 (1875)	0.25 (494)	0.50 (987)	0.75 (1481)	0.95 (1875)
1	RSS MSE	0.2456	0.2844	0.3698	0.69	0.1654	0.2094	0.3074	0.562	0.1272	0.213	0.3738	0.8034
	SRS MSE	0.3256	0.409	0.4772	0.706	0.1494	0.238	0.3308	0.6366	0.2696	0.3848	0.5956	0.9118
	$\frac{SRS\ MSE}{RSS\ MSE}$	1.3257	1.4381	1.2904	1.0232	0.9033	1.1366	1.0761	1.1327	2.1195	1.8066	1.5934	1.1349
2	Time point	0.25 (374)	0.50 (748)	0.75 (1122)	0.95 (1421)	0.25 (374)	0.50 (748)	0.75 (1122)	0.95 (1421)	0.25 (374)	0.50 (748)	0.75 (1122)	0.95 (1421)
	RSS MSE	0.2508	0.3534	0.48	0.8132	0.2416	0.379	0.6014	1.321	0.082	0.1792	0.3176	0.8194
	SRS MSE	0.31	0.4346	0.5014	0.8528	0.2286	0.3896	0.7086	1.5468	0.14	0.3308	0.6468	1.2396
	$\frac{SRS\ MSE}{RSS\ MSE}$	1.236	1.2298	1.0446	1.0487	0.9462	1.028	1.1783	1.1709	1.7073	1.846	2.0365	1.5128

*MSEs are values listed in the table time 10⁻³.

Table 9. MSEs of transition probabilities when sample size is 800 (set size: 16, cycle number: 50).

Censoring level	Time point	P ₀₀				P ₀₁				P ₀₂			
		0.25 (494)	0.50 (987)	0.75 (1481)	0.95 (1875)	0.25 (494)	0.50 (987)	0.75 (1481)	0.95 (1875)	0.25 (494)	0.50 (987)	0.75 (1481)	0.95 (1875)
1	RSS MSE	0.2324	0.2782	0.362	0.7098	0.1602	0.2282	0.2832	0.558	0.0984	0.1496	0.3526	0.9446
	SRS MSE	0.3256	0.409	0.4772	0.706	0.1494	0.238	0.3308	0.6366	0.2696	0.3848	0.5956	0.9118
	$\frac{SRS\ MSE}{RSS\ MSE}$	1.4010	1.4702	1.3182	0.9946	0.9326	1.0429	1.1681	1.1409	2.7398	2.5722	1.6892	0.9653
2	Time point	0.25 (374)	0.50 (748)	0.75 (1122)	0.95 (1421)	0.25 (374)	0.50 (748)	0.75 (1122)	0.95 (1421)	0.25 (374)	0.50 (748)	0.75 (1122)	0.95 (1421)
	RSS MSE	0.2504	0.3248	0.4156	0.8074	0.232	0.377	0.5526	1.193	0.0696	0.1248	0.255	1.009
	SRS MSE	0.31	0.4346	0.5014	0.8528	0.2286	0.3896	0.7086	1.5468	0.14	0.3308	0.6468	1.2396
	$\frac{SRS\ MSE}{RSS\ MSE}$	1.238	1.3381	1.2064	1.0562	0.9853	1.0334	1.2823	1.2966	2.0115	2.6506	2.5365	1.2285

*MSEs are values listed in the table time 10⁻³.

Table 10. MSEs of transition probabilities when sample size is 1600 (set size: 4, cycle number: 400).

Censoring level	Time point	P ₀₀				P ₀₁				P ₀₂			
		0.25 (494)	0.50 (987)	0.75 (1481)	0.95 (1875)	0.25 (494)	0.50 (987)	0.75 (1481)	0.95 (1875)	0.25 (494)	0.50 (987)	0.75 (1481)	0.95 (1875)
1	RSS MSE	0.1304	0.1608	0.2006	0.3228	0.0682	0.1016	0.1486	0.2652	0.089	0.133	0.2068	0.386
	SRS MSE	0.1776	0.2182	0.2228	0.3184	0.0708	0.116	0.1638	0.2774	0.1424	0.2164	0.2562	0.4034
	$\frac{SRS\ MSE}{RSS\ MSE}$	1.3620	1.3570	1.1107	0.9864	1.0381	1.1417	1.1023	1.0460	1.6	1.6271	1.2389	1.0451
2	Time point	0.25 (374)	0.50 (748)	0.75 (1122)	0.95 (1421)	0.25 (374)	0.50 (748)	0.75 (1122)	0.95 (1421)	0.25 (374)	0.50 (748)	0.75 (1122)	0.95 (1421)
	RSS MSE	0.1256	0.1892	0.2448	0.416	0.1204	0.1896	0.301	0.6958	0.0546	0.1	0.1926	0.4876
	SRS MSE	0.1706	0.2316	0.2666	0.382	0.1172	0.1912	0.3212	0.6808	0.0682	0.1692	0.2864	0.6042
	$\frac{SRS\ MSE}{RSS\ MSE}$	1.3583	1.2241	1.0891	0.9183	0.9734	1.0084	1.0671	0.9784	1.2491	1.692	1.487	1.2391

*MSEs are values listed in the table time 10⁻³.

Table 11. MSEs of transition probabilities when sample size is 1600 (set size: 8, cycle number: 200).

Censoring level	Time point	P ₀₀				P ₀₁				P ₀₂			
		0.25 (494)	0.50 (987)	0.75 (1481)	0.95 (1875)	0.25 (494)	0.50 (987)	0.75 (1481)	0.95 (1875)	0.25 (494)	0.50 (987)	0.75 (1481)	0.95 (1875)
1	RSS MSE	0.1216	0.1376	0.1806	0.3102	0.0842	0.1072	0.1436	0.2852	0.064	0.0968	0.1864	0.3792
	SRS MSE	0.1776	0.2182	0.2228	0.3184	0.0708	0.116	0.1638	0.2774	0.1424	0.2164	0.2562	0.4034
	$\frac{SRS\ MSE}{RSS\ MSE}$	1.4605	1.5858	1.2337	1.0264	0.8409	1.0821	1.1407	0.9727	2.225	2.2355	1.3745	1.0638
2	Time point	0.25 (374)	0.50 (748)	0.75 (1122)	0.95 (1421)	0.25 (374)	0.50 (748)	0.75 (1122)	0.95 (1421)	0.25 (374)	0.50 (748)	0.75 (1122)	0.95 (1421)
	RSS MSE	0.1242	0.1626	0.234	0.3918	0.1178	0.1802	0.3154	0.6758	0.0432	0.0862	0.1506	0.406
	SRS MSE	0.1706	0.2316	0.2666	0.382	0.1172	0.1912	0.3212	0.6808	0.0682	0.1692	0.2864	0.6042
	$\frac{SRS\ MSE}{RSS\ MSE}$	1.3736	1.4244	1.1393	0.975	0.9949	1.061	1.0184	1.0074	1.5787	1.9629	1.9017	1.4882

*MSEs are values listed in the table time 10⁻³.

Table 12. MSEs of transition probabilities when sample size is 1600 (set size: 16, cycle number: 100).

Censoring level	Time point	P ₀₀				P ₀₁				P ₀₂			
		0.25 (494)	0.50 (987)	0.75 (1481)	0.95 (1875)	0.25 (494)	0.50 (987)	0.75 (1481)	0.95 (1875)	0.25 (494)	0.50 (987)	0.75 (1481)	0.95 (1875)
1	RSS MSE	0.1174	0.1348	0.1898	0.314	0.079	0.1076	0.1388	0.285	0.047	0.0726	0.1568	0.4178
	SRS MSE	0.1776	0.2182	0.2228	0.3184	0.0708	0.116	0.1638	0.2774	0.1424	0.2164	0.2562	0.4034
	$\frac{SRS\ MSE}{RSS\ MSE}$	1.5128	1.6187	1.1739	1.014	0.8962	1.0781	1.1801	0.9733	3.0298	2.9807	1.6339	0.9655
2	Time point	0.25 (374)	0.50 (748)	0.75 (1122)	0.95 (1421)	0.25 (374)	0.50 (748)	0.75 (1122)	0.95 (1421)	0.25 (374)	0.50 (748)	0.75 (1122)	0.95 (1421)
	RSS MSE	0.1174	0.144	0.1824	0.3708	0.1268	0.1756	0.2388	0.5998	0.0314	0.0698	0.1224	0.4386
	SRS MSE	0.1706	0.2316	0.2666	0.382	0.1172	0.1912	0.3212	0.6808	0.0682	0.1692	0.2864	0.6042
	$\frac{SRS\ MSE}{RSS\ MSE}$	1.4532	1.4942	1.4616	1.0302	0.9242	1.0888	1.3451	1.135	2.172	2.4241	2.3399	1.3776

*MSEs are values listed in the table time 10⁻³.

Table 13. MSEs of transition probabilities when sample size is 1600 (set size: 32, cycle number: 50).

Censoring level	Time point	P ₀₀				P ₀₁				P ₀₂			
		0.25 (494)	0.50 (987)	0.75 (1481)	0.95 (1875)	0.25 (494)	0.50 (987)	0.75 (1481)	0.95 (1875)	0.25 (494)	0.50 (987)	0.75 (1481)	0.95 (1875)
1	RSS MSE	0.0946	0.1256	0.1838	0.3958	0.0768	0.1204	0.1384	0.2762	0.0348	0.067	0.1536	0.5648
	SRS MSE	0.1776	0.2182	0.2228	0.3184	0.0708	0.116	0.1628	0.2774	0.1424	0.2164	0.2562	0.4034
	$\frac{SRS\ MSE}{RSS\ MSE}$	1.8774	1.7373	1.2122	0.8044	0.9219	0.9635	1.1835	1.0043	4.0920	3.2299	1.6680	0.7142
2	Time point	0.25 (374)	0.50 (748)	0.75 (1122)	0.95 (1421)	0.25 (374)	0.50 (748)	0.75 (1122)	0.95 (1421)	0.25 (374)	0.50 (748)	0.75 (1122)	0.95 (1421)
	RSS MSE	0.106	0.153	0.2008	0.4198	0.1096	0.16	0.2408	0.5412	0.0234	0.0512	0.1012	0.4432
	SRS MSE	0.1706	0.2316	0.2666	0.382	0.1172	0.1912	0.3212	0.6808	0.0682	0.1692	0.2864	0.6042
	$\frac{SRS\ MSE}{RSS\ MSE}$	1.6094	1.5137	1.3277	0.9104	1.0693	1.195	1.3339	1.2579	2.9145	3.3047	2.83	1.3633

*MSEs are values listed in the table time 10⁻³.

Table 14. Estimated MSEs of distribution function estimators at some percentiles for $n = 200$, the first censoring level.

Transition Probabilities	Percentiles	N = 200		
		k		
		1	2	4
P ₀₀	25	1.3656	1.203	1.1194
	50	1.7638	1.416	1.2624
	75	2.0354	1.8094	1.6582
	95	3.1796	2.4926	2.46
P ₀₁	25	0.625	0.5722	0.6688
	50	0.909	0.9442	1.0306
	75	1.4436	1.3414	1.2398
	95	2.7198	2.3626	1.9136
P ₀₂	25	1.1448	0.9606	0.713
	50	1.6222	1.4108	1.1904
	75	2.4028	2.169	1.6038
	95	4.3172	3.6052	3.2912

*MSEs are values listed in the table time 10^{-3} .

Table 15. Estimated MSEs of distribution function estimators at some percentiles for $n = 400$, the first censoring level.

Transition Probabilities	Percentiles	N = 400			
		k			
		1	2	4	8
P ₀₀	25	0.658	0.6074	0.5676	0.4714
	50	0.9104	0.703	0.6406	0.5148
	75	0.9584	0.8576	0.831	0.7994
	95	1.4228	1.2196	1.1972	1.442
P ₀₁	25	0.285	0.2986	0.3238	0.323
	50	0.454	0.4882	0.4808	0.453
	75	0.6364	0.6682	0.6346	0.6784
	95	1.4032	1.202	1.118	1.0656
P ₀₂	25	0.5572	0.4808	0.3538	0.2812
	50	0.8574	0.7678	0.604	0.3752
	75	1.1706	1.0432	0.8662	0.7912
	95	2.1278	1.739	1.5456	1.8046

*MSEs are values listed in the table time 10^{-3} .

Table 16. Estimated MSEs of distribution function estimators at some percentiles for $n = 800$, the first censoring level.

Transition Probabilities	Percentile	N = 800				
		k				
		1	2	4	8	16
P ₀₀	25	0.3256	0.3322	0.2734	0.2456	0.2324
	50	0.409	0.3732	0.3416	0.2844	0.2782
	75	0.4772	0.4326	0.3698	0.3698	0.362
	95	0.706	0.6322	0.6246	0.69	0.7098
P ₀₁	25	0.1494	0.1514	0.132	0.1654	0.1602
	50	0.238	0.236	0.2354	0.2094	0.2282
	75	0.1654	0.3486	0.3404	0.3074	0.2832
	95	0.6366	0.5518	0.6134	0.562	0.558
P ₀₂	25	0.2696	0.3414	0.1748	0.1272	0.0984
	50	0.3848	0.3944	0.3046	0.213	0.1496
	75	0.5956	0.4814	0.4262	0.3738	0.3526
	95	0.9118	0.9222	0.797	0.8034	0.9446

*MSEs are values listed in the table time 10^{-3} .

Table 17. Estimated MSEs of distribution function estimators at some percentiles for $n = 1600$, the first censoring level.

Transition Probabilities	Percentile	N = 1600				
		k				
		1	4	8	16	32
P ₀₀	25	0.1776	0.1304	0.1216	0.1174	0.0946
	50	0.2182	0.1608	0.1376	0.1348	0.1256
	75	0.2228	0.2006	0.1806	0.1898	0.1838
	95	0.3184	0.3228	0.3102	0.314	0.3958
P ₀₁	25	0.0708	0.0682	0.0842	0.079	0.0768
	50	0.116	0.1016	0.1072	0.1076	0.1204
	75	0.1638	0.1486	0.1436	0.1388	0.1384
	95	0.2774	0.2652	0.2852	0.285	0.2762
P ₀₂	25	0.1424	0.089	0.064	0.047	0.0348
	50	0.2164	0.133	0.0968	0.0726	0.067
	75	0.2562	0.2068	0.1862	0.1568	0.1536
	95	0.4034	0.386	0.3972	0.4178	0.5648

*MSEs are values listed in the table time 10^{-3} .

Table 18. Estimated MSEs of distribution function estimators at some percentiles for $n = 200$, the second censoring level.

Transition Probabilities	Percentiles	N = 200		
		k		
		1	2	4
P ₀₀	25	1.2288	1.1386	1.0694
	50	1.7086	1.64	1.4876
	75	2.155	2.1342	1.8656
	95	4.2988	3.5782	3.051
P ₀₁	25	1.0092	0.989	1.0136
	50	1.6052	1.5298	1.7982
	75	2.609	2.2954	2.5184
	95	7.0796	5.7448	4.8754
P ₀₂	25	0.5586	0.4936	0.4792
	50	1.4554	1.0338	0.903
	75	2.4034	1.7634	1.763
	95	6.1004	4.716	3.5986

*MSEs are values listed in the table time 10^{-3} .

Table 19. Estimated MSEs of distribution function estimators at some percentiles for $n = 400$, the second censoring level.

Transition Probabilities	Percentiles	N = 400			
		k			
		1	2	4	8
P ₀₀	25	0.6214	0.5844	0.5568	0.4834
	50	0.8646	0.8768	0.7088	0.6744
	75	1.062	1.0922	0.9802	0.8824
	95	1.8	2.2054	1.5862	1.6454
P ₀₁	25	0.4518	0.4774	0.4982	0.473
	50	0.769	0.822	0.787	0.7204
	75	1.293	1.354	1.2174	1.2074
	95	3.2158	3.0468	2.5236	2.4508
P ₀₂	25	0.2676	0.264	0.2176	0.1586
	50	0.6798	0.5378	0.4766	0.3866
	75	1.2176	0.8876	0.8024	0.5998
	95	2.5728	2.394	2.0142	1.6136

*MSEs are values listed in the table time 10^{-3} .

Table 20. Estimated MSEs of distribution function estimators at some percentiles for $n = 800$, the second censoring level.

Transition Probabilities	Percentile	N = 800				
		k				
		1	2	4	8	16
P ₀₀	25	0.31	0.2866	0.2486	0.2508	0.2504
	50	0.4346	0.4286	0.3544	0.3534	0.3248
	75	0.5014	0.614	0.4704	0.48	0.4156
	95	0.8528	1.0384	0.8282	0.8132	0.8074
P ₀₁	25	0.2286	0.2412	0.2388	0.2416	0.232
	50	0.3896	0.4454	0.3582	0.379	0.377
	75	0.7086	0.6816	0.5548	0.6014	0.5526
	95	1.5468	1.6006	1.4196	1.321	1.193
P ₀₂	25	0.14	0.124	0.1138	0.082	0.0696
	50	0.3308	0.2904	0.2142	0.1792	0.1248
	75	0.6468	0.4634	0.3636	0.3176	0.255
	95	1.2396	1.2606	1.0166	1.8194	1.009

*MSEs are values listed in the table time 10^{-3} .

Table 21. Estimated MSEs of distribution function estimators at some percentiles for $n = 1600$, the second censoring level.

Transition Probabilities	Percentile	N = 1600				
		k				
		1	4	8	16	32
P ₀₀	25	0.1706	0.1256	0.1242	0.1174	0.106
	50	0.2316	0.1892	0.1626	0.155	0.153
	75	0.2666	0.2448	0.234	0.1824	0.2008
	95	0.382	0.416	0.3918	0.3708	0.4196
P ₀₁	25	0.1172	0.1204	0.1178	0.1268	0.1096
	50	0.1912	0.1896	0.1802	0.1756	0.16
	75	0.3212	0.301	0.3154	0.2388	0.2408
	95	0.6808	0.6958	0.6758	0.5998	0.5412
P ₀₂	25	0.0682	0.0546	0.0432	0.0314	0.0234
	50	0.1692	0.1	0.0862	0.0698	0.0512
	75	0.2864	0.1926	0.1506	0.1224	0.1012
	95	0.6042	0.4876	0.406	0.4386	0.4432

*MSEs are values listed in the table time 10^{-3} .

Table 22. Estimated bias of distribution function estimators at some percentiles for $n = 200$, the first censoring level.

Transition Probabilities	Percentiles	N = 200		
		k		
		1	2	4
P ₀₀	25	0.029	0.0271	0.0256
	50	0.033	0.0306	0.0303
	75	0.0343	0.0317	0.0322
	95	0.0425	0.0398	0.039
P ₀₁	25	0.0204	0.0195	0.0201
	50	0.0242	0.0243	0.0245
	75	0.027	0.0274	0.0271
	95	0.0384	0.0381	0.0366
P ₀₂	25	0.0256	0.0244	0.0216
	50	0.0351	0.0304	0.0262
	75	0.0395	0.0352	0.032
	95	0.0522	0.0503	0.0447

Table 23. Estimated bias of distribution function estimators at some percentiles for $n = 400$, the first censoring level.

Transition Probabilities	Percentiles	N = 400			
		k			
		1	2	4	8
P ₀₀	25	0.0207	0.0198	0.0177	0.017
	50	0.0231	0.0225	0.0208	0.0193
	75	0.0241	0.023	0.0228	0.0215
	95	0.0292	0.028	0.0295	0.0288
P ₀₁	25	0.0137	0.0146	0.0139	0.014
	50	0.0168	0.0167	0.0168	0.0171
	75	0.019	0.0203	0.0199	0.0201
	95	0.0277	0.026	0.0282	0.0265
P ₀₂	25	0.0187	0.0173	0.015	0.0127
	50	0.0237	0.0211	0.0186	0.0176
	75	0.0282	0.0246	0.0233	0.0217
	95	0.037	0.0319	0.0342	0.0333

Table 24. Estimated bias of distribution function estimators at some percentiles for $n = 800$, the first censoring level.

Transition Probabilities	Percentile	N = 800				
		k				
		1	2	4	8	16
P_{00}	25	0.0146	0.0144	0.0132	0.0117	0.0107
	50	0.0164	0.016	0.0147	0.0126	0.0127
	75	0.0167	0.0166	0.016	0.0159	0.0149
	95	0.0205	0.0204	0.0193	0.0208	0.0197
P_{01}	25	0.01	0.0103	0.0098	0.0096	0.0094
	50	0.0121	0.0121	0.0121	0.0117	0.0114
	75	0.0141	0.0143	0.0142	0.0142	0.0132
	95	0.0184	0.02	0.0188	0.0189	0.0178
P_{02}	25	0.0131	0.0122	0.0107	0.0091	0.0082
	50	0.0165	0.015	0.0131	0.0117	0.0107
	75	0.0188	0.0174	0.0165	0.0155	0.0143
	95	0.0239	0.0244	0.0226	0.0225	0.0244

Table 25. Estimated bias of distribution function estimators at some percentiles for $n = 1600$, the first censoring level.

Transition Probabilities	Percentile	N = 1600				
		k				
		1	4	8	16	32
P_{00}	25	0.0103	0.0094	0.0084	0.0079	0.0078
	50	0.0114	0.0099	0.0093	0.0093	0.0095
	75	0.012	0.0114	0.011	0.0106	0.0104
	95	0.0147	0.0139	0.0135	0.0139	0.0173
P_{01}	25	0.007	0.007	0.0066	0.0068	0.0072
	50	0.0087	0.0084	0.0081	0.0081	0.0086
	75	0.0106	0.0101	0.0101	0.0098	0.0101
	95	0.0134	0.0137	0.013	0.0132	0.0132
P_{02}	25	0.0091	0.0074	0.0063	0.0055	0.0044
	50	0.0112	0.0089	0.0083	0.0071	0.0063
	75	0.0132	0.0114	0.0105	0.0101	0.0093
	95	0.0166	0.0163	0.0148	0.0165	0.0207

Table 26. Estimated bias of distribution function estimators at some percentiles for $n = 200$, the second censoring level.

Transition Probabilities	Percentiles	N = 200		
		k		
		1	2	4
P ₀₀	25	0.0285	0.0269	0.0265
	50	0.0321	0.0324	0.0305
	75	0.0373	0.0366	0.0345
	95	0.0508	0.0475	0.0443
P ₀₁	25	0.0258	0.0253	0.0259
	50	0.0319	0.0318	0.0342
	75	0.0407	0.0378	0.0406
	95	0.0668	0.0601	0.0554
P ₀₂	25	0.0187	0.0178	0.0176
	50	0.0305	0.0262	0.0242
	75	0.0389	0.0332	0.0334
	95	0.0631	0.0544	0.0478

Table 27. Estimated bias of distribution function estimators at some percentiles for $n = 400$, the second censoring level.

Transition Probabilities	Percentiles	N = 400			
		k			
		1	2	4	8
P ₀₀	25	0.02	0.0193	0.0186	0.0176
	50	0.0237	0.0237	0.0217	0.0209
	75	0.0264	0.026	0.0252	0.0237
	95	0.033	0.0363	0.0319	0.0325
P ₀₁	25	0.017	0.0172	0.0178	0.0171
	50	0.0218	0.0227	0.023	0.0218
	75	0.0285	0.0281	0.028	0.0278
	95	0.0448	0.0428	0.0404	0.0389
P ₀₂	25	0.0131	0.013	0.0117	0.01
	50	0.0209	0.0186	0.0174	0.0158
	75	0.0277	0.0238	0.0225	0.0196
	95	0.0409	0.0391	0.0356	0.0325

Table 28. Estimated bias of distribution function estimators at some percentiles for $n = 800$, the second censoring level.

Transition Probabilities	Percentile	N = 800				
		k				
		1	2	4	8	16
P_{00}	25	0.0141	0.0135	0.0123	0.0128	0.0127
	50	0.0165	0.0165	0.015	0.015	0.0142
	75	0.0178	0.0198	0.0174	0.0175	0.0163
	95	0.0233	0.0257	0.0229	0.0228	0.0229
P_{01}	25	0.0121	0.0124	0.0122	0.0127	0.0121
	50	0.0155	0.017	0.0152	0.0155	0.0153
	75	0.021	0.0204	0.0192	0.0197	0.0186
	95	0.0314	0.0316	0.0299	0.0289	0.0289
P_{02}	25	0.0095	0.0091	0.0085	0.0073	0.0066
	50	0.0147	0.0136	0.0119	0.0106	0.009
	75	0.0205	0.0174	0.0152	0.0142	0.0127
	95	0.0275	0.0285	0.0251	0.0233	0.0262

Table 29. Estimated bias of distribution function estimators at some percentiles for $n = 1600$, the second censoring level.

Transition Probabilities	Percentile	N = 1600				
		k				
		1	4	8	16	32
P_{00}	25	0.0105	0.0089	0.0088	0.0085	0.0083
	50	0.0121	0.011	0.0103	0.0098	0.0099
	75	0.0132	0.0126	0.0122	0.0108	0.0114
	95	0.0159	0.0162	0.0157	0.0155	0.0166
P_{01}	25	0.0087	0.0088	0.0086	0.0088	0.0084
	50	0.011	0.0108	0.0107	0.0105	0.01
	75	0.0143	0.014	0.0139	0.0125	0.0123
	95	0.0209	0.0208	0.0209	0.0195	0.0187
P_{02}	25	0.0066	0.0058	0.0053	0.0045	0.0038
	50	0.0105	0.008	0.0076	0.0067	0.0057
	75	0.0133	0.0111	0.0098	0.0088	0.0079
	95	0.0197	0.0174	0.0162	0.017	0.0172

Table 30. Estimated MSEs of transition probabilities P_{11} , P_{12} at some percentiles for $n = 200$, the second censoring level.

Transition Probabilities	Percentiles	N = 200		
		k		
		1	2	4
P_{11}	25	0.0116	0.0125	0.0109
	50	0.0088	0.0068	0.0048
	75	0.006	0.0048	0.003
	95	0.0063	0.0054	0.0052
P_{12}	25	0.0116	0.0125	0.0109
	50	0.0088	0.0068	0.0048
	75	0.006	0.0048	0.003
	95	0.0063	0.0054	0.0052

Table 31. Estimated bias of transition probabilities P_{11} , P_{12} at some percentiles for $n = 200$, the second censoring level.

Transition Probabilities	Percentiles	N = 200		
		k		
		1	2	4
P_{11}	25	0.0841	0.0873	0.0861
	50	0.0752	0.0641	0.0554
	75	0.0624	0.0539	0.0433
	95	0.0636	0.0584	0.0591
P_{12}	25	0.0841	0.0873	0.0861
	50	0.0752	0.0641	0.0554
	75	0.0624	0.0539	0.0433
	95	0.0636	0.0584	0.0591

Table 32. Estimated MSEs of transition probabilities P_{11} , P_{12} at some percentiles for $n = 400$, the second censoring level.

Transition Probabilities	Percentiles	N = 400			
		k			
		1	2	4	8
P_{11}	25	0.0054	0.0074	0.0079	0.0085
	50	0.0041	0.0038	0.0026	0.002
	75	0.0029	0.0023	0.0014	0.0024
	95	0.003	0.0027	0.0026	0.0025
P_{12}	25	0.0054	0.0074	0.0079	0.0085
	50	0.0041	0.0038	0.0026	0.002
	75	0.0029	0.0023	0.0014	0.0024
	95	0.003	0.0027	0.0026	0.0025

Table 33. Estimated bias of transition probabilities P_{11} , P_{12} at some percentiles for $n = 400$, the second censoring level.

Transition Probabilities	Percentiles	N = 400			
		k			
		1	2	4	8
P_{11}	25	0.0583	0.0669	0.0726	0.0787
	50	0.0512	0.0482	0.0416	0.036
	75	0.0435	0.0371	0.0297	0.0367
	95	0.0421	0.0404	0.0417	0.038
P_{12}	25	0.0583	0.0669	0.0726	0.0787
	50	0.0512	0.0482	0.0416	0.036
	75	0.0435	0.0371	0.0297	0.0367
	95	0.0421	0.0404	0.0417	0.038

Cycle 1

$$\mathbf{X}_{(1)11} \leq X_{(2)11} \dots \leq X_{(k-1)11} \leq X_{(k)11} \rightarrow \mathbf{X}_{(1)1}$$

$$X_{(1)12} \leq \mathbf{X}_{(2)12} \dots \leq X_{(k-1)12} \leq X_{(k)12} \rightarrow \mathbf{X}_{(2)1}$$

...

$$X_{(1)1(k-1)} \leq X_{(2)1(k-1)} \dots \leq \mathbf{X}_{(k-1)1(k-1)} \leq X_{(k)1(k-1)} \rightarrow \mathbf{X}_{(k-1)1}$$

$$X_{(1)1k} \leq X_{(2)1k} \dots \leq X_{(k-1)1k} \leq \mathbf{X}_{(k)1k} \rightarrow \mathbf{X}_{(k)1}$$

Cycle 2

$$\mathbf{X}_{(1)21} \leq X_{(2)21} \dots \leq X_{(k-1)21} \leq X_{(k)21} \rightarrow \mathbf{X}_{(1)2}$$

$$X_{(1)22} \leq \mathbf{X}_{(2)22} \dots \leq X_{(k-1)22} \leq X_{(k)22} \rightarrow \mathbf{X}_{(2)2}$$

...

$$X_{(1)2(k-1)} \leq X_{(2)2(k-1)} \dots \leq \mathbf{X}_{(k-1)2(k-1)} \leq X_{(k)2(k-1)} \rightarrow \mathbf{X}_{(k-1)2}$$

$$X_{(1)2k} \leq X_{(2)2k} \dots \leq X_{(k-1)2k} \leq \mathbf{X}_{(k)2k} \rightarrow \mathbf{X}_{(k)2}$$

...

Cycle m

$$\mathbf{X}_{(1)m1} \leq X_{(2)m1} \dots \leq X_{(k-1)m1} \leq X_{(k)m1} \rightarrow \mathbf{X}_{(1)m}$$

$$X_{(1)m2} \leq \mathbf{X}_{(2)m2} \dots \leq X_{(k-1)m2} \leq X_{(k)m2} \rightarrow \mathbf{X}_{(2)m}$$

...

$$X_{(1)m(k-1)} \leq X_{(2)m(k-1)} \dots \leq \mathbf{X}_{(k-1)m(k-1)} \leq X_{(k)m(k-1)} \rightarrow \mathbf{X}_{(k-1)m}$$

$$X_{(1)mk} \leq X_{(2)mk} \dots \leq X_{(k-1)mk} \leq \mathbf{X}_{(k)mk} \rightarrow \mathbf{X}_{(k)m}$$

Figure 1. An RSS procedure of obtaining a sample size of $k \times m$.

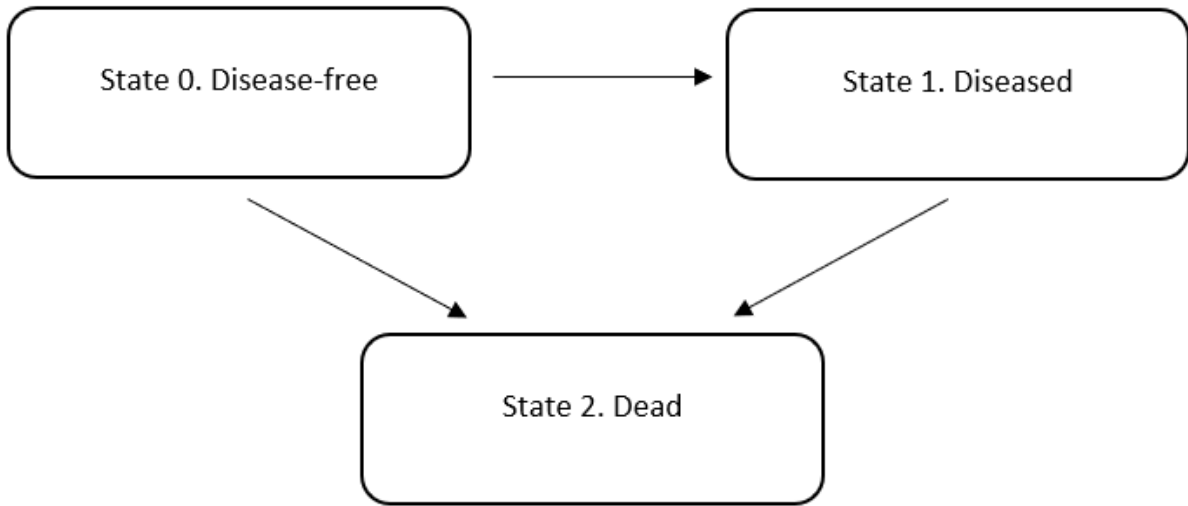


Figure 2. Illustration of progressive illness-death model.

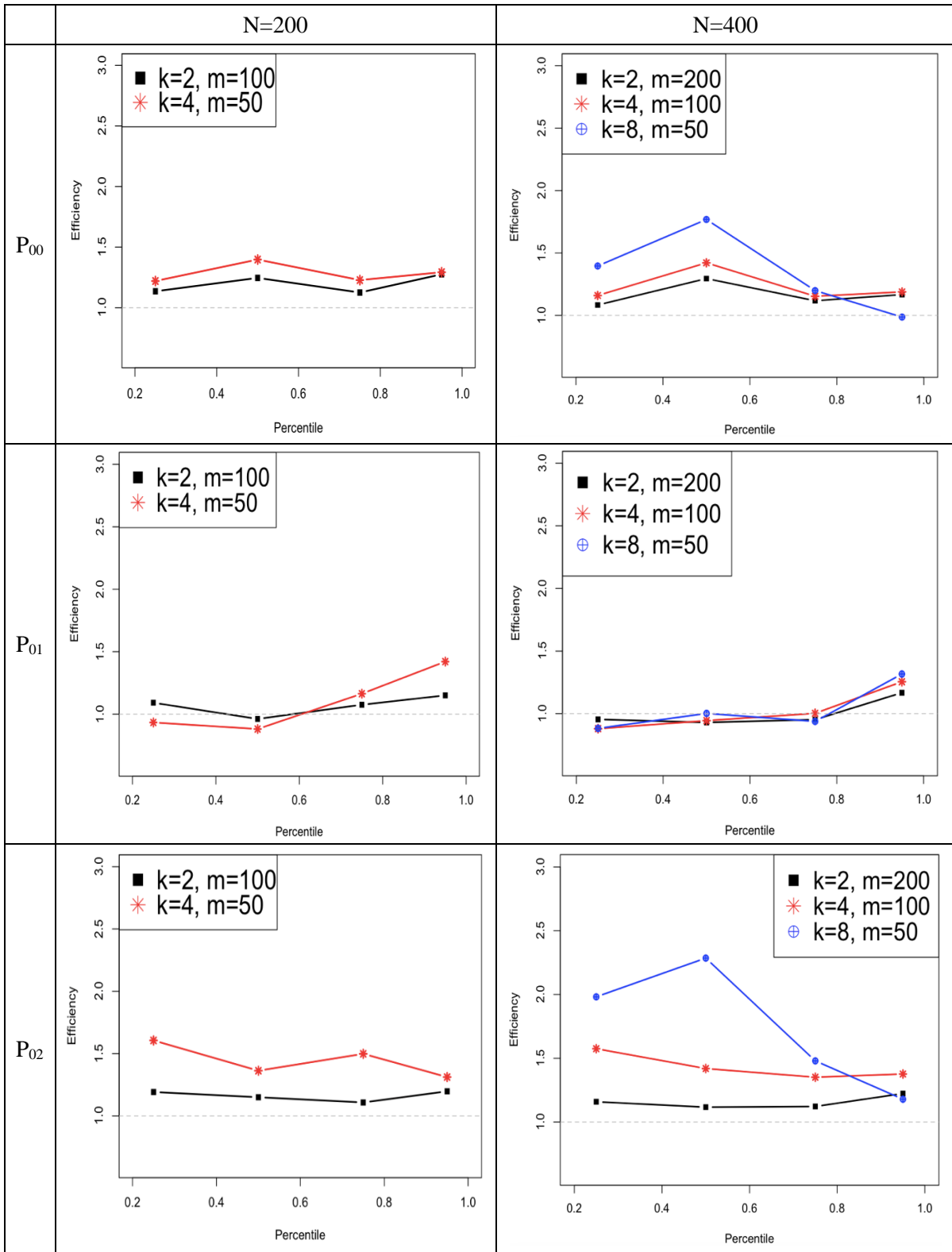


Figure 3. The efficiency of the AJ estimators based on RSS with respect to SRS counterparts at different percentiles for $n = 200, 400$, the first censoring level.

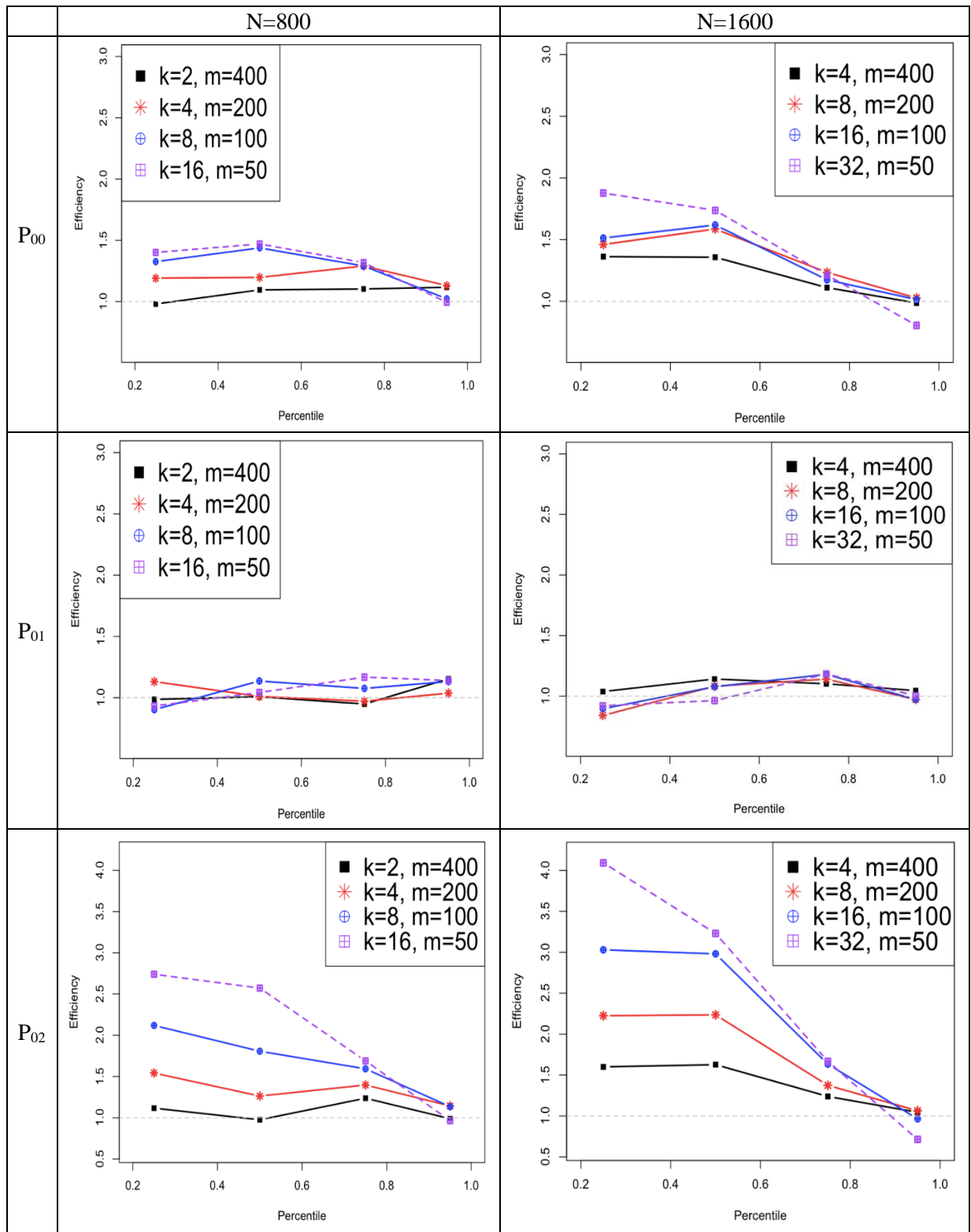


Figure 4. The efficiency of the AJ estimators based on RSS with respect to SRS counterparts at different percentiles for $n = 800, 1600$, the first censoring level.

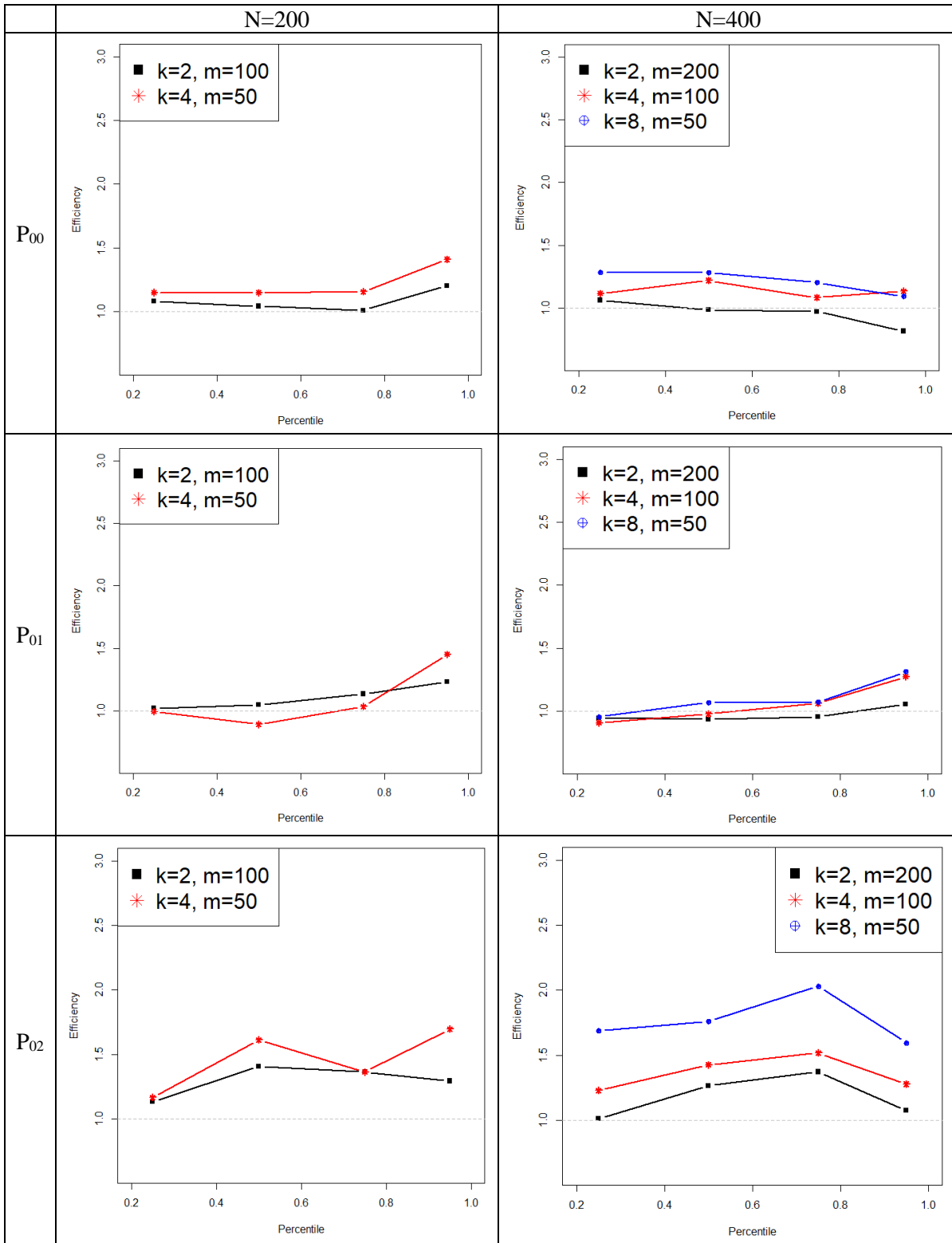


Figure 5. The efficiency of the AJ estimators based on RSS with respect to SRS counterparts at different percentiles for $n = 200, 400$, the second censoring level.

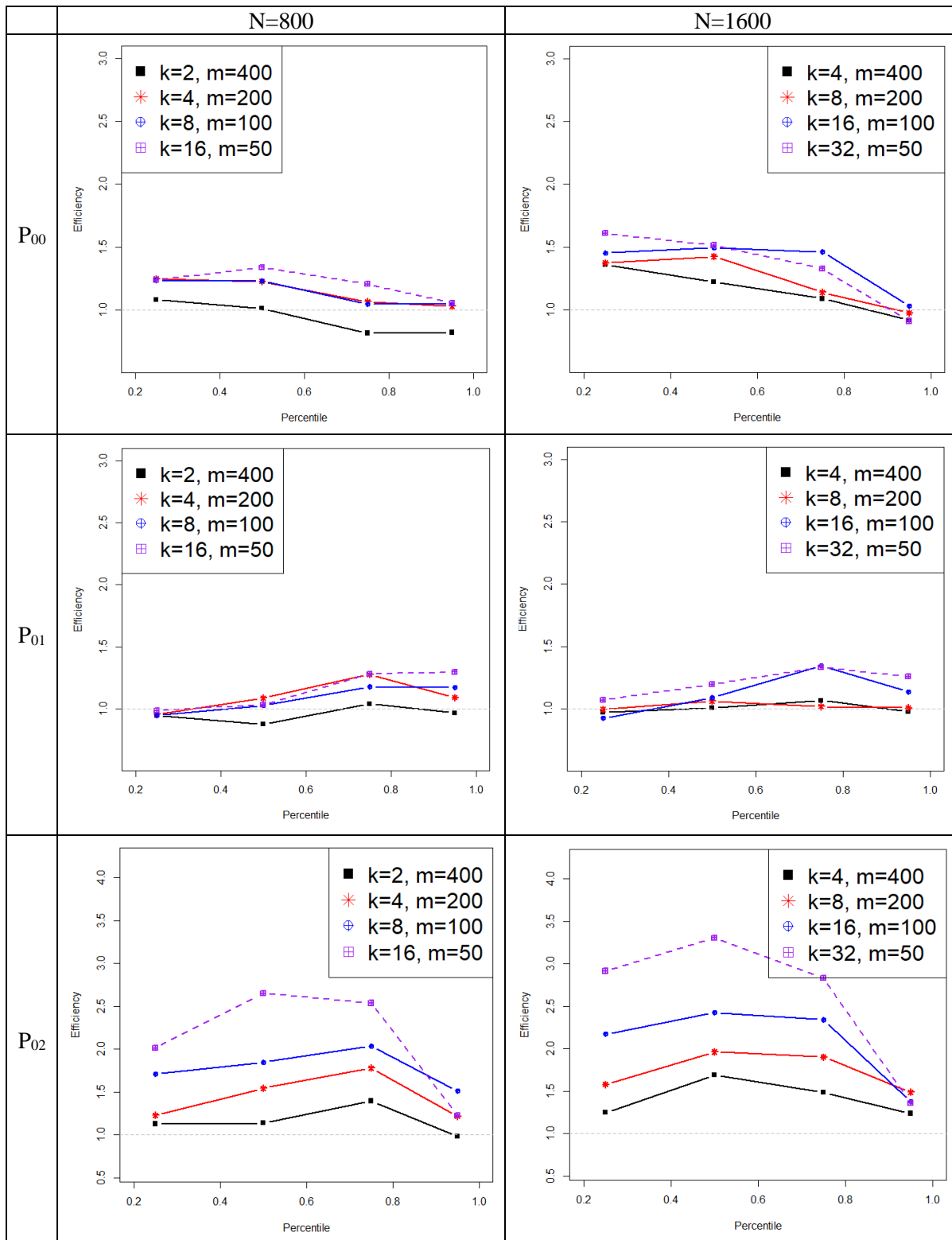


Figure 6. The efficiency of the AJ estimators based on RSS with respect to SRS counterparts at different percentiles for $n = 800, 1600$, the second censoring level.

3 Chapter II: Application of Aalen Johansen Estimator in Healthy Start Dataset

3.1 Introduction

3.1.1 Small for Gestational Age

Small for gestational age (SGA) is defined as an infant weighing less than the 10th percentile at birth.^{62,63} Some reference define SGA as a neonate whose weight and/or crown-heel length at birth is at least 2 standard deviation below the mean for the infant's gestational age, which is equivalent to the 2.3 percentile.⁶⁴ Appropriate for gestational age (AGA) is defined as infants between ninetieth and tenth percentile weight at birth.⁶⁵ There are three categories of SGA: a low birth, a low birth length or both low birth weight and length.⁶⁶ Infants born with SGA have increased risk of respiratory complications, hypotension, hypoglycemia, necrotizing enterocolitis, perinatal and later life morbidity and mortality.⁶⁴

Factors contributing to SGA include: maternal high blood pressure, chronic kidney disease, advanced diabetes, heart or respiratory disease, malnutrition, infection, substance use, cigarette smoking.⁶⁷ Kramer review 895 publications on risk factors of low birth weight.⁶⁸ Factors with direct causal impact on intrauterine growth retardation (IUGR) include infant sex, racial/ethnic origin, maternal height, pre-pregnancy weight, paternal weight and height, maternal birth weight, parity, history of prior low-birth-weight infants, gestational weight gain and caloric intake, general morbidity and episodic illness, malaria, cigarette smoking, alcohol consumption, and tobacco chewing. Black or Indian racial origin is one of the major determinants of IUGR for developing countries, while cigarette smoking is the leading factor for developed countries. For most of SGA infants, multi-factors contribute the result, such as women who smoke tend to be younger, thinner and from lower socioeconomic classes.^{68,69} The impact of socioeconomic level is suggested to be indirect, since women from lower socioeconomic groups tend to have increased smoking level.⁷⁰ Some socioeconomic indicator may have independent impact on SGA outcome, such as education. Studies have found that in both East and West Germany, maternal education

was associated with SGA delivery.⁷¹ Improved education stands for better self-care, sufficient utilization of health care system and more knowledge of health-related behavior. It is reported that pregnancy outside marriage had small but significant increase in adverse birth outcome including SGA, low birth weight and preterm.⁷² Several studies have found women older than 35 years old have increased risk of SGA, but not significant after adjusting for confounders.⁶³ Birth weight was reduced for first and second time teen mothers compared with reference group in England.⁷³ However, emotional support and intimacy for the teen mothers from family and friends represent highly preventive intervention for adverse birth outcomes.⁷⁴

3.1.2 Healthy Start Project

The term disparity is interpreted from racial or ethnic perspective for most of the time, but disparities present in many dimensions in healthcare system. Healthy People describes that “Health disparities adversely affect groups of people who have systematically experienced greater obstacles to health based on their racial or ethnic group; religion; socioeconomic status; gender; age; mental health; cognitive, sensory, or physical disability; sexual orientation or gender identity; geographic location; or other characteristics historically linked to discrimination or exclusion.”⁷⁵ Objectives in Healthy People 2010 include designing interventions to reduce illness, disability and premature death, with a goal to eliminate health disparities.⁷⁶ The nationwide Federal Government Healthy Start (HS) Project is a widely known program to reduce infant mortality and morbidity among disadvantage populations. The Central Hillsborough Healthy Start Project (CHHS) is a community-based program funded through the Maternal and Child Health Bureau’s Healthy Start Initiative by the Federal Government in the 17-census tract area of Tampa, Florida.⁷⁷ Many of the mothers are young, black, out of marriage and undereducated, which contributes to the poor birth outcomes in the area.⁷⁸

Perinatal risk reduction services were provided collaboratively from CHHS and Florida Department of Health. To identify those who will benefit from the service, risk screens were offered to pregnant

women and newborns. Mothers who were interested in having HS services took the screen test and were referred to local HS programs.

Healthy Start services are mainly composed of three parts: 1) initial contact 2) initial assessment 3) care co-ordination.⁷⁷ Within 5 days after screening and referral, initial contact is performed with participants. In this contact, HS team members review the potential risk factors of adverse birth outcomes, explain concerns, and determine future services with the client. Resource utilization issues, client's ability to access health care are also considered in this stage. Within 10 days after initial contact, initial assessment is attempted, which is a face-to-face contact. During this stage, professional assessment of social and biological risk factors is carried out. Assistance to overcome the risks are determined; health education about fetal development, normal/abnormal pregnancy, maternal nutrition, child spacing, preterm labor signs are performed; how clients access qualified resources are discussed; a comprehensive plan to ensure a smooth pregnant process till delivery is offered. Care co-ordination is a prenatal assessment and care utilization package addressing specific risk factors to prevent adverse maternal health and birth outcomes.

Previous study found CHHS Health Start program reduce the risk for low birth weight (OR=0.7; 95% CI=0.5-1.0) and preterm births (OR=0.7; 95% CI=0.5-0.9) for service recipients as compared to non-recipients significantly.⁷⁷ The study measured program effectiveness using odds ratios from logistic regression and number needed to treat (NNT). Another study involved central Hillsborough Health Start data from 2000 to 2007 indicates women exposed to air particulate pollutants had elevated risk for low birth weight (AOR=1.24; 95% CI=1.07–1.43), very low birth weight (AOR=1.58; 95% CI=1.09–2.29) and preterm (AOR=1.18; 95% CI=1.03–1.34).⁷⁹ And this adverse effect of air particulate pollutants was improved by Central Hillsborough Federal Healthy Start Project. The Federal Healthy Start program in collaboration with community partners in east Tampa contributed to the decline in primary teen pregnancy (10-19 years), but no impact on repeated pregnancy.⁸⁰ A study involved Healthy Start data in the above region from 2002-2009 linked with local vital record showed that pregnant women with the shortest interpregnancy interval (0-5 months: AOR=1.39, 95% CI 1.23-1.56) and longest interpregnancy

(IPI) (≥ 60 months: AOR=1.13, 95% CI 1.03-1.23) have an increased risk of adverse birth outcomes compared with IPI (18-23 months) group.⁸¹ Even though the role of Healthy Start in the association between interpregnancy interval and fetoinfant morbidities is unclear, the study has important public health implications. Providing education and counseling services for women to make optimized decision regarding birth space is a recommended public health strategy to reduce fetoinfant morbidities.⁸²

3.1.3 Illness-Death Model Applied to Healthy Start Project for SGA Outcome

The Markov models are a class of stochastic models that assume a finite number of health states (clusters) and allows movement or transition from one state to the other.³⁸ The rate of movement from one state to the next are measured in terms of transition probabilities.³⁹ The transition probability from state h to state j ($0 \leq s < t$) is represented mathematically as

$$P_{hj}(s, t) = P(X(t) = j | X(s) = h, H_{s-})$$

H_{s-} represents all the historical information from the data along the interval $[0, s)$. The model in equation 1 is assumed to be independent on H_{s-} . This is the memoryless property which must be satisfied by all Markov models. The property states that any future evolution of the Markov process depends only on its current state and is independent of the previously visited states.⁴⁰

The model applied to this study is similar to the progressive illness death/disability model.⁴¹ The model assumes that individuals from an initial state may transit to an intermediate state and may finally enter a terminal state. This model has been widely used in medical research to study the course of disease.⁴³ The setting of this study is similar to that of the illness death/disability model in which we investigate the likelihood of having a small for gestational age (SGA)⁶⁴⁻⁶⁶ infant (terminal state) in pregnant women (initial state). These women had the opportunity to participate in the healthy start program before delivery (intermediate state).

Previous studies on this topic have utilized a cross sectional design to investigate the association between socio-demographic characteristics and adverse birth outcomes.^{62,77} The longitudinal form of the

illness death/disability model is very attractive because it describes the progress of subjects from an initial state to a final state and have not been previously applied to maternal and child health settings.

In this study, we apply the illness-death model to study the progression of SGA^{65,83} infants in pregnant women who may participate in the healthy start program during follow-up. The impact of risk factors such as teen-mothers, obesity, race, marital status, smoking⁸⁴ and education on the effect of SGA was also investigated.^{63,67,85-88}

3.2 Methods

3.2.1 A Markov Model to Estimate Transition Probabilities for Healthy Start Study and Small for Gestational Age

A Markov model was used in this study to address possible transitions between states that include participant pregnancy (state 0), choosing to have healthy start service (state 1) and delivering a SGA infant (state 2), which corresponds with a typical three states progressive illness-death model well. There are three possible transitions among them: $0 \rightarrow 1$, $0 \rightarrow 2$, $1 \rightarrow 2$. At initial time, all subjects are in state 0, and they are supposed to reach the final absorbing state 2 at future time point, along the process, they may experience or not an intermediate state (state 1). In this study the intermediate state represents choosing to have healthy start service, the time spent in state 0 is referred to as no service and no SGA infant time. In this study, Aalen Johansen⁵³ approach based on Markov assumptions was exploited to estimate the transition probabilities.

Two sets of transition probabilities are to be estimated: for $0 \leq s < t$, $\{P_{0j}(s, t), j = 0,1,2\}$ and, for $0 < s < t$, $\{P_{1j}(s, t), j = 1,2\}$. N independent trajectories corresponding to n individuals are supposed to be observed.

3.2.2 Statistical Analysis

A Markov model is a stochastic model⁸⁹ with two properties: for each pair of states i, j at each instant of time n and $n+1$, the transition probability p_{ij} depends only on the current state i , no matter what the previous states are.⁹⁰ Therefore:

$$p_{ij} = P(X_{n+1} = j | X_n = i), \quad i, j \in S$$

$$P(X_{n+1} = j | X_n = i, X_{n-1} = i_{n-1}, \dots, X_0 = i_0) = P(X_{n+1} = j | X_n = i) = p_{ij}$$

If there are m states total, then the transition probability matrix which is also called stochastic matrix could be listed as:⁹⁰

$$\begin{bmatrix} p_{00} & p_{01} & \dots & p_{0m} \\ p_{10} & p_{11} & \dots & p_{1m} \\ \dots & \dots & \dots & \dots \\ p_{m0} & p_{m1} & \dots & p_{mm} \end{bmatrix}$$

For any s, t with 0 ≤ s < t, for Markov models we have

$$P_{hj}(s, t) = P(X(t) = j | X(s) = h, H_{s^-}) = P(X(t) = j | X(s) = h)$$

The future of the process after time s depends only on the state at time s, which is an important property made the dream of efficient estimation of transition probabilities come true.

Andersen et al.⁴⁰ defined the integrated hazard matrix:

$$A = (A_{hj})$$

$$\text{where } A_{hj}(t) = \int_0^t \alpha_{hj}(s) ds \text{ for all } h, j \text{ with } \alpha_{hh}(t) = - \sum_{h \neq j} \alpha_{hj}(t).$$

Then the transition probability matrix is given as:

$$P(s, t) = P_{hj}(s, t) = \prod_{(s,t]} (I + dA(u))$$

Where A = (A_{hij}). The Nelson-Aalen estimator of \hat{A}_{hj} is defined as:

$$\hat{A}_{hj}(t) = \begin{cases} \int_0^t J_h(u) (Y_h(u))^{-1} dN_{hj}(u), & h \neq j, \\ - \sum_{h \neq j} \hat{A}_{hj}(t), & h = j, \end{cases}$$

Where $J_h(u) = I(Y_h(u) > 0)$. $Y_h(u)$, the number of individuals observed in state h just prior time u.

$N_{hj}(u)$, the number of observed direct transitions from h to j in the time interval [0, u]. The estimated

N*N transition probability matrix is:⁴³

$$\hat{P}(s, t) = \prod_{(s,t]} (I + d\hat{A}(u))$$

This is called Aalen-Johansen estimator.⁵³ The variance of Aalen Johansen estimator could be calculated by a Greenwood-type formula. The covariance matrix $\hat{P}(s, t)$ is:^{40,43}

$$\overline{COV}(\hat{P}(s, t)) = \int_s^t \hat{P}(u, t)^T \otimes \hat{P}(s, u -) \overline{COV}(d\hat{A}(u)) \hat{P}(u, t) \otimes \hat{P}(s, u -)^T$$

Where $\overline{COV}(d\hat{A})$ is the covariance of the matrix $d\hat{A}$.

Linear confidence interval for $\hat{P}_{hj}(s, t)$ is defined:

$$\hat{P}_{hj}(s, t) \pm z_{\frac{\alpha}{2}} * \hat{\sigma}_{hj}(s, t)$$

Where $\hat{\sigma}_{hj}(s, t)$ is the empirical standard error, and $z_{\frac{\alpha}{2}}$ is the upper $\alpha/2$ quantile of the standard normal distribution. Log transformation, log log transformation and complementary log log transformation for small sample size were listed as following:^{91,92}

$$\begin{aligned} & \hat{P}_{hj}(s, t) \exp \left\{ \frac{\pm z_{\frac{\alpha}{2}} * \hat{\sigma}_{hj}(s, t)}{\hat{P}_{hj}(s, t)} \right\} \\ & \hat{P}_{hj}(s, t) \exp \left\{ \frac{\pm z_{\frac{\alpha}{2}} * \hat{\sigma}_{hj}(s, t)}{\hat{P}_{hj}(s, t) \log(\hat{P}_{hj}(s, t))} \right\} \\ & 1 - (1 - \hat{P}_{hj}(s, t)) \exp \left\{ \frac{\pm z_{\frac{\alpha}{2}} * \hat{\sigma}_{hj}(s, t)}{(1 - \hat{P}_{hj}(s, t)) \log(1 - \hat{P}_{hj}(s, t))} \right\} \end{aligned}$$

3.2.3 Population Examined in the Study

We conducted a retrospective cohort study using a statewide enhanced maternal-infant database that contains socio-demographic and perinatal information. Birth and death vital records from the Florida Department of Health for children born from January 1, 2009, through December 31, 2017 were considered. This dataset contains information on 25,161 pregnant women. Gestational age was calculated in weeks and computed by taking the interval between the date of last menstrual period reported by the mother at her first prenatal visit and the date of delivery. In the Florida database,

the menstrual estimates of gestational age of 19.7% pregnancies were inconsistent with birth weight (for example very low birth weight at term). For these cases, a clinical estimate computed by the physician was used.⁹³ The covariates examined as possibly affecting transition probabilities include: teen-mothers, marriage status, race, obesity, smoking and education. All cases with missing values on the above covariates were excluded from the dataset. Mothers age less than 19 years are considered teen mothers, otherwise are non-teen mothers. There are two race groups: white and non-white. Mothers whose pre-pregnancy BMI was larger than 30 is considered obese otherwise non-obese. Smoking groups are divided by maternal smoking or not. Mothers who had one year, two years, somewhat college, earned degree and above are considered college, otherwise non-college.

3.2.4 Setting

To model the above data, we consider a three-state Markov model consisting of pregnancy (state 0), choosing healthy start service (state 1) and delivery (state 2) (Figure 7). At the start of the study (state 0), all subjects are pregnant. Between pregnancy and delivery, mothers could decide whether to have healthy start service (state 1) or not. At delivery (state 2), these mothers will either have SGA or AGA infants. Study subjects can have transition between states, however, once a subject enters state 2 (absorption state), that subject will remain in this state with probability 1.⁹⁴

From the study design, five possible transition probabilities could be identified from this model:

- P_{01} , the probability of transition from state 0 to state 1 (from pregnancy to receiving healthy start service);
- P_{02} , the probability of transition from state 0 to state 2 (from pregnancy to delivering an SGA infant without receiving healthy start service);
- P_{12} , the probability of transition from state 1 to state 2 (received healthy start service and delivered an SGA infant);
- P_{00} , the probability of staying in state 0. Note that this is the transition probability that a woman did not choose healthy start service and delivered an AGA infant;

- P_{11} , the probability of staying in state 1, which is the probability that a woman had healthy start service and delivered AGA infant.

From transition probability matrix, the sum of P_{00} , P_{01} , P_{02} is equal to one. The sum of P_{11} and P_{12} is one. Once in state 2, the probability of getting out is 0 since it is an absorption state. Therefore, P_{22} is equal to one.

		From		
		State 0	State 1	State 2
	State 0	P_{00}	P_{01}	P_{02}
To	State 1	0	P_{11}	P_{12}
	State 2	0	0	P_{22}

3.3 Results and Discussion

3.3.1 Results

The data used for this study included 25,161 pregnant women. Of these, 3,409 (13.55%) received healthy start service. A total number of 3,607 (14.34%) deliveries were SGA (Table 34).

Table 35 shows demographic characteristics of all study participants. The mean age of study participants was 26 years with a minimum of 12 and a maximum age of 51 years. Additionally, 7.44% of study participants were below 19 years old. The overall mean pre-pregnancy BMI of study subjects was 26.5(SD=6.78) with a mean gestational age of 38.21(SD=2.38). Also, 29.14% of study participants were married, 40.46% were white, 24.31% were obese, 4.63% were smokers and 33.78% had one year or more college education.

Table 36 presents transition probabilities at 2, 3, 4, 5, 6, 7, 8, 9 months from pregnancy to delivery. As can be observed from this table, prior to 5 months (no delivery) the only transition was from state 0 (pregnancy) to state 1 (HS). There was no transition between state 0 and state 2 or state 1 and state 2 (SGA infant), meaning no delivery took place during this time. After 5 months (as study participants

deliver) the transition probability from states 0 to 2 and from states 1 to 2 increases with gestational age. At the end of the study, 28.1% of study subjects who did not receive HS service had SGA infant whereas 26.2% of subjects that received HS services had SGA infants with a risk reduction of about 2% (Table 36).

Table 37 shows transition probabilities categorized by selected risk factors to investigate the impact of HS services on SGA. Overall, women that received HS service had lower probabilities of having SGA infants compared to women that did not have HS service. Six risk factors were considered: teen mothers, marriage status, race, obesity, smoking status, and education. Additionally, for all risk factors considered, except for obese mothers, receiving HS service narrows the gap in the transition probability of having an SGA baby in high risk mothers compared to low risk mothers (Figure 8). As can be observed from Figure 8, teenage mothers that did not receive HS service had 5.6% higher probabilities of having SGA infants compared to non-teenage mothers. None married mothers that did not receive HS service had 12.9% higher probabilities of having SGA infants compared to married mothers. Non-white mothers that did not receive HS services have 16% higher probability compared to white mothers for having SGA infants. This rate reduces to 9% after receiving HS services. Smoking mothers that did not receive HS service had 13% higher probabilities of having SGA infants compared to non-smoking mothers. However, after receiving HS service, this difference reduces to 10%. Non-college mothers that did not receive HS services had 9% higher probability compared to college mothers for having SGA infants. This proportion reduced to 8% after receiving service.

Figure 9 to figure 20 present transition probabilities of teen mothers, marriage groups, race groups, obesity status, smoking status, and education status. These plots, which are visual display of results from Table 37, conforms well with this table. For all the plots, transition probabilities P_{00} , P_{11} showed decreasing trend along time, while P_{02} , P_{12} showed increasing trend. For plots of marriage, race, smoking and education statuses, there were obvious differences between the risk groups in transition probabilities P_{00} . However, for all the above groups, the plots of P_{11} for two groups entangled together, which indicates similarity in transition probabilities between groups.

3.3.2 Discussion

In this study we have successfully applied the Markov model to investigate the effectiveness of the healthy start program in relation to the delivery of an SGA infant. Our results show that the risk of delivering an SGA infant was greatly reduced when a mother received healthy start services compared to mothers that did not receive healthy start services. Using a similar data, Salihu et al.⁷⁷ observed the healthy start program to be effective in reducing preterm birth but ineffective in impacting SGA. To our knowledge, this is the only study that has shown the health start program in reducing SGA.

Another important observation in this study is that the gap in the rate of SGA in high risk mothers that participated in the healthy start program compared to low risk mothers that also participated in the healthy start program is filled. Additionally, the rate of SGA was almost similar in high risk mothers that participated in the healthy start program compared to low risk mothers that did not take part in the healthy start program. To our knowledge, these observations have not been previously reported by others on the effectiveness of the healthy start.

A previous study on the effectiveness of the healthy start program did not find the Hillsborough healthy start program to be effective in reducing the rates of SGA. Since 2009, structural changes have occurred in the healthy start program in Hillsborough County, Florida. The University of South Florida, Tampa was sub-contracted to take care of the evaluation part of the program, as well as the revamping of the database might have resulted in quality data. Another reason why the rate of SGA is significantly lower for program participants might be attributed to the large sample size that this study possessed compared to previous study. This study tracked nine years of high-risk population data from 2009 to 2017 resulting to almost doubling of the power compared to previous studies on this program.

Another advantage that this study commands is the use of the Markov model to display the transition probabilities over a period of time. Most study on this topic have been cross sectional. The longitudinal design of this study ensures that we followed women over time and observe their movement throughout their pregnancy period. Unlike traditional longitudinal models focusing on trends over time, the method used in this paper provides a novel angle to analyze longitudinal data as we have characterized discrete

sequences based on their transition patterns. The model demonstrates a satisfactory and stable performance.

Additionally, the whole study process was expressed in a three states Markov model clearly and naturally. The number of states the system requires is not large, which makes the system reliable and stable.⁹⁵ Markov models have advantage in sequence dependent behavior. For example, certain events have to take place until other events occur.⁹⁵ In this study, all three states are followed by time sequence, with the first state happened at pregnancy time, the second state at choosing service time, the third state as delivery time. The number of states in this study is not large, therefore, computation won't require too much memory and execution time. Correctly specifying states is not a challenge for this research. Markov assumption is restrictive, but this study generally meets the assumption.

Table 34. Number of records for percentage of women choosing services and percentage of newborns with SGA in healthy start data set.

Delivery Status	Number of records	Percentage
Total	25,161	100.00%
Women choosing services	3,409	13.55%
Newborns with SGA	3,607	14.34%

Table 35. Summary of demographic statistics in healthy start data set.

Demographics	Statistics	Value
	N	25161
	Mean (SD)	26.05 (5.87)
Age (yrs)	Median	25.00
	Min.	12.00
	Max.	51.00
	N	23603
	Mean (SD)	26.50 (6.78)
Pre-pregnancy BMI	Median	24.90
	Min.	13.70
	Max.	76.00
	N	25161
	Mean (SD)	38.21 (2.38)
Gestational weeks	Median	39.00
	Min.	20.00
	Max.	43.00
Age (<19)	N (%)	1873 (7.44)
Married (Yes)	N (%)	7327 (29.14)
White (Yes)	N (%)	10179 (40.46)
Obese (Yes)	N (%)	5736 (24.31)
Smoking (Yes)	N (%)	1164 (4.63)
Somewhat College (Yes)	N (%)	8470 (33.78)

Table 36. Transition Probability at 2, 3, 4, 5, 6, 7, 8, 9, and last day* from pregnancy in healthy start data set.

Time (days)	State 0 to 0	State 0 to 1	State 0 to 2	State 1 to 2	State 1 to 2
60	97.112%	2.888%	0.000%	100%	0.000%
90	93.576%	6.424%	0.000%	100%	0.000%
120	91.655%	8.345%	0.000%	100%	0.000%
150	90.269%	9.714%	0.016%	100%	0.000%
180	89.267%	10.648%	0.084%	99.964%	0.036%
210	88.445%	11.378%	0.177%	99.899%	0.101%
240	87.078%	11.967%	0.955%	99.210%	0.790%
270	79.821%	11.471%	8.707%	91.952%	8.048%
294*	62.591%	9.353%	28.057%	73.809%	26.191%

*: the last woman giving birth in the study is 294 days from getting pregnancy.

Table 37. Transition probabilities by treatment groups (%) in healthy start data set.

Covariates		State 0 to 0	State 0 to 1	State 0 to 2	State 1 to 1	State 1 to 2
Age	Teen mother (Age < 19)	63.70 (59.35, 68.05)	2.91 (1.15, 4.66)	33.40 (0.00, 100.00)	73.43 (64.19, 82.67)	26.57 (17.33, 35.81)
	Non-teen mother (Age ≥ 19)	69.90 (66.32, 73.48)	2.32 (2.04, 2.60)	27.78 (0.00, 64.41)	74.52 (71.10, 77.94)	25.48 (22.06, 28.90)
Marriage	Not Married	64.73*** (59.02, 70.45)	2.51 (2.14, 2.87)	32.76 (0.00, 100.00)	72.67 (68.74, 76.60)	27.33 (23.40, 31.26)
	Married	78.19*** (76.09, 80.30)	1.99 (0.00, 4.81)	19.82 (0.00, 54.33)	80.50 (76.08, 84.92)	19.50 (15.08, 23.92)
Race	Non-white	62.46*** (57.39, 67.53)	2.37 (1.80, 2.95)	35.17 (0.00, 75.08)	71.09*** (67.77, 74.40)	28.91*** (25.60, 32.23)
	White	78.05*** (76.40, 79.69)	2.55 (2.25, 2.86)	19.40 (0.00, 68.89)	79.95*** (74.48, 85.42)	20.05*** (14.58, 25.52)
BMI	Obese (BMI>30)	74.35 (71.50, 77.21)	1.66 (0.86, 2.45)	23.99 (0.00, 86.14)	75.88 (65.68, 86.07)	24.12 (13.93, 34.32)
	Non-obese (BMI≤30)	67.54 (62.91, 72.17)	2.59 (2.23, 2.94)	29.87 (0.00, 82.39)	73.73 (70.32, 77.15)	26.27 (22.85, 29.68)
Smoking	Smoking	56.90*** (49.07, 64.73)	1.84 (0.80, 2.88)	41.26 (0.00, 100.00)	64.55 (44.71, 84.39)	35.45 (15.61, 55.29)
	Non-smoking	69.81*** (66.17, 73.44)	2.38 (2.07, 2.68)	27.81 (0.00, 64.78)	75.43 (72.46, 78.41)	24.57 (21.59, 27.54)
Education	Non-college	66.14*** (62.02, 70.25)	2.58 (2.15, 3.00)	31.28 (0.00, 71.77)	72.12 (67.93, 76.32)	27.88 (23.68, 32.07)
	College	76.32*** (74.28, 78.36)	1.93 (1.65, 2.22)	21.74 (0.00, 100.00)	80.00 (76.04, 83.95)	20.00 (23.68, 32.07)

***: significant difference between treatment groups.

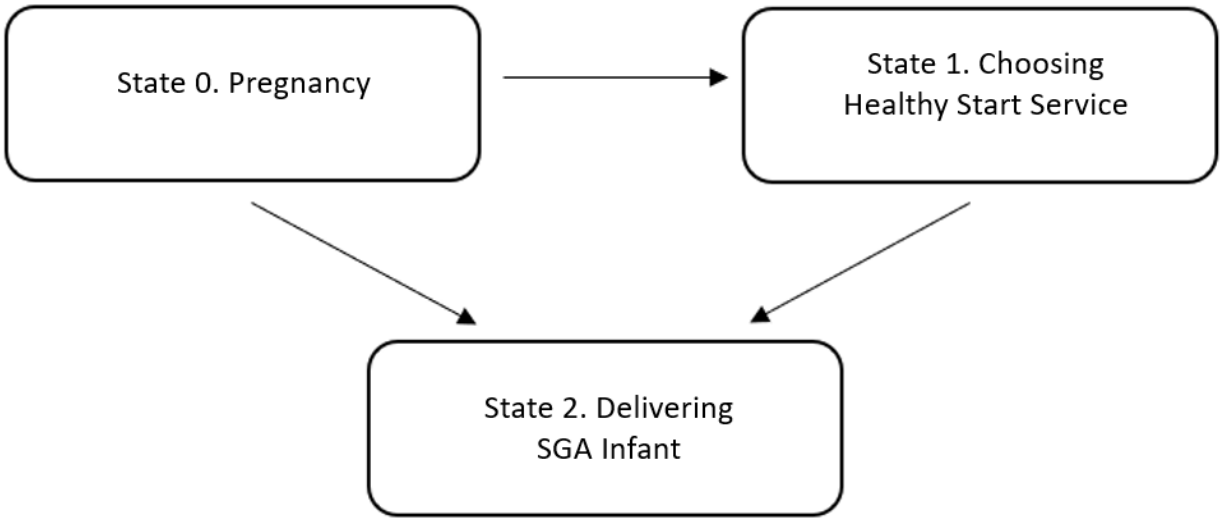


Figure 7. Illustration of progressive illness-death model applied to Healthy Start project.

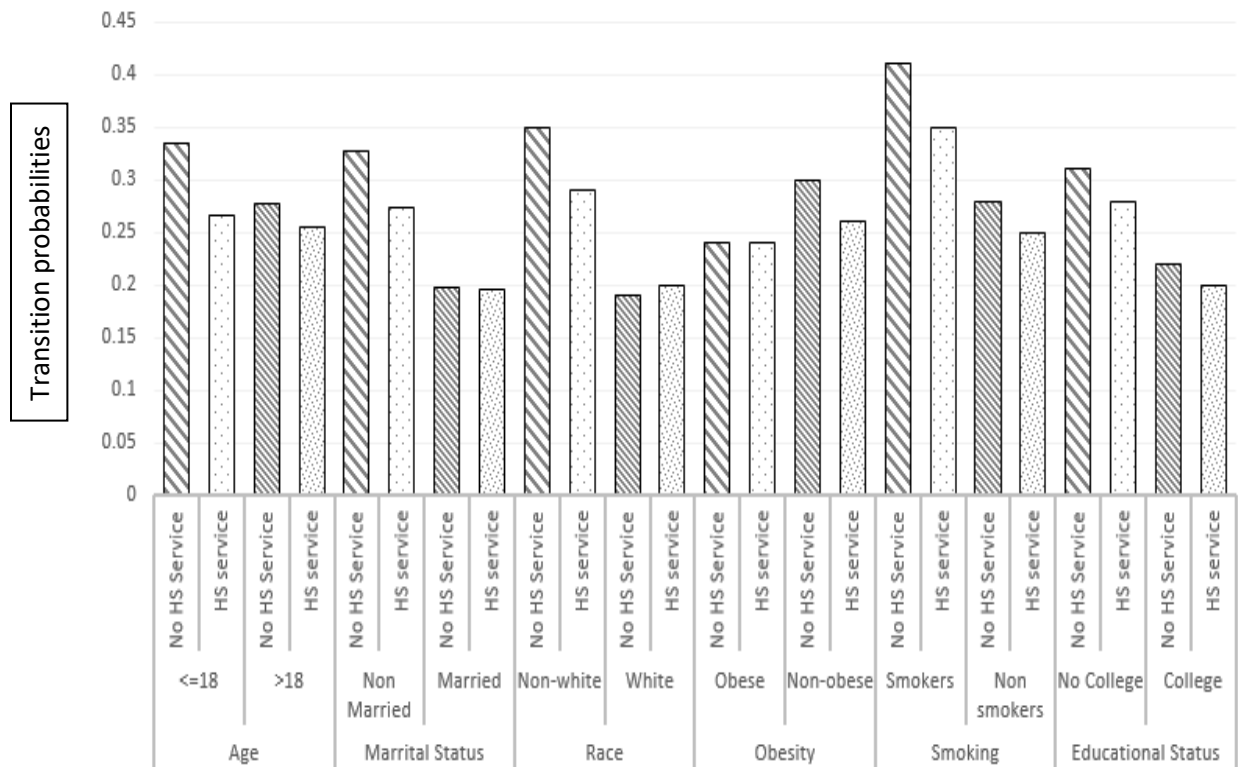


Figure 8. Transition probabilities of having an SGA infant by risk groups in healthy start data set.

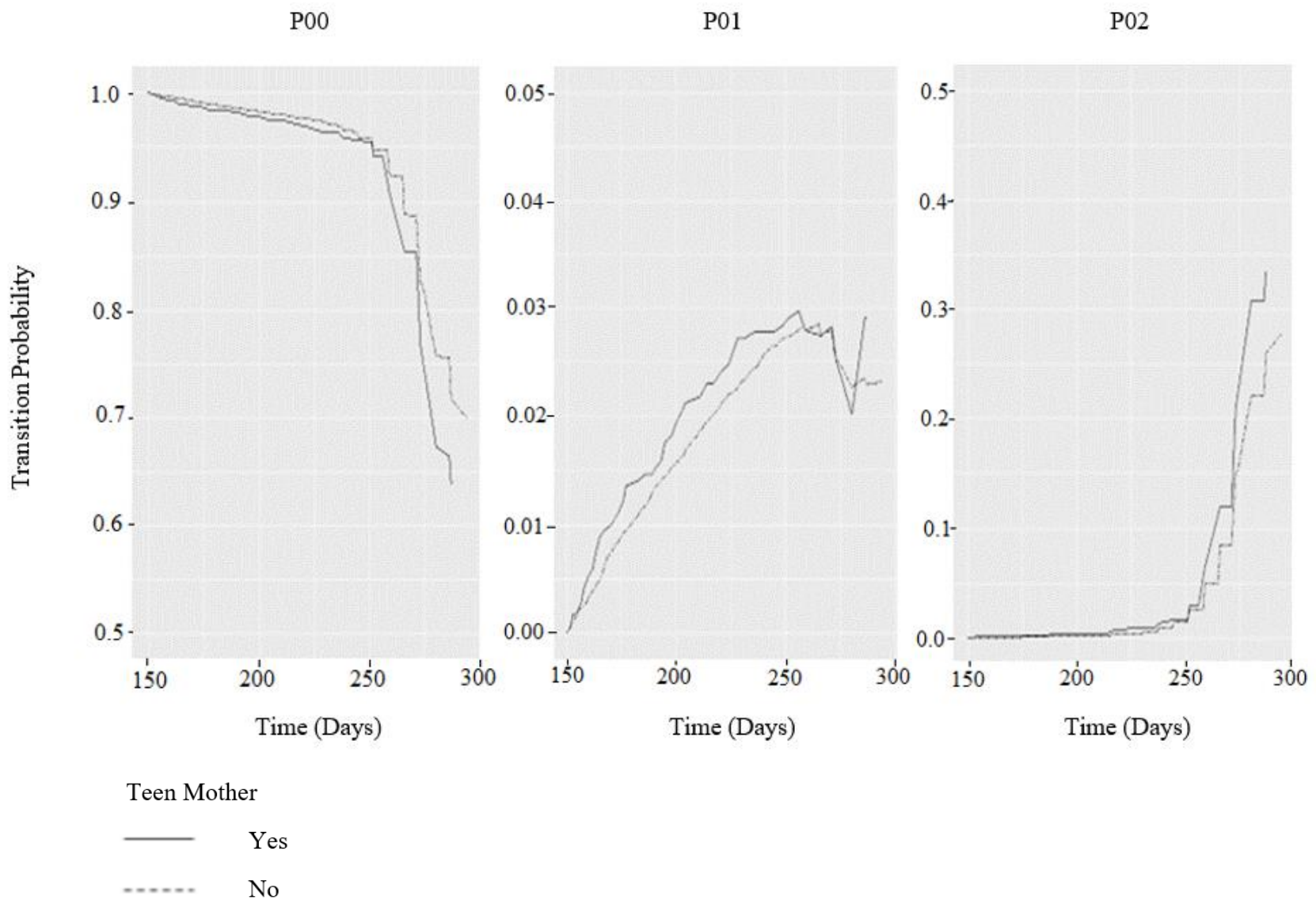


Figure 9. Transition probabilities P00, P01 and P02 plots for teen mothers in healthy start data set.

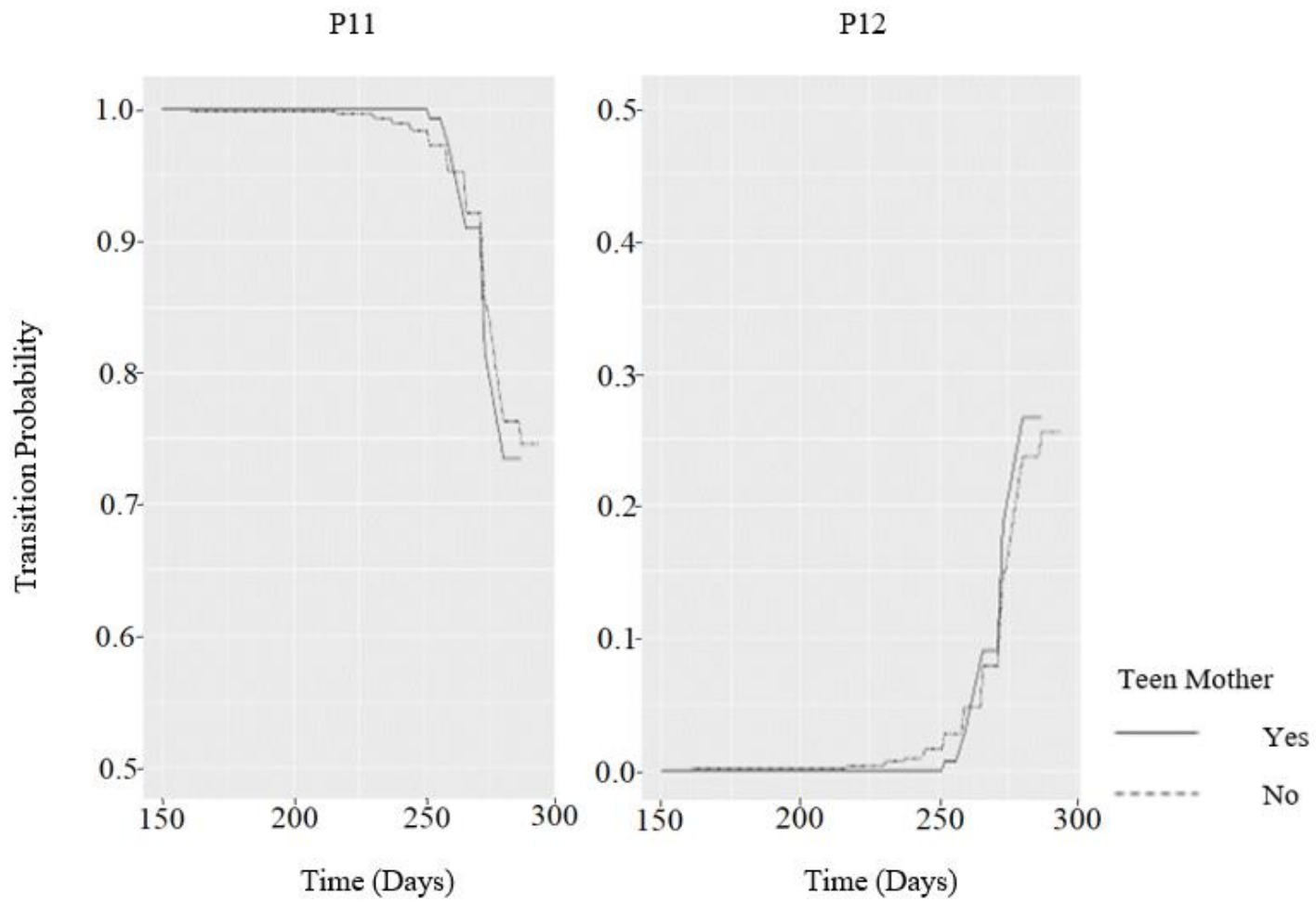


Figure 10. Transition probabilities P11 and P12 plots for teen mothers in healthy start data set.

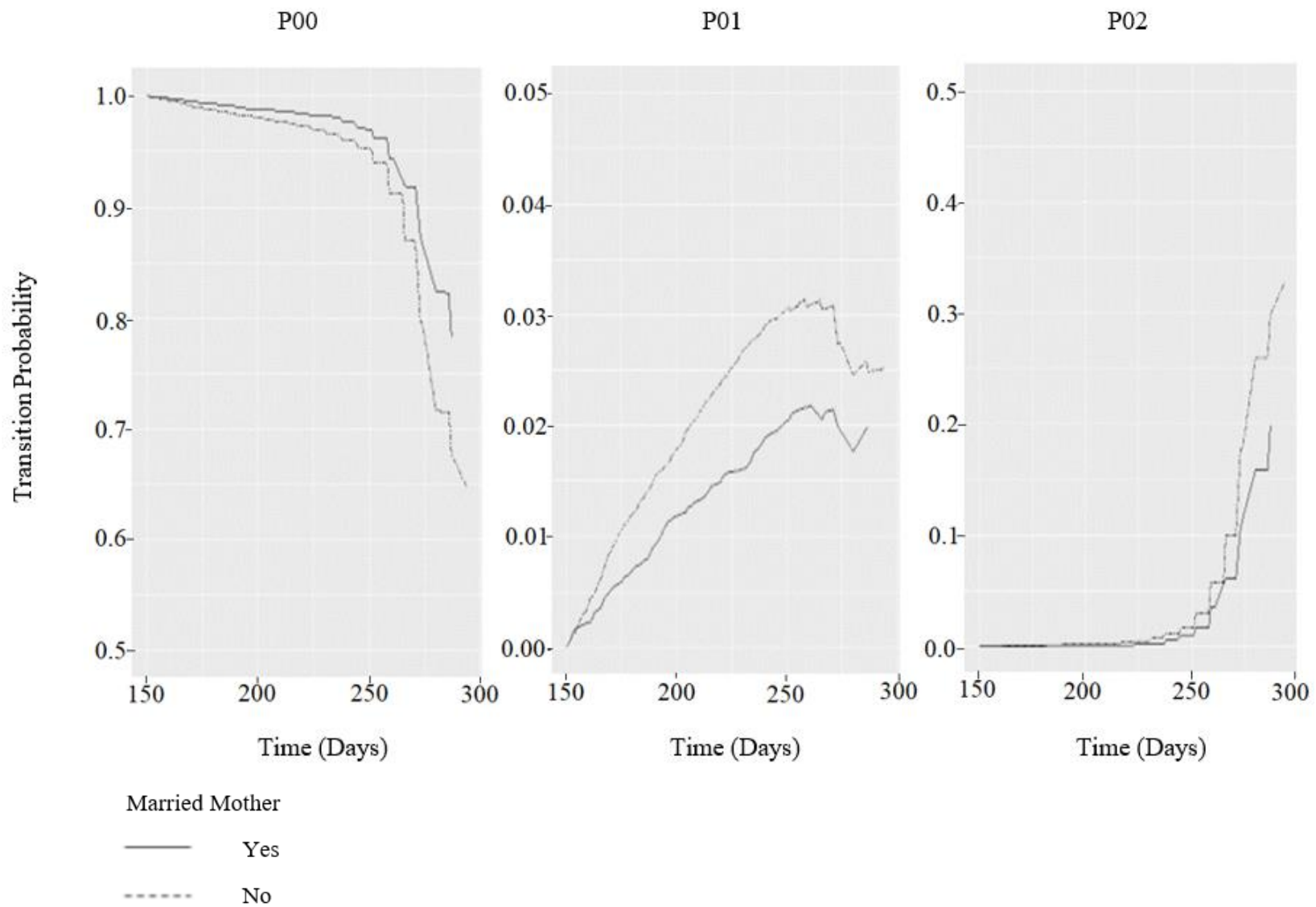


Figure 11. Transition probabilities P00, P01 and P02 plots for marriage groups in healthy start data set.

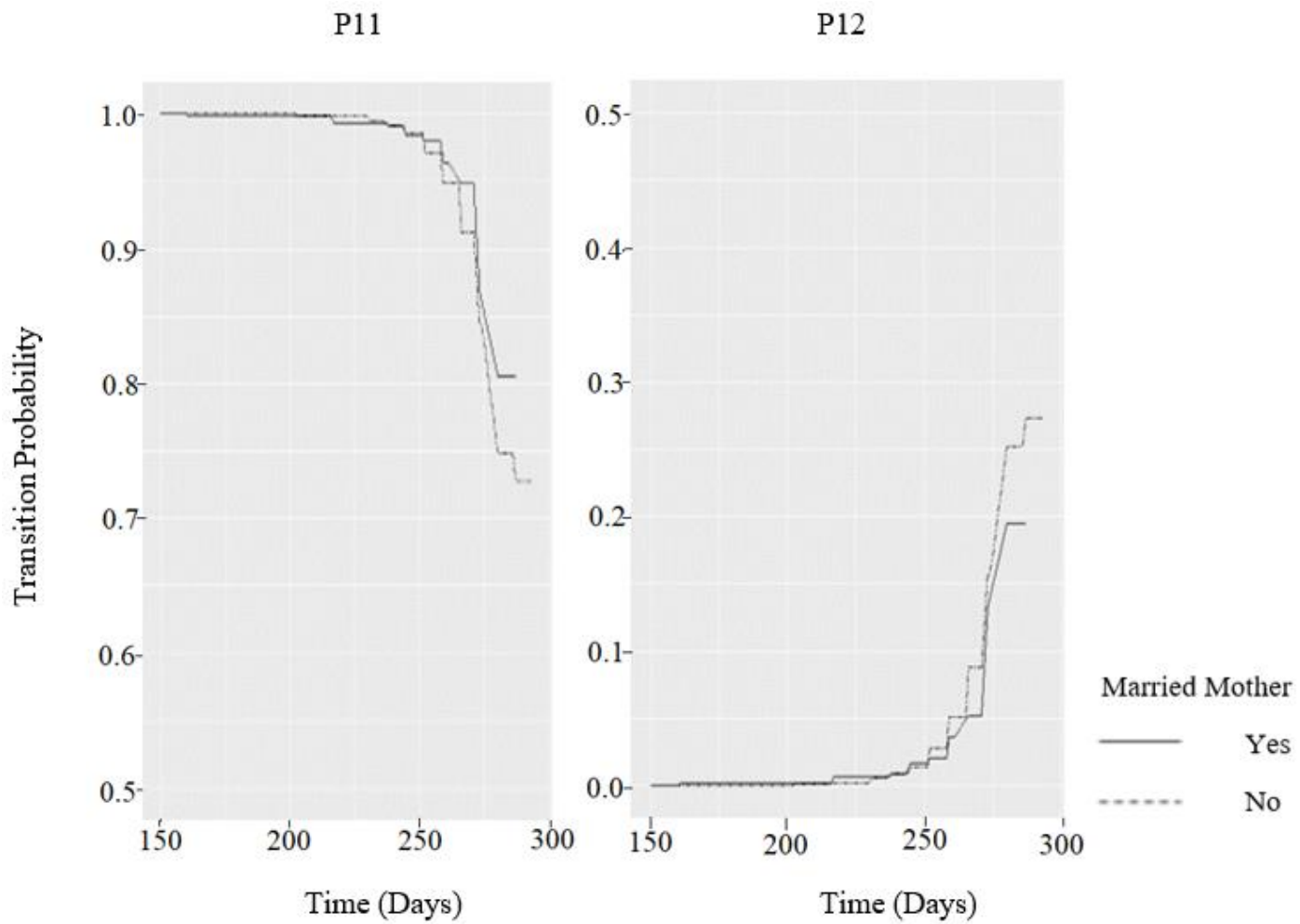


Figure 12. Transition probabilities P11 and P12 plots for marriage groups in healthy start data set.

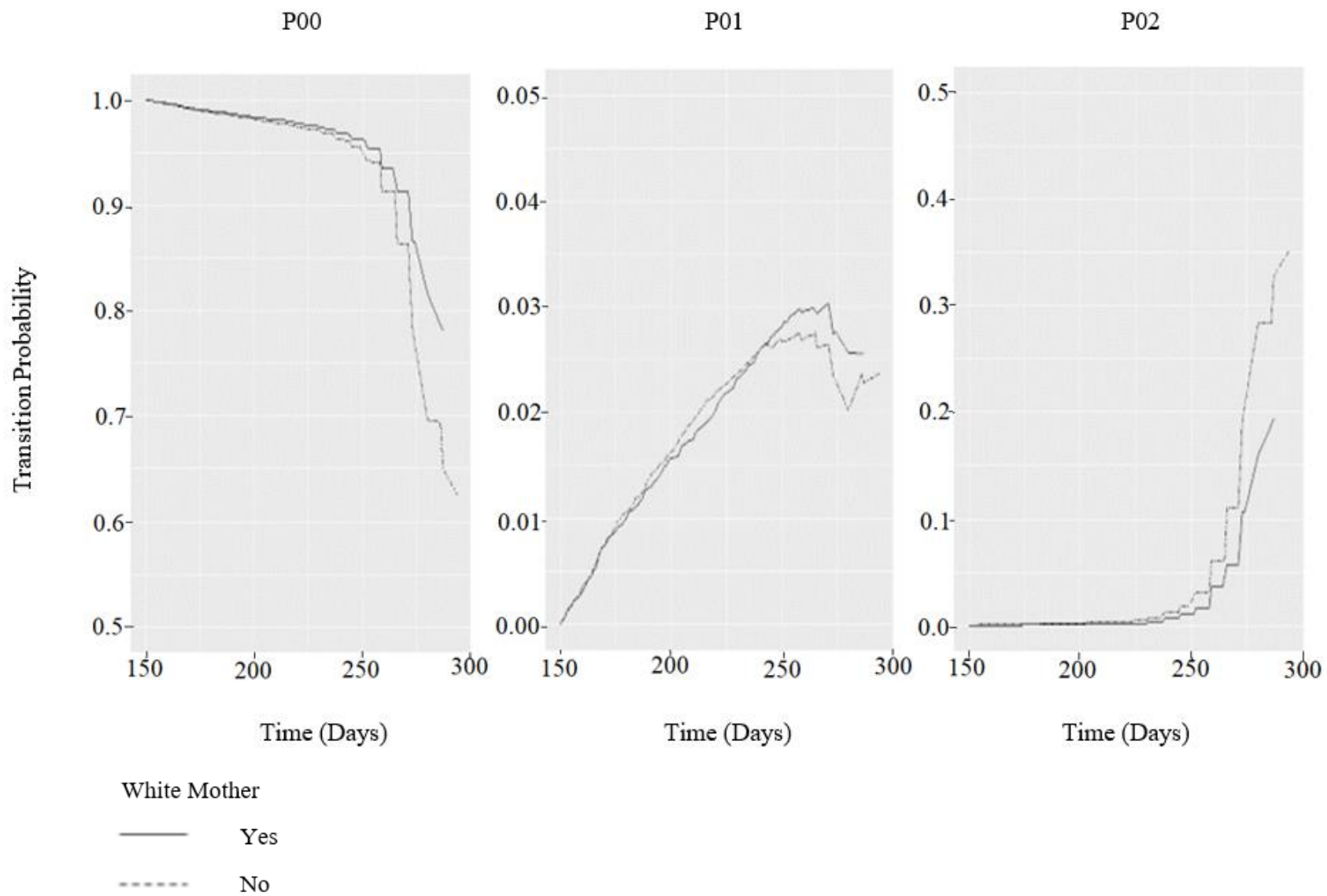


Figure 13. Transition probabilities P00, P01 and P02 plots for race groups in healthy start data set.

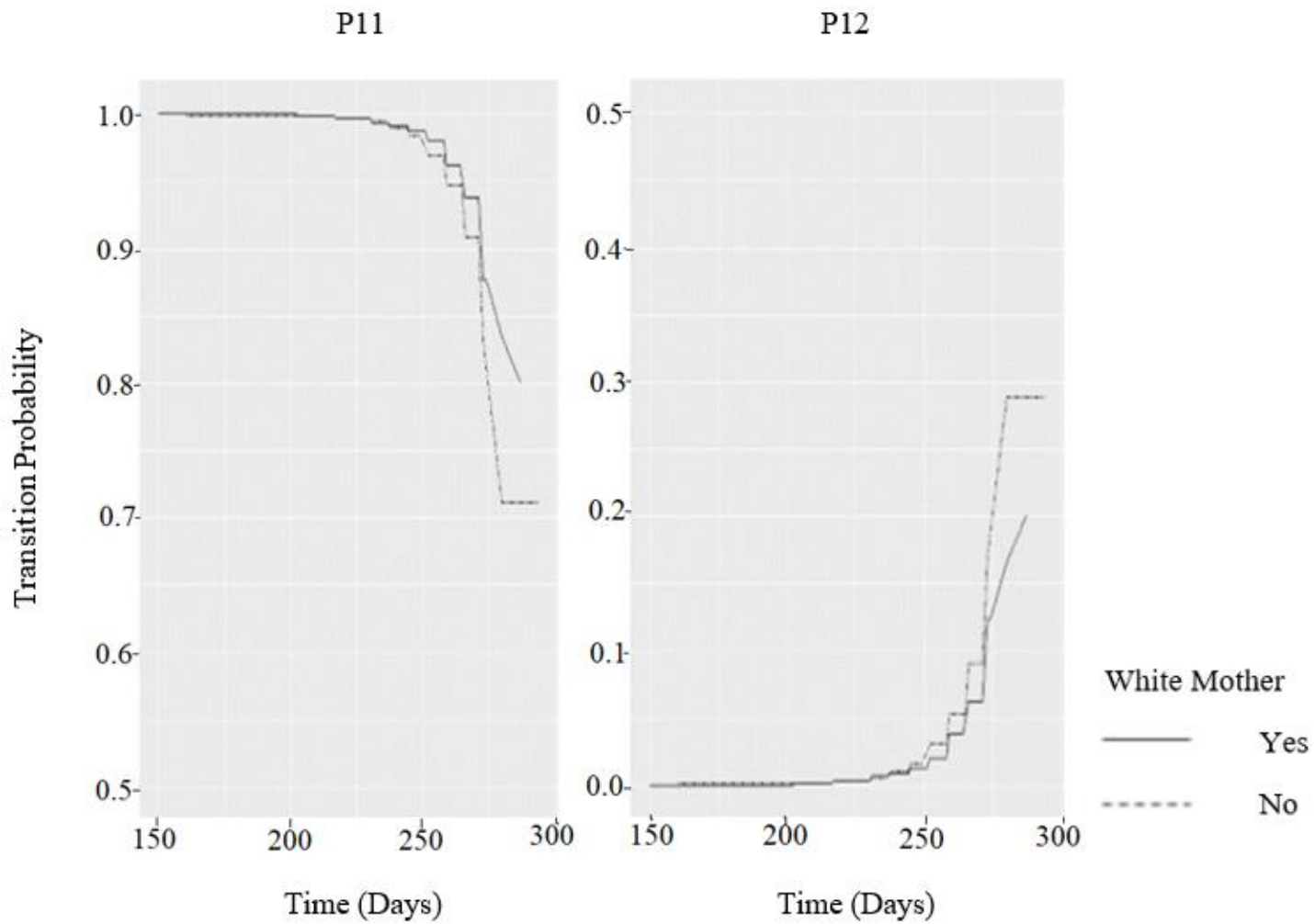


Figure 14. Transition probabilities P11 and P12 plots for race groups in healthy start data set.

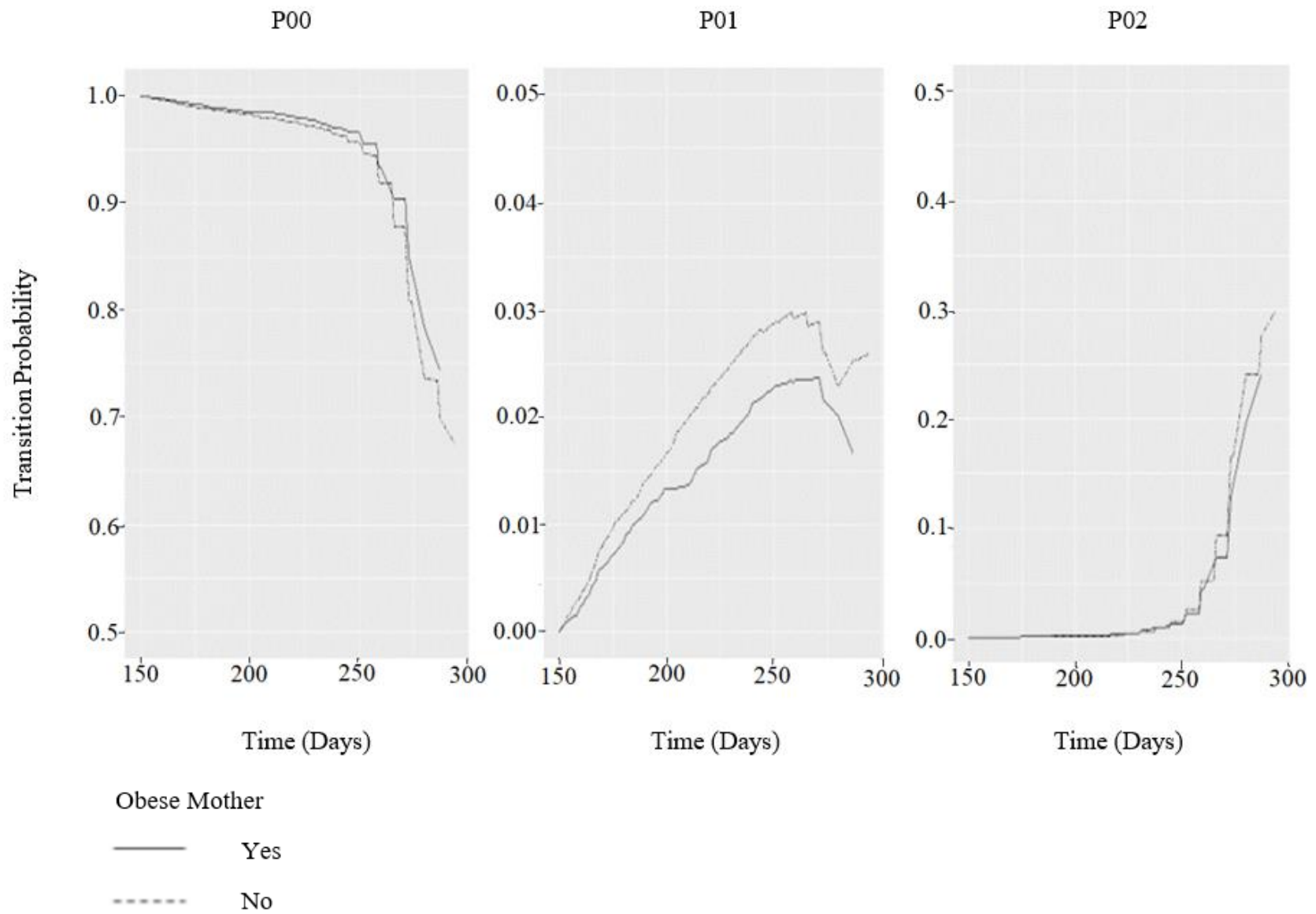


Figure 15. Transition probabilities P00, P01 and P02 plots for obese status in healthy start data set.

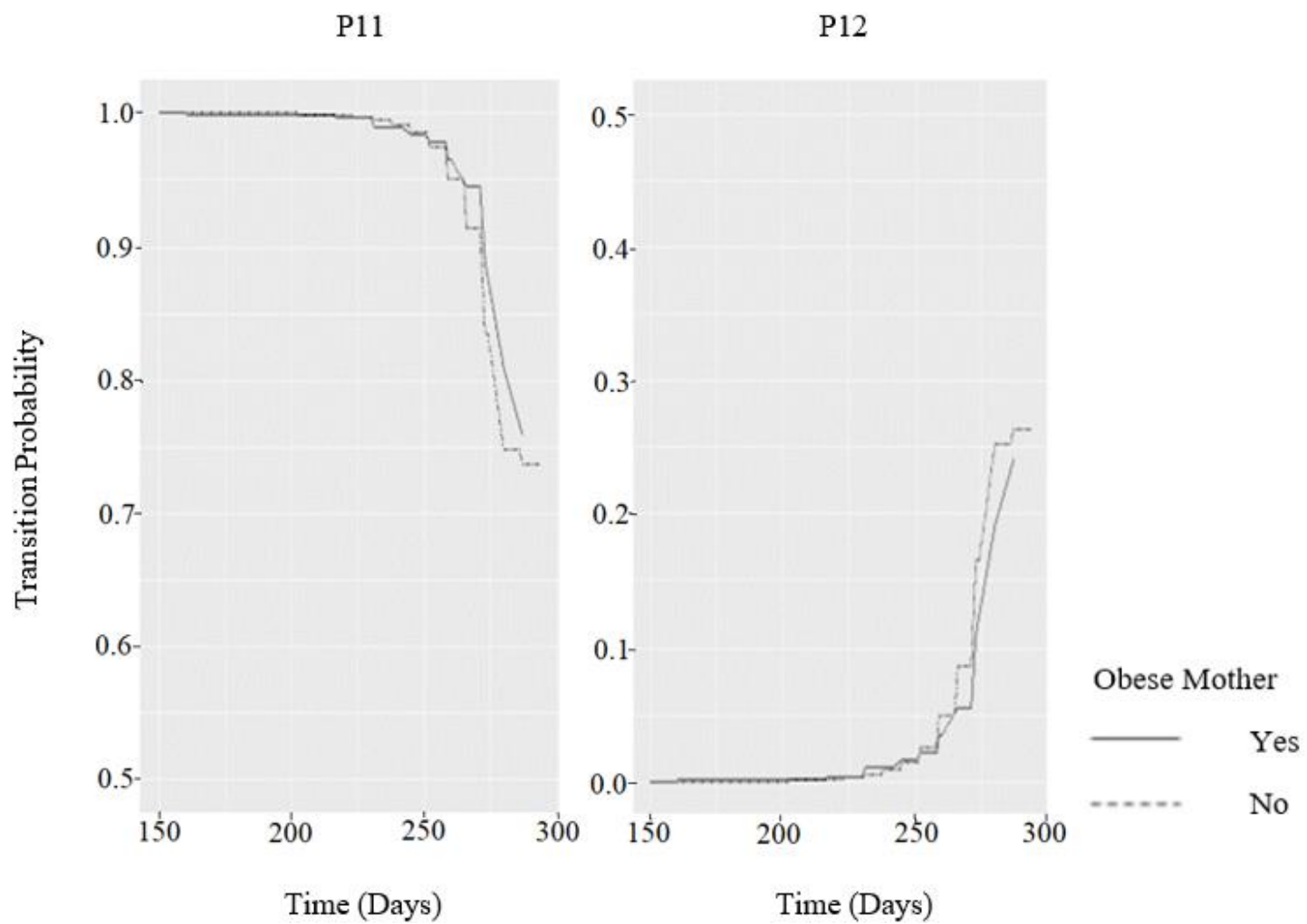


Figure 16. Transition probabilities P11 and P12 plots for obese status in healthy start data set.

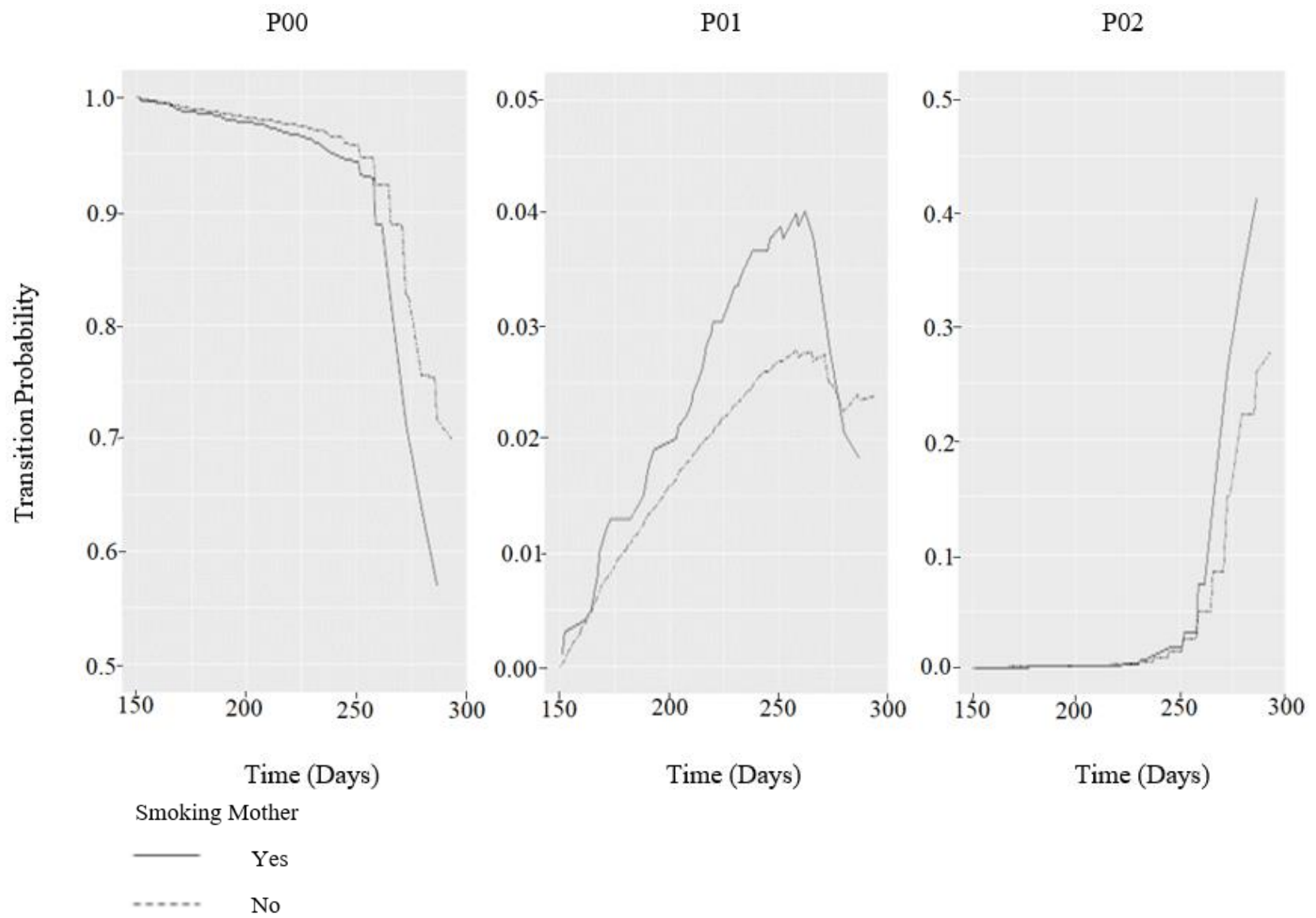


Figure 17. Transition probabilities P00, P01 and P02 plots for smoking groups in healthy start data set.

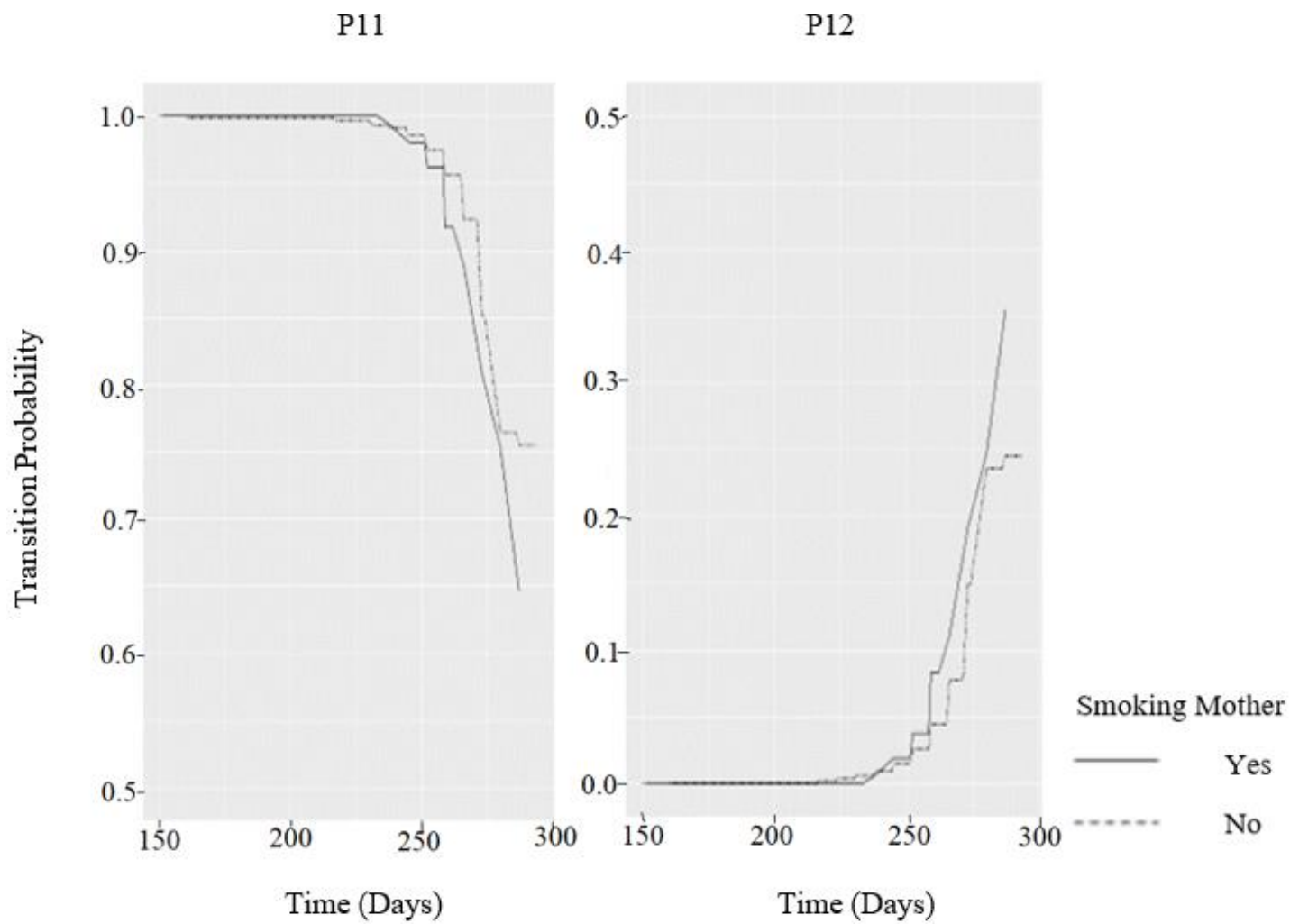


Figure 18. Transition probabilities P11 and P12 plots for smoking groups in healthy start data set.

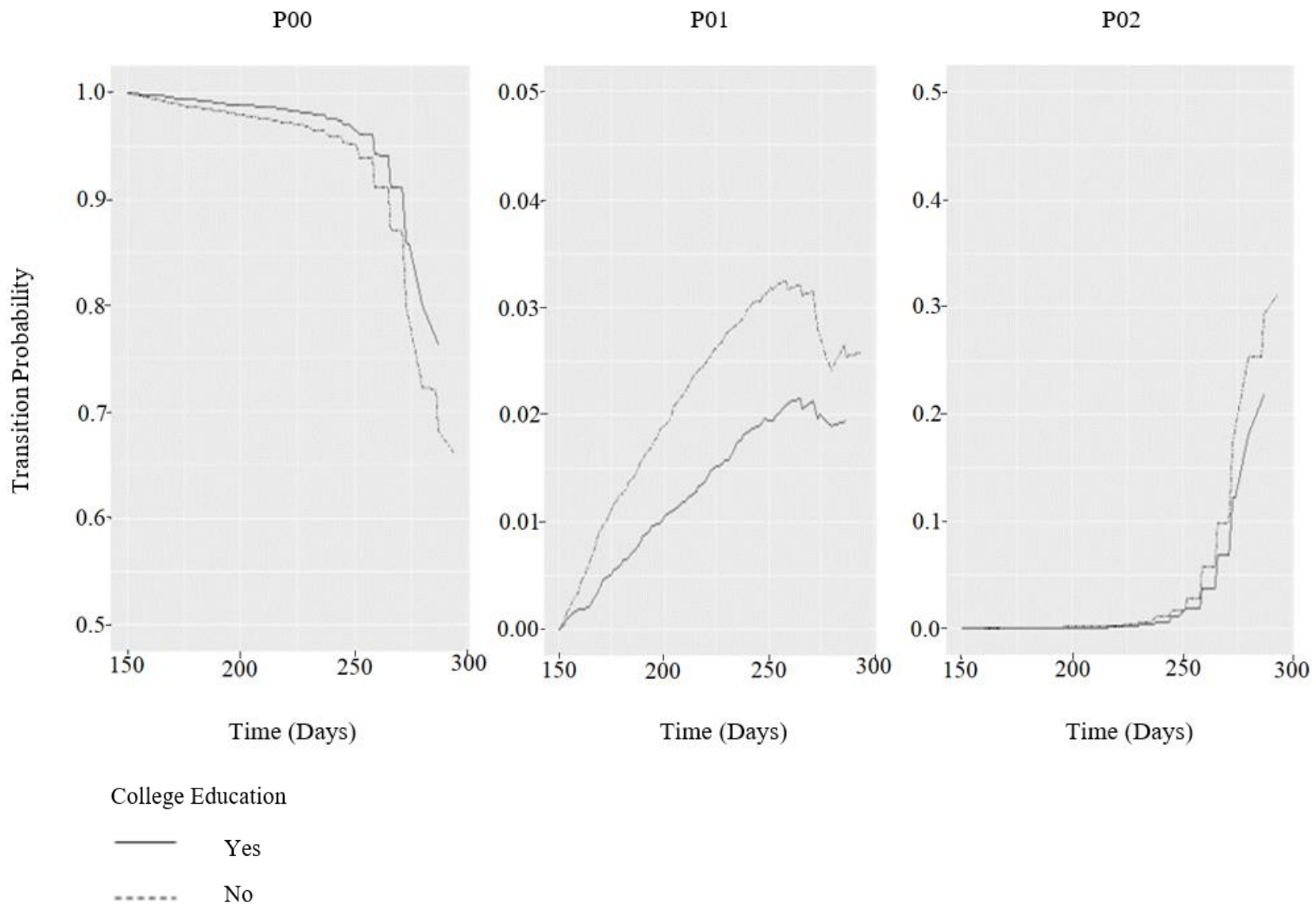


Figure 19. Transition probabilities P00, P01 and P02 plots for education groups in healthy start data set.

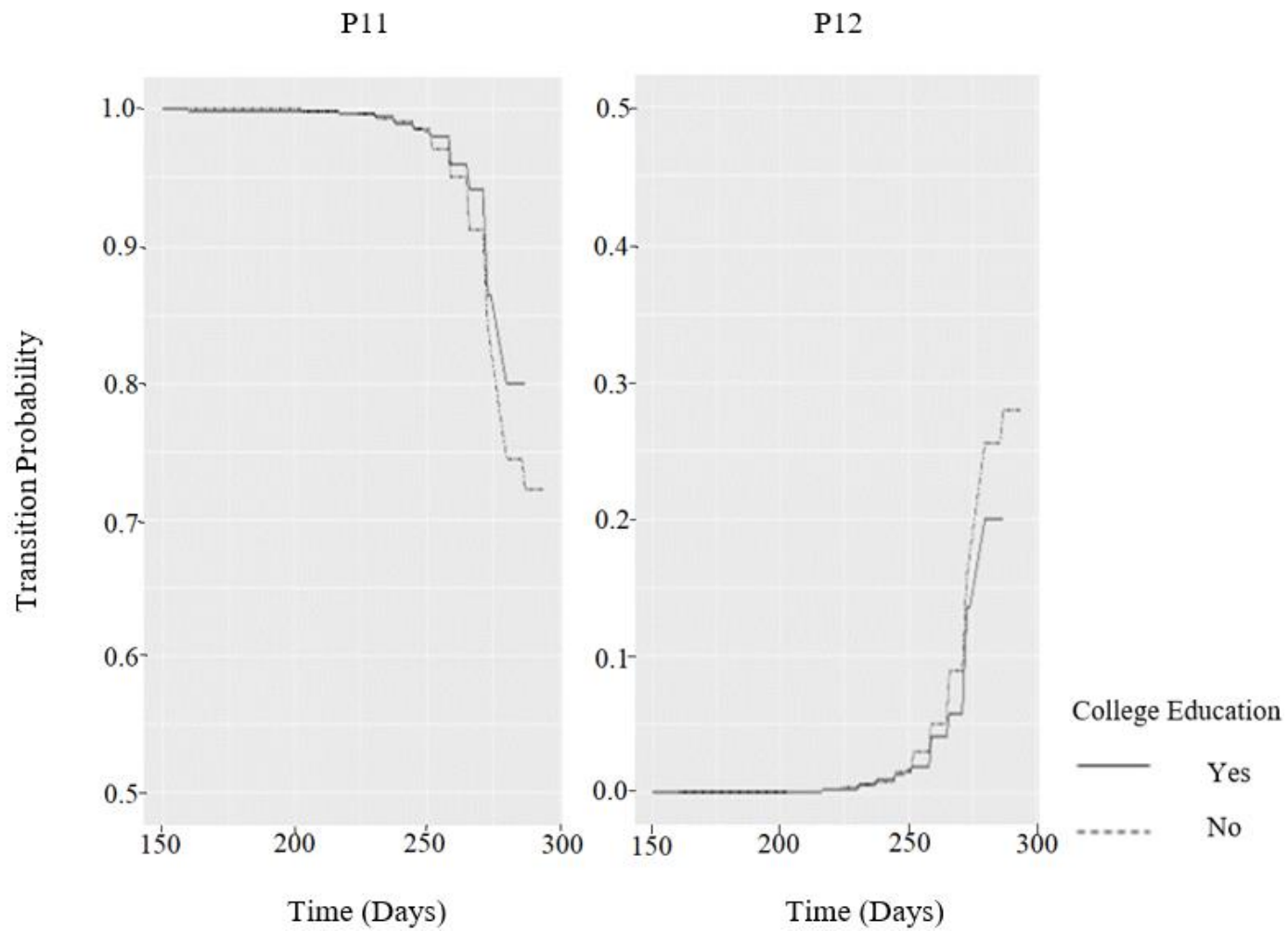


Figure 20. Transition probabilities P11 and P12 plots for education groups in healthy start data set.

4 Chapter III: Application of Aalen–Johansen Estimator Based on RSS Design to Colon

Cancer Dataset

4.1 Data Description

“ColonTP” data include colon adenocarcinoma patients who have received en bloc resection. This data set was first introduced by Laurie et al.⁹⁶ A complete study report was described by Moertel et al.⁹⁷ After enrollment was completed in October 1987, eligible patients were randomly assigned to control group, therapy with levamisole group, or therapy with levamisole plus fluorouracil.⁹⁷ Participants were followed to record cancer recurrence as well as survival up to 5 years. The data set includes 929 subjects and 15 variables as following:⁴³

time1 Disease free survival time (time to recurrence, death, or censoring, whichever occurs first)

event1 Disease free survival indicator (1=dead or relapsed, 0=alive disease free)

Stime Time to death or censoring.

event Death indicator (1=dead, 0=alive).

rx Treatment (Obs=observation, Lev=Levamisole, Lev+5-FU=Levamisole+fluorouracil).

sex 1=male.

age Age in years.

obstruct Obstruction of colon by tumour.

perfor Perforation of colon.

adhere Adherence to nearby organs.

nodes Number of lymph nodes with detectable cancer.

differ Differentiation of tumour (1=well, 2=moderate, 3=poor).

extent Extent of local spread (1=submucosa, 2=muscle, 3=seros, 4=contiguous structures).

surg Time from surgery to registration (0=short, 1=long).

nodes More than 4 positive lymph nodes.

Even though originally the data was designed for colon cancer survival analysis with Kaplan-Meier method, the data later became an classic real world example for illness death model cited by many literatures.^{43,60,98,99} In an illness death model, when subjects entered the study they were considered in State 0. Those who encountered cancer recurrence entered State 1. Death is State 2, which is an absorption state. Since the data set includes the following variables: “event1 (Disease free survival indicator)”, “event (Death indicator)”, “time1 (Disease free survival time” and “Stime (Time to death or censoring)”, it could be a perfect illustration for illness-death model.

4.2 Methods

4.2.1 An RSS modified Aalen Johansen Estimator for Colon Cancer Dataset

An illness-death model was used in this study to address possible transitions between states that include health (state 0), cancer recurrence (state 1) and death (state 2), which corresponds with a typical three states progressive illness-death model well. There are three possible transitions among them: $0 \rightarrow 1$, $0 \rightarrow 2$, $1 \rightarrow 2$. At initial time, all subjects are in state 0, and they are supposed to reach the final absorbing state 2 at future time point, along the process, they may experience or not an intermediate state (state 1). In this study the intermediate state represents cancer recurrence, the time spent in state 0 is referred to as healthy with no cancer recurrence. Aalen Johansen⁵³ estimator based on Markov assumptions under RSS sampling design was exploited to estimate the transition probabilities.

To perform RSS sampling method, a covariable correlated with the interested outcome variable needs to be identified for ranking purpose. Previous literatures indicate that age is correlated with colon cancer mortality rate.^{100,101} To investigate the relationship between survival time and all available variable, firstly, all censored variables were excluded. Because after excluding the censored subjects, “Stime” is the uncensored variable indicating survival time until death. There are 452 subjects remaining. A simple linear regression was carried out with “Stime” as dependent variable, other 11 variables as independent variables. The regression coefficients and p values are listed in Table 38. A simple random

sample of sample size 200 was selected from colonTP dataset. Transition probabilities P_{00}, P_{01}, P_{02} were calculated. Squared errors of distribution functions for the above estimators at 25%, 50%, 75% and 95% percentiles were calculated to show the performance. A ranked set sample of set number 5 and cycle number 40, which is equal to sample size 200 was selected from colonTP dataset. Transition probabilities P_{00}, P_{01}, P_{02} were calculated. Squared errors of distribution functions for the above estimators at 25%, 50%, 75% and 95% percentiles were calculated to show the performance. Table 39 presents the study results.

4.2.2 Statistical Analysis

The proposed Aalen–Johansen (AJ) estimator for transition probabilities based on the $RSS(k, m)$ sampling in illness death model $(\tilde{Z}, \tilde{T}, \delta_0, \delta_1, \beta)$ is applied to calculate the transition probability.

For transition probability from state 0 to state 0, we have

$$\hat{P}_{00RSS}^{AJ}(s, t) = \frac{1}{k} \sum_{r=1}^k \hat{P}_{00[r]}(s, t),$$

$$\hat{P}_{00[r]}(s, t) = \prod_{s < \tilde{Z}_l \leq t, l=1}^m \left(1 - \frac{\delta_{0l}}{R_0(\tilde{Z}_l)}\right)^{1_{\{s < Y_{[r]l}^* \leq t\}}} \quad (r = 1, \dots, k),$$

where $Y_{[r]1}^*, \dots, Y_{[r]m}^*$ are ordered values of the units of the r^{th} rank and $R_0(t) = \sum_{l=1}^m I(\tilde{Z}_l \geq t)$.

For transition probability from state 1 to state 1, we have

$$\hat{P}_{11RSS}^{AJ}(s, t) = \frac{1}{k} \sum_{r=1}^k \hat{P}_{11[r]}(s, t),$$

$$\hat{P}_{11[r]}(s, t) = \prod_{s < \tilde{T}_l \leq t, \beta_l=1, l=1}^m \left(1 - \frac{\delta_{1l}}{R_1(\tilde{T}_l)}\right)^{1_{\{s < Y_{[r]l}^* \leq t\}}} \quad (r = 1, \dots, k),$$

where $R_1(t) = \sum_{l=1}^m I(\tilde{Z}_l < t \leq \tilde{T}_l)$.

Then the modified transition probability from state 0 to state 1 is proposed as

$$\hat{P}_{01RSS}^{AJ}(s, t) = \frac{1}{k} \sum_{r=1}^k \hat{P}_{01[r]}(s, t),$$

$$\hat{P}_{01[r]}(s, t) = \sum_{l=1}^m \hat{P}_{00[r]}(s, \tilde{Z}_l^-) \hat{P}_{11[r]}(\tilde{Z}_l, t) I(s < \tilde{Z}_l \leq t) \frac{\beta_l}{R_0(\tilde{Z}_l)} \quad (r = 1, \dots, k),$$

Finally, it is obvious to propose the following transition probabilities from state 0 to state 2, and from state 1 to state 2. Since in Aalen–Johansen (AJ) transition probabilities, $P_{00} + P_{01} + P_{02} = 1$ and $P_{11} + P_{12} = 1$.

$$\begin{aligned}\hat{P}_{02RSS}^{AJ}(s, t) &= \frac{1}{k} \sum_{r=1}^k \hat{P}_{02[r]}(s, t), \\ \hat{P}_{02[r]}(s, t) &= 1 - \hat{P}_{00[r]}(s, t) - \hat{P}_{01[r]}(s, t) \quad (r = 1, \dots, k), \\ \hat{P}_{12RSS}^{AJ}(s, t) &= \frac{1}{k} \sum_{r=1}^k \hat{P}_{12[r]}(s, t), \\ \hat{P}_{12[r]}(s, t) &= 1 - \hat{P}_{11[r]}(s, t) \quad (r = 1, \dots, k)\end{aligned}$$

4.3 Results

It is clearly showed in Table 38 that age ($p = 0.0438$), obstruction of colon by tumor ($p = 0.0177$), differentiation of tumor ($p = 0.0049$), extent of local spread and more than 4 positive lymph nodes ($p = 0.0044$) have significant p values ($p < 0.05$). Previous literature research indicates that age is correlated with colon cancer mortality rate.^{100,101} Age is a continuous demographic variable, which is uncostly to measure. All above properties make age a qualified covariant for ranking in RSS sampling design.

In Table 39, among 12 estimators, except for P_{00} at 95 percentile and P_{01} at 50 percentile, 10 RSS estimators perform equal to or better than SRS competitors. For transition probability P_{00} at 25 percentiles, the squared error from SRS design is 476.19 times as high as the squared error from RSS design. For transition probability P_{02} at 25 percentiles, the squared error from SRS design is 232.43 times as high as the squared error from RSS design. Only the estimators for P_{00} at 95 percentiles ($1.74e-4$) and P_{01} at 50 percentiles ($5.49e-4$) have higher squared error for RSS design than for corresponding SRS design ($1.14e-4$, $1.15e-5$). Interestingly, the dominant superiority of RSS estimator over SRS estimator corresponds well with simulation results with transition probabilities P_{00} and P_{02} having more improved efficiency than transition probability P_{01} .

In Table 40, it is shown that when RSS sample size is 200 with a set number of 4, the estimators are very close to the real transition probabilities of full colonTP dataset except for transition probability P_{01} at

50 percentiles. For transition probability P_{00} , all four RSS estimators are close to the true values with differences within 5%. For transition probability P_{02} , three out of four RSS estimators are highly close to the corresponding true values with differences within 2%. The only exception is at 75%, however, the difference is still not large (6.36%) and the performance is still better than SRS counterpart. When transition probability is P_{01} , the RSS estimators are not as close to the true values as for other transition probabilities, except for 25% (99.67%).

4.4 Discussion

The application results represent simulation study well with excellent estimation of the real transition probabilities for P_{00} and P_{02} based on RSS sampling design. However, the efficiency of the RSS estimator is compromised for transition probability P_{01} . As we previously discussed, this is due to the fact that progressive illness-death model is a stochastic process, random variables T_{01} (time from state 0 to state 1), T_{12} (time from state 1 to state 2) and T_{02} (time from state 0 to state 2) are independent.⁶⁰ When ranking variable is correlated with T_{02} , the transition probabilities that are influenced are P_{00} and P_{02} , which justifies the Markov property of illness-death model.

Additionally, in this real-world application RSS modified AJ estimator still shows dominant superiority than its SRS counterpart. Even for transition probability P_{01} which has the worst performance for RSS design, 3 out of 4 RSS estimators show efficiency advantage than their SRS competitors. This dominant advantage of RSS modified AJ estimator over SRS design also correspond with the simulation results.

A limitation of this application study is that the RSS study we performed here did not design and collect data from a population, which is not a real RSS sample. The colon cancer dataset itself is a random sample from population from which we draw an RSS sample. Carrying out research like this underestimates the difficulty of conducting an RSS study directly from population. However, there is no reported clinical study based on RSS design. Hopefully the situation could be improved in future.

Table 38. Regression coefficients and p values of linear regression between survival time and various variables in colonTP dataset.

Variable name	Regression coefficient	p-value
Rx (Lev)	-77.051	0.24771
Rx (Lev+5FU)	-9.699	0.89251
sex	46.112	0.42484
age	-4.774	0.04377*
obstruct	-165.165	0.01773*
perfor	83.796	0.59073
adhere	-1.801	0.98120
nodes	-3.635	0.72026
differ	-157.241	0.00494*
extent	-202.991	0.00440*
surg	-10.905	0.86051
node4	-232.771	0.00818*

Table 39. Squared errors of distribution function estimators at some percentiles for $n = 200$ in colonTP dataset.

Time point	P_{00}				P_{01}				P_{02}			
	0.25 (728)	0.50 (1455)	0.75 (2183)	0.95 (2765)	0.25 (728)	0.50 (1455)	0.75 (2183)	0.95 (2765)	0.25 (728)	0.50 (1455)	0.75 (2183)	0.95 (2765)
RSS squared error	2.31e-7	4.2e-6	3.23e-4	1.74e-4	3.29e-7	5.49e-4	1.6e-4	9.09e-5	1.11e-6	6.47e-4	2.81e-5	1.32e-5
SRS squared error	1.1e-4	4.86e-4	6.75e-4	1.14e-4	3.11e-5	1.15e-5	1.61e-4	6.36e-4	2.58e-4	6.47e-4	1.77e-4	2.11e-4
Ratio of the two	476.19	115.71	2.09	0.66	94.53	0.02	1.01	7.00	232.43	1.00	6.30	15.98

Table 40. Percentage of RSS estimators of the real transition probabilities in colonTP dataset with a 200 sample size.

Transition probability	Time point (percentile)	RSS (%)
P ₀₀	25	99.92
	50	99.60
	75	96.11
	95	96.99
P ₀₁	25	99.67
	50	73.54
	75	118.45
	95	119.56
P ₀₂	25	100.47
	50	106.36
	75	101.13
	95	100.71

5 Limitations

The limitations of the study include not considering the cost of sampling and ranking, which is assumed to be minimum. However, in practice though the cost of sampling and ranking is much less than full measurement, but it cannot be totally ignored. There are studies focusing on comparing cost of RSS plus ranking expenses to the cost of SRS design and conclude that RSS design with optimum set number is still more efficient than SRS.¹⁰² Another limitation of the study is that properties of the proposed transition probabilities P_{11} and P_{12} are partially discussed. Since the ranking variable is T_{02} (time from state 0 to state 2), the simulation results show that the efficiency of the RSS modified transition probability P_{01} was not improved. But after ranking T_{12} (time from state 1 to state 2) transition probabilities P_{11} and P_{12} under RSS design show improvement under several simulation scenarios.

Future studies could consider ranking variable T_{01} (time from state 0 to state 1) to see if the performance of transition probability P_{01} will improve compared to its SRS competitor. Under the same design, if there is no improvement of transition probability P_{02} , then the Markov property of illness death model is justified again, since T_{01} (time from state 0 to state 1) and T_{02} (time from state 0 to state 2) are independent. Efforts could also be put on continuing studying transition probabilities P_{11} and P_{12} under RSS design for larger sample sizes and other censoring levels. We suspect that after ranking T_{12} , the efficiency improvement of transition probabilities P_{11} and P_{12} will be like Kaplan-Meier estimator, since they are the last two states of Markov process, and it is one direction with no branches. Finally, we would recommend the use of novel derivatives of ranked set sampling to extend the current study, such as Partially Rank-Ordered Set (PROS),³⁴ even order ranked set sampling (EORSS)¹⁰³ and quartile pair ranked set sampling (QPRSS).¹⁰⁴

6 Contributions

This is the first time that the AJ nonparametric estimator using RSS sampling design was proposed and compared with its SRS counterpart through simulation study. The effect of RSS design on two censoring indicators were investigated simultaneously, which is an extension of previous study regarding survival analysis with only one censoring event. The simulation study finds the dominant superiority of RSS modified AJ estimator for transition probabilities P_{00} and P_{02} in illness-death Markov model over its counterpart from SRS design for various sample sizes. The simulation study also indicates that when sample size is fixed, as set number increases the efficiency improvement for the estimator from RSS design compared with its SRS competitor becomes more significant, which proves the efficiency gain comes from ranking process.

It is the first time that illness-death model based on conventional AJ estimator was applied to investigate the effectiveness of the healthy start program in relation to the delivery of an SGA infant. Our results show that the risk of delivering an SGA infant was greatly reduced when a mother received healthy start services compared to mothers that did not receive healthy start services. Previously, Salihu et al.⁷⁷ observed the healthy start program to be effective in reducing preterm birth but ineffective in impacting SGA. To our knowledge, this is the only study that has shown the health start program in reducing SGA.

In this study, it is the first time RSS modified AJ estimator was applied to a real-world colon cancer dataset. The RSS estimator presents dominant advantage over its SRS counterparts in approximating transition probabilities P_{00} , P_{01} and P_{02} . The application results correspond well with the simulation outcome.

References

1. Olken F, Rotem D. Simple random sampling from relational databases. 1986.
2. Banerjee A, Chaudhury S. Statistics without tears: Populations and samples. *Industrial psychiatry journal*. 2010;19(1):60.
3. Martínez-Mesa J, González-Chica DA, Duquia RP, Bonamigo RR, Bastos JL. Sampling: how to select participants in my research study? *Anais brasileiros de dermatologia*. 2016;91:326-330.
4. Sampling Essentials: Practical Guidelines for Making Sampling Choices. In. Thousand Oaks, California: SAGE Publications, Inc.; 2012.
5. Singh S. Simple Random Sampling. In: *Advanced Sampling Theory with Applications: How Michael 'selected' Amy Volume I*. Dordrecht: Springer Netherlands; 2003:71-136.
6. Elfil M, Negida A. Sampling methods in clinical research; an educational review. *Emergency*. 2017;5(1).
7. Daniel J. Choosing the type of probability sampling. *Sampling essentials: Practical guidelines for making sampling choices*. 2012:125-175.
8. Acharya AS, Prakash A, Saxena P, Nigam A. Sampling: Why and how of it. *Indian Journal of Medical Specialties*. 2013;4(2):330-333.
9. McIntyre G. A method for unbiased selective sampling, using ranked sets. *Australian Journal of Agricultural Research*. 1952;3(4):385-390.
10. Patil G. ranked set sampling. *Environmental and Ecological Statistics*. 1995;2(4):271-285.
11. Halls LK, Dell TR. Trial of Ranked-Set Sampling for Forage Yields. *Forest Science*. 1966;12(1):22-26.
12. McCrum-Gardner E. Sample size and power calculations made simple. *International Journal of Therapy and Rehabilitation*. 2010;17(1):10-14.
13. Hazra A. Using the confidence interval confidently. *J Thorac Dis*. 2017;9(10):4125-4130.
14. Wolfe DA. Ranked Set Sampling: Its Relevance and Impact on Statistical Inference. *ISRN Probability and Statistics*. 2012;2012:568385.
15. Ozturk O. Sampling from partially rank-ordered sets. *Environmental and Ecological statistics*. 2011;18(4):757-779.
16. Hatefi A, Jozani MJ. Information content of partially rank-ordered set samples. *AStA Advances in Statistical Analysis*. 2017;101(2):117-149.
17. Ozturk O. Quantile inference based on partially rank-ordered set samples. *Journal of Statistical Planning and Inference*. 2012;142(7):2116-2127.
18. Dell TR, Clutter JL. Ranked Set Sampling Theory with Order Statistics Background. *Biometrics*. 1972;28(2):545-555.
19. Ozturk O, Bilgin OC, Wolfe DA. Estimation of population mean and variance in flock management: a ranked set sampling approach in a finite population setting. *Journal of Statistical Computation and Simulation*. 2005;75(11):905-919.
20. Aragon E, Gore S, Patil G. Ranked set sampling: a bibliography. *Environmental and Ecological Statistics*. 1999;6:91-98.
21. Stokes SL, Sager TW. Characterization of a ranked-set sample with application to estimating distribution functions. *Journal of the American Statistical Association*. 1988;83(402):374-381.
22. Stokes SL. Estimation of Variance Using Judgment Ordered Ranked Set Samples. *Biometrics*. 1980;36(1):35-42.
23. Stokes SL. Inferences on the correlation coefficient in bivariate normal populations from ranked set samples. *Journal of the American Statistical Association*. 1980;75(372):989-995.

24. Martin WL, Sharik TL, Oderwald RG, Smith DW. Evaluation of ranked set sampling for estimating shrub phytomass in Appalachian oak forests. 1980.
25. Cobby J, Ridout M, Bassett P, Large R. An investigation into the use of ranked set sampling on grass and grass-clover swards. *Grass and Forage Science*. 1985;40(3):257-263.
26. Nelson LE, Switzer GL, Lockaby BG. Nutrition of *Populus deltoides* plantations during maximum production. *Forest Ecology and Management*. 1987;20(1):25-41.
27. Mode NA, Conquest LL, Marker DA. Ranked set sampling for ecological research: accounting for the total costs of sampling. *Environmetrics*. 1999;10(2):179-194.
28. Mode NA, Conquest LL, Marker DA. Incorporating prior knowledge in environmental sampling: ranked set sampling and other double sampling procedures. *Environmetrics*. 2002;13(5-6):513-521.
29. Murray R, Ridout M, Cross J. The use of ranked set sampling in spray deposit assessment. *Aspects of Applied Biology*. 2000;57:141-146.
30. Al-Saleh MF, Al-Shrafat K. Estimation of average milk yield using ranked set sampling. *Environmetrics*. 2001;12(4):395-399.
31. Kvam PH. Ranked set sampling based on binary water quality data with covariates. *Journal of Agricultural, Biological, and Environmental Statistics*. 2003;8(3):271.
32. Yu PLH, Tam CYC. Ranked set sampling in the presence of censored data. *Environmetrics*. 2002;13(4):379-396.
33. Strzalkowska-Kominiak E, Mahdizadeh M. On the Kaplan–Meier estimator based on ranked set samples. *Journal of Statistical Computation and Simulation*. 2014;84(12):2577-2591.
34. Nematollahi S, Nazari S, Shayan Z, Ayatollahi SMT, Amanati A. Improved Kaplan-Meier Estimator in Survival Analysis Based on Partially Rank-Ordered Set Samples. *Computational and Mathematical Methods in Medicine*. 2020;2020:7827434.
35. Zamanzade E, Parvardeh A, Asadi M. Estimation of mean residual life based on ranked set sampling. *Computational Statistics & Data Analysis*. 2019;135:35-55.
36. Samawi HM, Helu A, Rochani H, Yin J, Yu L, Vogel R. Reducing sample size needed for accelerated failure time model using more efficient sampling methods. *Journal of Statistical Theory and Practice*. 2018;12(3):530-541.
37. Ata Tutkun N, Koyuncu N, Karabey U. Discrete-time survival analysis under ranked set sampling: an application to Turkish motor insurance data. *Journal of Statistical Computation and Simulation*. 2019;89(4):660-667.
38. Sonnenberg FA, Beck JR. Markov Models in Medical Decision Making: A Practical Guide. *Medical Decision Making*. 1993;13(4):322-338.
39. SCHEIKE TH, ZHANG M-J. Direct Modelling of Regression Effects for Transition Probabilities in Multistate Models. *Scandinavian Journal of Statistics*. 2007;34(1):17-32.
40. Andersen P, Borgan O, Gill R, Keiding N. *Statistical Models Based on Counting Processes*. New York: Springer-Verlag; 1993.
41. Hougaard P. *Analysis of Multivariate Survival Data*. Springer-Verlag; 2000.
42. Touraine C, Helmer C, Joly P. Predictions in an illness-death model. 2016;25(4):1452-1470.
43. Balboa V, de Uña-Álvarez J. Estimation of Transition Probabilities for the Illness-Death Model: Package TP.idm. 2018. 2018;83(10):19.
44. Kodell RL, Nelson CJ. An Illness-Death Model for the Study of the Carcinogenic Process Using Survival/Sacrifice Data. *Biometrics*. 1980;36(2):267-277.
45. Commenges D, Joly P, Letenneur L, Dartigues J. Incidence and mortality of Alzheimer's disease or dementia using an illness-death model. 2004;23(2):199-210.
46. Harezlak J, Gao S, Hui SL. An illness–death stochastic model in the analysis of longitudinal dementia data. 2003;22(9):1465-1475.
47. Frydman H, Szarek M. Estimation of overall survival in an ‘illness–death’ model with application to the vertical transmission of HIV-1. 2010;29(19):2045-2054.

48. Andersen PK, Esbjerg S, Sørensen TIA. Multi-state models for bleeding episodes and mortality in liver cirrhosis. 2000;19(4):587-599.
49. Putter H, Fiocco M, Geskus RB. Tutorial in biostatistics: competing risks and multi-state models. 2007;26(11):2389-2430.
50. Frydman H. Semiparametric Estimation in a Three-State Duration-Dependent Markov Model from Interval-Censored Observations with Application to AIDS Data. *Biometrics*. 1995;51(2):502-511.
51. Aalen OJTAoS. Nonparametric inference for a family of counting processes. 1978:701-726.
52. Aalen OO. *Statistical inference for a family of counting processes*. University of California, Berkeley; 1975.
53. Aalen OO, Johansen S. An Empirical Transition Matrix for Non-Homogeneous Markov Chains Based on Censored Observations. *Scandinavian Journal of Statistics*. 1978;5(3):141-150.
54. Meira-Machado L, de Uña-Álvarez J, Cadarso-Suárez C, Andersen PK. Multi-state models for the analysis of time-to-event data. 2009;18(2):195-222.
55. Datta S, Satten GAJB. Estimation of integrated transition hazards and stage occupation probabilities for non-Markov systems under dependent censoring. 2002;58(4):792-802.
56. Meira-Machado L, De Una-Alvarez J, Cadarso-Suarez CJLDA. Nonparametric estimation of transition probabilities in a non-Markov illness–death model. 2006;12(3):325-344.
57. Meira-Machado L, Sestelo MJB. Estimation in the progressive illness-death model: A nonexhaustive review. 2019;61(2):245-263.
58. Machado LM. Presmoothed Landmark estimators of the transition probabilities. 2016.
59. Chen Z, Bai Z, Sinha BK. *Ranked Set Sampling: Theory and Applications*. New York: Springer-Verlag; 2004.
60. de Uña-Álvarez J, Meira-Machado L. Nonparametric estimation of transition probabilities in the non-Markov illness-death model: A comparative study. *Biometrics*. 2015;71(2):364-375.
61. Sevinc B, Cetintav B, Esemem M, Gurler S. RSSampling: A Pioneering Package for Ranked Set Sampling. *R J*. 2019;11(1):401.
62. ALEXANDER GR, KOGAN M, MARTIN J, PAPIERNIK E. What Are the Fetal Growth Patterns of Singletons, Twins, and Triplets in the United States? *Clinical Obstetrics and Gynecology*. 1998;41(1):115-125.
63. McCowan L, Horgan RP. Risk factors for small for gestational age infants. *Best Practice & Research Clinical Obstetrics & Gynaecology*. 2009;23(6):779-793.
64. Saenger P, Czernichow P, Hughes I, Reiter EO. Small for Gestational Age: Short Stature and Beyond. *Endocrine Reviews*. 2007;28(2):219-251.
65. Battaglia FC, Lubchenco LO. A practical classification of newborn infants by weight and gestational age. *The Journal of Pediatrics*. 1967;71(2):159-163.
66. Verkauskienė R, Wikland KA, Niklasson A. Variation in size at birth in infants born small for gestational age in Lithuania. *Acta Paediatrica*. 2002;91(3):329-334.
67. Hokken-Koelega ACS, Chernausek SD, Kiess W. *Small for Gestational Age : Causes and Consequences*. Basel: Karger; 2009.
68. Kramer MS. Determinants of low birth weight: methodological assessment and meta-analysis. *Bull World Health Organ*. 1987;65(5):663-737.
69. Ash S, Fisher C, Truswell A, Allen J, Irwig L. Maternal weight gain, smoking and other factors in pregnancy as predictors of infant birth-weight in Sydney women. *Aust NZ J Obstet Gynaecol*. 1989(29):212-219.
70. Thompson J, Clark P, Robinson E, et al. Risk factors for small-for-gestational-age babies: The Auckland Birthweight Collaborative Study. *Journal of Paediatrics and Child Health*. 2001;37(4):369-375.
71. Raum E, Arabin B, Schlaud M, Walter U, Schwartz FW. The impact of maternal education on intrauterine growth: a comparison of former West and East Germany. 2001;30(1):81-87.

72. Raatikainen K, Heiskanen N, Heinonen S. Marriage still protects pregnancy. *BJOG: An International Journal of Obstetrics & Gynaecology*. 2005;112(10):1411-1416.
73. Khashan AS, Baker PN, Kenny LC. Preterm birth and reduced birthweight in first and second teenage pregnancies: a register-based cohort study. 2010;10(1):36.
74. Turner RJ, Grindstaff CF, Phillips N. Social Support and Outcome in Teenage Pregnancy. *Journal of Health and Social Behavior*. 1990;31(1):43-57.
75. US Department of Health and Human Services. The Secretary's Advisory Committee on National Health Promotion and Disease Prevention Objectives for 2020. Phase I report: recommendations for the framework and format of Healthy People 2020. Section IV. Advisory Committee findings and recommendations. In:2010.
76. US Department of Health and Human Services. Healthy People 2010. Understanding and improving health. In. Vol 1 and 2. 2nd ed. Washington, DC: US GPO2000.
77. Salihu HM, Mbah AK, Jeffers D, Alio AP, Berry L. Healthy start program and feto-infant morbidity outcomes: evaluation of program effectiveness. *Matern Child Health J*. 2009;13(1):56-65.
78. Jevitt C, Zapata L, Harrington M, Berry E. Screening for Perinatal Depression With Limited Psychiatric Resources. *Journal of the American Psychiatric Nurses Association*. 2005;11(6):359-363.
79. Salihu H, August E, Mbah A, et al. Effectiveness of a Federal Healthy Start Program in Reducing the Impact of Particulate Air Pollutants on Feto-Infant Morbidity Outcomes. *Maternal & Child Health Journal*. 2012;16(8):1602-1611.
80. Salihu HM, August EM, Jeffers DF, Mbah AK, Alio AP, Berry E. Effectiveness of a Federal Healthy Start Program in Reducing Primary and Repeat Teen Pregnancies: Our Experience over the Decade. 2011;24(3):153-160.
81. Salihu HM, August EM, Mbah AK, et al. The Impact of Birth Spacing on Subsequent Feto-Infant Outcomes among Community Enrollees of a Federal Healthy Start Project. 2012;37(1):137-142.
82. Badura M, Johnson K, Hench K, Reyes M. Healthy Start: Lessons Learned on Interconception Care. *Women's Health Issues*. 2008;18(6, Supplement):S61-S66.
83. Malloy MH. Size for gestational age at birth: impact on risk for sudden infant death and other causes of death, USA 2002. *Arch Dis Child Fetal Neonatal Ed*. 2007;92(6):F473-478.
84. Mitchell E, Thompson J, Robinson E, et al. Smoking, nicotine and tar and risk of small for gestational age babies. *Acta Paediatrica*. 2002;91(3):323-328.
85. Cogswell ME, Yip R. The influence of fetal and maternal factors on the distribution of birthweight. *Semin Perinatol*. 1995;19(3):222-240.
86. Ford JH. Preconception risk factors and SGA babies: Papilloma virus, omega 3 and fat soluble vitamin deficiencies. *Early Human Development*. 2011;87(12):785-789.
87. Luke B, Williams C, Minogue J, Keith L. The changing pattern of infant mortality in the US: the role of prenatal factors and their obstetrical implications. *Int J Gynaecol Obstet*. 1993;40(3):199-212.
88. Steer P. Small for Gestational Age: Causes and Consequences. *J Anat*. 2009;215(2):224-224.
89. Kijima M. Continuous-time Markov chains. In: *Markov Processes for Stochastic Modeling*. Boston, MA: Springer US; 1997:167-241.
90. Haigh J. *Probability models*. [electronic resource]. 2nd ed. ed: Springer; 2013.
91. Thomas DR, Grunkemeier GL. Confidence Interval Estimation of Survival Probabilities for Censored Data. *Journal of the American Statistical Association*. 1975;70(352):865-871.
92. Kalbfleisch J, Prentice R. *Statistical Analysis of Failure Time Data*. 2 ed: John Wiley & Sons; 2002.
93. Taffel S, Heuser RL, Johnson DP. *A method of imputing length of gestation on birth certificates*. Hyattsville, Md. : Washington, D.C. : U.S. Dept. of Health and Human Services, Public Health Service, Office of Health Research, Statistics, and Technology, National Center for Health Statistics ; For sale by the Supt. of Docs., U.S. G.P.O., 1982; 1982.

94. Bartolomeo N, Trerotoli P, Moretti A, Serio G. A Markov model to evaluate hospital readmission. *BMC Med Res Methodol.* 2008;8:23.
95. Boyd MA, Lau S. An Introduction to Markov Modeling: Concepts and Uses. In. NASA Technical Reports Server1998.
96. Laurie JA, Moertel CG, Fleming TR, et al. Surgical adjuvant therapy of large-bowel carcinoma: an evaluation of levamisole and the combination of levamisole and fluorouracil. The North Central Cancer Treatment Group and the Mayo Clinic. *Journal of Clinical Oncology.* 1989;7(10):1447-1456.
97. Moertel CG, Fleming TR, Macdonald JS, et al. Levamisole and Fluorouracil for Adjuvant Therapy of Resected Colon Carcinoma. *New England Journal of Medicine.* 1990;322(6):352-358.
98. Balboa-Barreiro V, de Una-Alvarez J, Meira-Machado L, Balboa-Barreiro MV. Package ‘TP. idm’. 2018.
99. Meira-Machado L, de Uña-Álvarez J, Datta S. Nonparametric estimation of conditional transition probabilities in a non-Markov illness-death model. *Computational Statistics.* 2015;30(2):377-397.
100. van Leeuwen BL, Pählman L, Gunnarsson U, Sjövall A, Martling AJC. The effect of age and gender on outcome after treatment for colon carcinoma: A population-based study in the Uppsala and Stockholm region. 2008;67(3):229-236.
101. Aquina CT, Mohile SG, Tejani MA, et al. The impact of age on complications, survival, and cause of death following colon cancer surgery. *British Journal of Cancer.* 2017;116(3):389-397.
102. Nahhas RW, Wolfe DA, Chen H. Ranked Set Sampling: Cost and Optimal Set Size. 2002;58(4):964-971.
103. Noor-ul-Amin M, Tayyab M, Hanif MJPUJoM. Mean estimation using even order ranked set sampling. 2019;51(1):91-99.
104. Tayyab M, Noor-ul-Amin M, Hanif MJPotnaos, India section A: physical sciences. Quartile pair ranked set sampling: development and estimation. 2021;91(1):111-116.