

1 **INTEGRATIVE SURVIVAL ANALYSIS WITH UNCERTAIN**
2 **EVENT TIMES IN APPLICATION TO A SUICIDE RISK**
3 **STUDY**

4 BY WENJIE WANG[†] ROBERT ASELTINE[‡] KUN CHEN^{*,†,‡} AND JUN
5 YAN^{†,‡}

6 *University of Connecticut[†] and University of Connecticut Health Center[‡]*

7 The concept of integrating data from disparate sources to ac-
8 celerate scientific discovery has generated tremendous excitement in
9 many fields. The potential benefits from data integration, however,
10 may be compromised by the uncertainty due to incomplete/imperfect
11 record linkage. Motivated by a suicide risk study, we propose an ap-
12 proach for analyzing survival data with uncertain event times arising
13 from data integration. Specifically, in our problem deaths identified
14 from the hospital discharge records together with reported suicidal
15 deaths determined by the Office of Medical Examiner may still not
16 include all the death events of patients, and the missing deaths can
17 be recovered from a complete database of death records. Since the
18 hospital discharge data can only be linked to the death record data
19 by matching basic patient characteristics, a patient with a censored
20 death time from the first dataset could be linked to multiple po-
21 tential event records in the second dataset. We develop an integra-
22 tive Cox proportional hazards regression, in which the uncertainty in
23 the matched event times is modeled probabilistically. The estimation
24 procedure combines the ideas of profile likelihood and the expecta-
25 tional conditional maximization algorithm (ECM). Simulation studies
26 demonstrate that under realistic settings of imperfect data linkage,
27 the proposed method outperforms several competing approaches in-
28 cluding multiple imputation. A marginal screening analysis using the
29 proposed integrative Cox model is performed to identify risk fac-
30 tors associated with death following suicide-related hospitalization
31 in Connecticut. The identified diagnostics codes are consistent with
32 existing literature and provide several new insights on suicide risk
33 prediction and prevention.

34 **1. Introduction.** In many fields of science, engineering, and medicine,
35 combining multiple datasets from disparate sources has made it possible to
36 tackle important problems at an accelerated rate through integrative statis-
37 tical learning. These datasets cover overlapped or interrelated measurements
38 from individuals. In an ideal situation, the multi-source data should pertain
39 to the same set of fully identified individuals. For example, in a cancer study,
40 multi-platform genetic data such as mRNA gene expression, DNA methy-
41 lation, and copy number variation are available from each patient (Zhao

*Corresponding author; kun.chen@uconn.edu

Keywords and phrases: Cox model, Data linkage, ECM algorithm, Integrative learning, Suicide prevention

1 [et al., 2015](#)); an integrative analysis then ensures a comprehensive coverage
2 of genetic perspectives to understand the disease mechanism. In practice,
3 however, more than often, a unique identifier is not provided or does not
4 even exist to link multi-source or multi-platform datasets. This gives rise
5 to the so-called “data/record linkage” problem, i.e., matching records from
6 different sources that belong to the same person or entity based on available
7 characteristics of the entity (e.g., [Winglee, Valliant and Scheuren, 2005](#)); see
8 [Harron, Goldstein and Dibben \(2015\)](#) for a recent review. Matching errors
9 are bound to occur ([Bohensky et al., 2010](#)), and the potential benefits from
10 data integration may be compromised. Therefore, in statistical analysis with
11 integrated data, it is important to take into account the uncertainty due to
12 imperfect linkage.

13 Our research was motivated by the survival analysis of youth and young
14 adult patients in the State of Connecticut who were at elevated risk of sui-
15 cide because of having been hospitalized for suicide attempt or intentional
16 self-injury. Data from diagnosis were available from the Connecticut Hospi-
17 tal Inpatient Discharge Data (HIDD). Deaths by suicide were determined
18 from the Office of the Connecticut Medical Examiner (OCME). It has been
19 revealed, however, that suicidal death is often underreported in key Western
20 countries ([Pritchard and Hansen, 2015](#); [Tøllefsen et al., 2016](#)). The death
21 records identified from the OCME for this group are incomplete because,
22 first suicide deaths may be underreported, and second they do not include
23 deaths due to other causes. Hence, some patients with censored suicide times
24 might have died. While the missing deaths may possibly be recovered from a
25 complete mortality database of the state, the HIDD data can only be linked
26 to the death records by matching basic patient characteristics such as date
27 of birth, gender, race, and residential zip code, because there is no unique
28 identifier to join the two datasets even before the data were de-identified
29 in order to protect patient privacy. Consequently, in the integrated data, a
30 censored death time before matching could be linked to multiple possible
31 death times in the mortality data; see details in Section 2.

32 Figure 1 illustrates the data matching patterns in a general integrated
33 survival analysis setup similar to that in our suicide risk study. In dataset I,
34 a positive number of subjects’ event times are observed and known to be
35 accurate (Case 1). For those subjects whose event times are censored in
36 dataset I, their event times might be captured in dataset II. After the link-
37 age process with partial identifiers, the event time of any subject who does
38 not find a match in dataset II is still censored (Case 2). As such, Case 1 and
39 Case 2 consist of non-censored and censored subjects, respectively, in a stan-
40 dard right censored data setting. Challenges are brought by those subjects

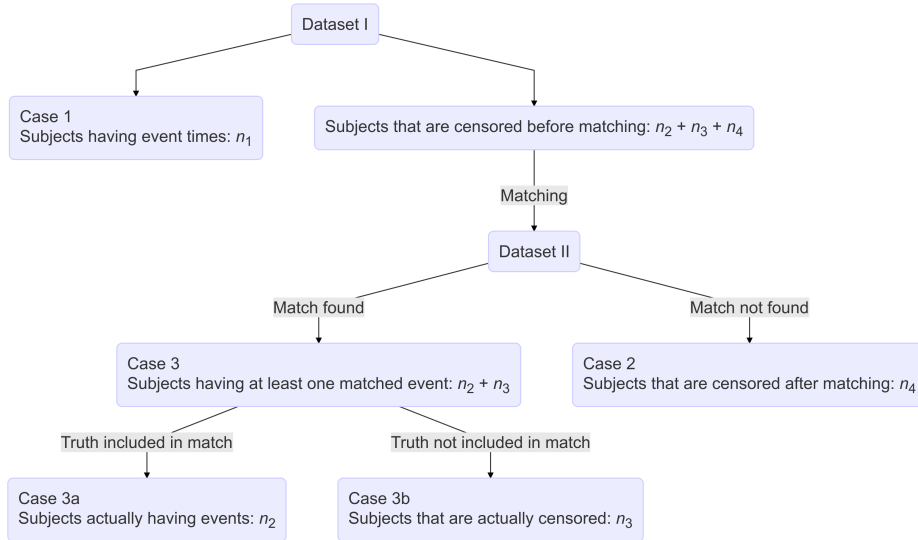


FIG 1. Illustration of the data matching patterns for studies with event time outcomes.

1 with one or more matches (Case 3): we are not sure which one, if any, of
 2 the matched event times is the truth. The subjects in Case 3 can be further
 3 classified into two types: Case 3a contains subjects whose true event time
 4 is included in the matched records, and Case 3b contains subjects whose true
 5 event time is not included in the matched records, and, hence, is actually
 6 censored. This classification is unknown and has to be inferred from the
 7 data. The task can be regarded as a missing data problem, in which the
 8 indicators of whether each matched record is true are missing.

9 Some efforts have been made to similar problems of mis-measured out-
 10 comes or uncertain endpoints. [Snapinn \(1998\)](#) proposed a modification of
 11 the Cox proportional hazard model ([Cox, 1972](#)) for nonfatal uncertain end-
 12 points by assigning weights that represent the likelihood of each potential
 13 endpoint being true. The determination of the weights, however, requires an
 14 additional diagnostic score and depends on a subjective estimation of the
 15 relative frequency of true endpoints to false endpoints suggested by the end-
 16 point committee or experts in the therapeutic area. [Richardson and Hughes](#)
 17 [\(2000\)](#) proposed an estimation procedure for the product limit estimate of
 18 survival function with no covariate based on the expectation maximization
 19 (EM) algorithm ([Dempster, Laird and Rubin, 1977](#)) when a binary diag-
 20 nosis outcome was measured with uncertainty. The method was designed
 21 for discrete-time contexts where the time points of outcome testing were

1 predetermined. Meier, Richardson and Hughes (2003) extended the discrete
2 proportional hazard model (Kalbfleisch and Prentice, 2002) to mis-measured
3 outcomes under a setting similar to Richardson and Hughes (2000) but al-
4 lowed covariate effects. In a more general setting, regression methods have
5 been developed for linked data where the response and covariates come from
6 two databases (e.g., Hof and Zwiderman, 2012, 2015; Tancredi and Liseo,
7 2015). None of the existing works was designed to handle the data integra-
8 tion problem in a survival analysis like ours.

9 We propose an integrative Cox proportional hazard model for data with
10 uncertain event time points. The uncertainty in the integrated survival data
11 is modeled probabilistically, where the probabilities depend only on the rel-
12 ative hazards from the Cox model itself. The model reduces to the regular
13 Cox model when there is no uncertain record. In contrast to the method of
14 Snapinn (1998), our method does not require any extra diagnostic variable or
15 prior knowledge on the initial probabilities indicating the true outcomes. The
16 estimation procedure combines the ideas of profile likelihood and the expect-
17 ation conditional maximization (ECM) algorithm. The proposed method is
18 shown to outperform naive approaches in simulation studies under realistic
19 settings similar to the real data example. We apply the proposed approach to
20 identifying risk factors associated with patient survival after suicide-related
21 hospitalization, using data obtained by integrating the HIDD/OCME data
22 and the mortality record data of the period 2005–2012 in Connecticut. The
23 identified diagnostic codes are mostly consistent with existing results and
24 provide several new insights on suicide risk prediction and prevention.

25 The rest of this paper is organized as follows. The settings for integrated
26 survival data for the Connecticut suicide risk analysis and the associated
27 challenges are presented in Section 2. In Section 3, we present the inte-
28 grative Cox regression modeling framework. The estimation procedure is
29 developed in Section 4. The simulation studies are presented in Section 5.
30 A marginal screening analysis using the proposed integrative model for the
31 Connecticut suicide risk study is reported in Section 6. Section 7 concludes
32 with a discussion. Implementation of the proposed methods is available in
33 a package named `intsurv` for R (R Development Core Team, 2017), which
34 can be accessed at <https://github.com/wenjie2wang/intsurv>.

35 **2. Integrated Survival Data of a Patient Group with Elevated**
36 **Suicide Risk.** Suicide is a serious public health problem in the US. Death
37 by suicide is increasing among all age groups in the US, with a 24% increase
38 in suicide rates observed from 1999 to 2014. There is a strong tendency for
39 suicide attempters to make additional attempts after the initial suicide at-

1 tempt (Suominen et al., 2004), and suicide attempt is a strong predictor of
2 suicidal death (Bostwick et al., 2015). Understanding factors associated with
3 suicide for patients hospitalized due to suicide attempt is critical to a better
4 allocation of selected prevention efforts among those at elevated risk. An im-
5 mediate challenge in statistical modeling is that attributing death to suicide
6 is not easy as suicidal death is often under-reported. For example, Pritchard
7 and Hansen (2015) showed that undetermined and accidental death was a
8 main source of the under-reported-suicides across different countries includ-
9 ing the US; Tøllefsen et al. (2016) reported that from re-evaluations of 1800
10 deaths in Scandinavia, 9% of the natural deaths and accidents were reclassi-
11 fied as suicides in the Norwegian data, and 21% of the undetermined deaths
12 were reclassified as suicides in the Swedish data.

13 We focused on patients of age 15–30 with high suicide risk in Connecticut.
14 This group of patients consisted of those who were admitted to a hospital in
15 Connecticut due to suicide attempt or self-inflicted injury, survived, and were
16 discharged, during fiscal years 2005–2012. The entry time of each patient into
17 the study is the time of last such discharge. The event time is the time to
18 death from all causes, including suicide, since the entry time. The cutoff date
19 of the HIDD is September 30, the end of fiscal year of 2012, which means
20 that the patients were followed up until this time. The OCME provided
21 data on suicide deaths of this period, which included a field for reporting
22 source that allowed accurate identification of the corresponding patients in
23 HIDD. Since the HIDD and OCME data only captured reported suicide
24 deaths, we acquired the complete mortality data of the same period from
25 the Connecticut Department of Public Health, aiming to recover the missing
26 deaths through record linkage using basic patient characteristics. The HIDD
27 and OCME data lead to Dataset I while the mortality data is Dataset II in
28 Figure 1. We stress that here we set the terminal event as death from all
29 causes rather than only due to suicide. This is mainly because the cause of
30 death is not available in the mortality data so that it can not be recovered
31 from data integration. On the other hand, without data integration, ignoring
32 unreported suicidal deaths and deaths due to other causes would jeopardize
33 the validity of statistical results. Because suicide is a major cause of death
34 among young suicide attempters, death due to all causes stands as a valid
35 terminal event to study in our problem.

36 A total of 7,304 patients were followed up until September 30, 2012.
37 Among them, 4,981 were white (2,775 female and 2,206 male) and 2,323
38 were non-white (1,304 female and 1,019 male). Before matching, Case 1
39 consisted of 133 patients with confirmed suicide death from the OCME, a
40 censoring rate of 98.2%. For the 7,171 patients with censored event times,

1 we made record linkage with the Connecticut state mortality database by
 2 date of birth, gender, and race. Since the death time had to happen after
 3 the discharge, we excluded any matched event before the discharge date of
 4 each patient during the matching process. After matching, Case 2 consisted
 5 of 6,546 patients with no matched record, while Case 3 consisted of 625 pa-
 6 tients with at least one matched records. In Case 3, 584 patients had one
 7 match, 39 patients had two matches, and two patients had four matches,
 8 it was possible for each patient to be still alive on September 30, 2012, in
 9 which case, the true death time is censored.

10 The HIDD data contained a large number of records on the characteris-
 11 tics of patients and their previous hospital admissions. The research inter-
 12 est was to identify important diagnostic categories associated with patient
 13 death. The diagnostics were recorded as ICD-9 diagnosis codes, or more
 14 formally ICD-9-CM (International Classification of Diseases, 9th Revision,
 15 Clinical Modification). We grouped the ICD-9 codes by their three leading
 16 characters that define the major diagnosis categories. Suicide attempts were
 17 identified by both ICD-9 external cause of injury codes and other ICD-9
 18 code combinations indicative of suicidal behavior (Patrick et al., 2010; Chen
 19 and Aseltine, 2017). Other ICD-9 codes during the inpatient hospitalization
 20 fell into 167 major diagnosis categories, which led to 167 indicator variables.
 21 Not all 167 indicators, however, can be used as covariates. Among them,
 22 51 ICD-9 indicators had quasi-complete separation (Albert and Anderson,
 23 1984) in our data; that is, there was no death event among those whose
 24 diagnosis included any of these ICD-9 categories. Although they could be
 25 potentially useful in predicting survival and thus merit further investiga-
 26 tion, they cannot be considered as covariates in a Cox regression framework
 27 adopted in this work, since their coefficient estimates would tend to be neg-
 28 ative infinite. To focus on the main idea, we further filtered out another 58
 29 ICD-9 indicators by restricting every cell of the cross table of the diagnosis
 30 indicator and event indicator to be at least three. The remaining 58 ICD-9
 31 codes were used in a marginal screening analysis; see Section 6.

32 **3. Integrative Cox Model.** Consider a random sample of n subjects
 33 who fall into the three cases as illustrated in Figure 1. Let I_1 , I_2 , and I_3
 34 be the indices of the subjects in Case 1, 2, and 3, respectively. For subject
 35 $j \in I_1$, we observe the event time V_j . For subject $j \in I_2$, we observe the
 36 censoring time C_j . For subject $j \in I_3$, the true event time V_j has $s_j \geq 2$
 37 possibilities, $0 < V_{j,1} < \dots < V_{j,s_j-1} < V_{j,s_j}$, but we only observe $0 < V_{j,1} <$
 38 $\dots < V_{j,s_j-1} < C_j$, where C_j is the censoring time such that $C_j < V_{j,s_j}$. The
 39 reason for $C_j < V_{j,s_j}$ is case 3b in Figure 1, where none of the matches is

1 correct, so the actual death time must be after C_j . Regarding subjects in
 2 Case 1–2 as having only $s_j = 1$ possibility with $V_{j,1} = V_j$, we use a unified
 3 notation for the observed data from subject j

$$(T_{j,k}, \Delta_{j,k}, \mathbf{x}_j) : k \in \{1, \dots, s_j\},$$

4 where \mathbf{x}_j is a p -dimensional vector of predictors, $T_{j,k} = \min(V_{j,k}, C_j)$, $\Delta_{j,k} =$
 5 $\mathbf{1}(V_{j,k} \leq C_j)$, and C_j is the censoring time. For cases 1–2, $\Delta_{j,1}$ is the event
 6 indicator and the notation is the same as in standard right-censored data.
 7 For Case 3, we have $s_j \geq 2$; $\Delta_{j,1} = \dots = \Delta_{j,s_j-1} = 1$ and $\Delta_{j,s_j} = 0$ are
 8 indicators denoting that all the matches before C_j are possible events and
 9 the last possibility is always censored. These notations will be used in the
 10 estimation procedure.

11 The true event time V_j of subject j , $j \in \{1, \dots, n\}$, is assumed to follow
 12 a Cox model with hazard function

$$(1) \quad h_j(t) = h_0(t) \exp(\mathbf{x}_j^\top \boldsymbol{\beta}),$$

13 where $h_0(\cdot)$ is an unspecified baseline function, and $\boldsymbol{\beta}$ is a vector of unknown
 14 coefficient of the covariate vector \mathbf{x}_j . Let $S_j(t) = \exp\{-H_0(t) \exp(\mathbf{x}_j^\top \boldsymbol{\beta})\}$,
 15 where $H_0(t) = \int_0^t h(s) ds$, be the survival function of subject j . The density
 16 function is then $f_j(t) = h_j(t) S_j(t)$. In addition, we assume that the censoring
 17 time C_j has an unknown density function $g(t)$, distribution function $G(t)$,
 18 survival function $\bar{G}(t) = 1 - G(t)$, does not depend on the covariates \mathbf{x}_j ,
 19 and is independent of the event times conditional on the covariates \mathbf{x}_j . The
 20 conditional independence assumption of the censoring time is justified for
 21 our study because the censoring was administrative.

22 We propose to model the uncertain records in a probabilistic way by
 23 introducing a vector of truth indicator for each subject. For subject j , let
 24 $\mathbf{Z}_j = (Z_{j,1}, \dots, Z_{j,s_j})$ be a random vector from multinomial distribution
 25 $\text{Multi}(\mathbf{1}, \boldsymbol{\pi}_j)$,

$$Z_{j,k} = \begin{cases} 1, & V_j = V_{j,k}, \text{ or } (T_{j,k}, \Delta_{j,k}) \text{ is the truth} \\ 0, & \text{otherwise} \end{cases},$$

26 where $k \in \{1, \dots, s_j\}$, $\sum_{k=1}^{s_j} Z_{j,k} = 1$, $0 \leq \pi_{j,k} \leq 1$ and $\sum_{k=1}^{s_j} \pi_{j,k} = 1$.
 27 As such, for each subject j , $j \in \{1, \dots, n\}$, $\boldsymbol{\pi}_j = (\pi_{j,1}, \dots, \pi_{j,s_j})$ is the
 28 probability vector where $\pi_{j,k} = \Pr(V_j = V_{j,k})$ (i.e., probability of the k -th
 29 record being true). Clearly, for $j \in I_1 \cup I_2$, we have $s_j = 1$ and $\pi_{j,1} = 1$,
 30 i.e., $Z_{j,1} = 1$ with probability 1. For $j \in I_3$, however, the truth indicators
 31 can be regarded as missing. That $Z_{j,k} = 1$, $k \in \{1, \dots, s_j - 1\}$, corresponds
 32 Case 3a, while $Z_{j,s_j} = 1$ suggests Case 3b.

1 Let $\mathbf{T}_j = (T_{j,1}, \dots, T_{j,s_j})$ and $\mathbf{\Delta}_j = (\Delta_{j,1}, \dots, \Delta_{j,s_j})$, with realizations
 2 $\mathbf{t}_j = (t_{j,1}, \dots, t_{j,s_j})$ and $\mathbf{\delta}_j = (\delta_{j,1}, \dots, \delta_{j,s_j})$, respectively. Let the set of
 3 all model parameters be $\boldsymbol{\theta} = \{\boldsymbol{\beta}, \boldsymbol{\pi}, h_0(\cdot), g(\cdot)\}$, where $\boldsymbol{\pi} = (\boldsymbol{\pi}_1, \dots, \boldsymbol{\pi}_n)$.
 4 Let \mathbf{z}_j be a realization of \mathbf{Z}_j . Given the truth indicators, we assume that
 5 the distribution of the fake records is independent of the true record and
 6 degenerates to a point mass at the point of the observed fake records. This
 7 assumption allows us to get away with modeling the intractable distribution
 8 of the fake records (e.g., the fake death times produced from imperfect data
 9 matching in our suicide risk study), so that the likelihood of $(\mathbf{T}_j, \mathbf{\Delta}_j)$ given
 10 \mathbf{Z}_k only depends on the likelihood of the true record. The complete-data
 11 likelihood of $(\mathbf{T}_j, \mathbf{\Delta}_j, \mathbf{Z}_j)$ from subject j turns out to be

$$(2) \quad L_j^C(\boldsymbol{\theta}) = \prod_{k=1}^{s_j} \left\{ \pi_{j,k} [f_j(t_{j,k}) \overline{G}(t_{j,k})]^{\delta_{j,k}} [g(t_{j,k}) S_j(t_{j,k})]^{1-\delta_{j,k}} \right\}^{z_{j,k}}.$$

12 The derivation detail is available in Section 1 of the Supplementary Ma-
 13 terials (Wang et al., 2019). All the possible realizations of \mathbf{Z}_j are $\mathbf{z}_j =$
 14 $(1, 0, 0, \dots, 0)$, $(0, 1, 0, \dots, 0)$, \dots , $(0, 0, \dots, 0, 1)$. The observed-data likeli-
 15 hood contribution from subject j is then obtained by summing out \mathbf{z}_j in (2):

$$(3) \quad L_j^O(\boldsymbol{\theta}) = \sum_{k=1}^{s_j} \pi_{j,k} [f_j(t_{j,k}) \overline{G}(t_{j,k})]^{\delta_{j,k}} [g(t_{j,k}) S_j(t_{j,k})]^{1-\delta_{j,k}}.$$

16 Let $\mathbf{Y}_{\text{obs}} = \{(\mathbf{t}_1, \boldsymbol{\delta}_1, \mathbf{x}_1), \dots, (\mathbf{t}_n, \boldsymbol{\delta}_n, \mathbf{x}_n)\}$ denote the observed data of the n
 17 independent subjects. The likelihood for the observed data is then given by
 18 $L^O(\boldsymbol{\theta}) = \prod_{j=1}^n L_j^O(\boldsymbol{\theta})$.

19 Thus far the observed-date likelihood in (3) is derived from a missing data
 20 perspective, but it can also be understood in several different ways. Intu-
 21 itively, for subject j , each of its s_j records leads to a likelihood of the event
 22 time and the censoring time, i.e., $[f_j(t_{j,k}) \overline{G}(t_{j,k})]^{\delta_{j,k}} [g(t_{j,k}) S_j(t_{j,k})]^{1-\delta_{j,k}}$ for
 23 $k \in \{1, \dots, s_j\}$, and the $L_j^O(\boldsymbol{\theta})$, the contribution of subject j to $L^O(\boldsymbol{\theta})$, is
 24 then constructed as a weighted sum with weights $\pi_{j,k}$ satisfying $0 \leq \pi_{j,k} \leq 1$
 25 and $\sum_{k=1}^{s_j} \pi_{j,k} = 1$. From the perspective of finite mixture model, the $\pi_{j,k}$'s
 26 are the mixing probabilities, and the above likelihood form of each mixture
 27 component is a direct consequence of our assumption that given the truth
 28 indicator the distribution of the fake records degenerates such that the dis-
 29 tribution of $(\mathbf{T}_j, \mathbf{\Delta}_j)$ only depends on the true record. Interestingly, the pro-
 30 posed method is also connected to a trimmed likelihood approach (e.g., Hadi
 31 and Luceño, 1997; Neykov et al., 2007), for which, however, the optimiza-
 32 tion problem is combinatorial in nature and a naive exhaustive search is not

1 feasible; see Section 4.4 for details. In contrast, the proposed probabilistic
 2 formulation allows us to develop an ECM algorithm to conduct maximum
 3 likelihood estimation. We remark that our approach may allow potential in-
 4 corporation of certain known missing mechanism of the true label, through
 5 imposing more structures on $\boldsymbol{\pi}_j$ or modeling them using covariates. For in-
 6 stance, in some applications it may be reasonable to assume that the prior
 7 probability of being censored is the same for all the subjects with uncertain
 8 records. In this work, however, we focus on the unconstrained situation.

9 4. Model Estimation via an ECM Algorithm.

10 4.1. *Estimation Procedure.* The ECM algorithm is a variation of the
 11 powerful EM algorithm for dealing with incomplete data (Meng and Rubin,
 12 1993). It replaces the M-step of an EM algorithm with multiple conditional
 13 maximization (CM) steps which are often computationally easier to handle.
 14 We propose a maximum likelihood estimation procedure for the integrative
 15 Cox model following the architecture of the ECM, in which the CM-steps
 16 utilize a profile likelihood similar to the partial likelihood (Cox, 1975).

17 The complete-data loglikelihood can be decomposed into two parts which
 18 involve two exclusive sets of parameters, respectively. Let $\mathbf{Y}_{\text{mis}} = (z_1, \dots, z_n)$
 19 and $\mathbf{Y} = \{\mathbf{Y}_{\text{obs}}, \mathbf{Y}_{\text{mis}}\}$. From (2), the complete-data loglikelihood is

$$(4) \quad \ell(\boldsymbol{\theta} \mid \mathbf{Y}) = \ell(\boldsymbol{\beta}, \boldsymbol{\pi}, h_0(\cdot) \mid \mathbf{Y}) + \ell_c(g(\cdot) \mid \mathbf{Y}),$$

20 where

$$(5) \quad \begin{aligned} & \ell(\boldsymbol{\beta}, \boldsymbol{\pi}, h_0(\cdot) \mid \mathbf{Y}) \\ &= \sum_{j=1}^n \sum_{k=1}^{s_j} z_{j,k} \{ \log \pi_{j,k} + \delta_{j,k} \log f_j(t_{j,k}) + (1 - \delta_{j,k}) \log S_j(t_{j,k}) \}, \end{aligned}$$

21 and

$$(6) \quad \ell_c(g(\cdot) \mid \mathbf{Y}) = \sum_{j=1}^n \sum_{k=1}^{s_j} z_{j,k} \{ \delta_{j,k} \log \bar{G}(t_{j,k}) + (1 - \delta_{j,k}) \log g(t_{j,k}) \}.$$

22 The second part $\ell_c(g(\cdot) \mid \mathbf{Y})$ only involves the nuisance distribution of the
 23 censoring time.

24 We compute the conditional expectations of the complete-data loglikeli-
 25 hood (4) given the observed data \mathbf{Y}_{obs} and the set of parameter estimates
 26 $\boldsymbol{\theta}^{(i)} = \{\boldsymbol{\beta}^{(i)}, \boldsymbol{\pi}^{(i)}, h_0^{(i)}(\cdot), g^{(i)}(\cdot)\}$ at i -th iteration ($i = 0, 1, \dots$), where $\boldsymbol{\theta}^{(0)}$ is
 27 the initial/starting estimate. Define

$$w_{j,k}(\boldsymbol{\theta}^{(i)}) := P(Z_{j,k} = 1, \mathbf{T}_j, \boldsymbol{\Delta}_j \mid \boldsymbol{\theta}^{(i)}) = \pi_{j,k}^{(i)} \left(h_{j,k}^{(i)} S_{j,k}^{(i)} \bar{G}_{j,k}^{(i)} \right)^{\delta_{j,k}} \left(g_{j,k}^{(i)} S_{j,k}^{(i)} \right)^{1-\delta_{j,k}},$$

1 where $h_{j,k}^{(i)} = h_0^{(i)}(t_{j,k}) \exp(\mathbf{x}_j^\top \boldsymbol{\beta}^{(i)})$ and $S_{j,k}^{(i)} = \exp\{-H_0^{(i)}(t_{j,k}) \exp(\mathbf{x}_j^\top \boldsymbol{\beta}^{(i)})\}$,
 2 $\bar{G}_{j,k}^{(i)} = \bar{G}^{(i)}(t_{j,k})$, and $g_{j,k}^{(i)} = g^{(i)}(t_{j,k})$. By Bayes rule, we have

$$(7) \quad p_{j,k}(\boldsymbol{\theta}^{(i)}) := \mathbb{P}(Z_{j,k} = 1 \mid \mathbf{T}_j, \boldsymbol{\Delta}_j, \boldsymbol{\theta}^{(i)}) = \frac{w_{j,k}(\boldsymbol{\theta}^{(i)})}{\sum_{k=1}^{s_j} w_{j,k}(\boldsymbol{\theta}^{(i)})}.$$

3 Plugging (7) into (5) and (6), we obtain the E-step that involves two
 4 separate parts:

$$(8) \quad \mathbb{E} \ell\{\boldsymbol{\beta}, \boldsymbol{\pi}, h_0(\cdot) \mid \mathbf{Y}_{\text{obs}}, \boldsymbol{\theta}^{(i)}\} \\ = \sum_{j=1}^n \sum_{k=1}^{s_j} p_{j,k}(\boldsymbol{\theta}^{(i)}) \left\{ \log(\pi_{j,k}) + \delta_{j,k} \log f_j(t_{j,k}) + (1 - \delta_{j,k}) \log S(t_{j,k}) \right\},$$

5 and

$$(9) \quad \mathbb{E} \ell_c\{g(\cdot) \mid \mathbf{Y}_{\text{obs}}, \boldsymbol{\theta}^{(i)}\} \\ = \sum_{j=1}^n \sum_{k=1}^{s_j} p_{j,k}(\boldsymbol{\theta}^{(i)}) \left\{ \delta_{j,k} \log \bar{G}(t_{j,k}) + (1 - \delta_{j,k}) \log g(t_{j,k}) \right\}.$$

6 The separation of the two terms in parameters facilitates the M-step. The
 7 first term (8) can be handled by profiling out the nuisance parameters. Note
 8 that, for fixed $\boldsymbol{\beta}$ and $\boldsymbol{\pi}$, the $h_0(t)$ maximizing the conditional expectation (8)
 9 is a discrete function that is positive only at possible event times and zero
 10 anywhere else. Let $Y_{j,k}(t) = \mathbf{1}(t_{j,k} \geq t)$ and $N_{j,k}(t) = z_{j,k} \mathbf{1}(t_{j,k} \leq t, \delta_{j,k} = 1)$.
 11 Then the true number of events by time t is $N(t) = \sum_{j=1}^n \sum_{k=1}^{s_j} N_{j,k}(t)$.
 12 Let $dN(t)$ denote the number of true events at time t . Let $\tilde{N}_{j,k}(t; \boldsymbol{\theta}^{(i)}) =$
 13 $p_{j,k}(\boldsymbol{\theta}^{(i)}) \mathbf{1}(t_{j,k} \leq t, \delta_{j,k} = 1)$ and $\tilde{N}(t; \boldsymbol{\theta}^{(i)}) = \sum_{j=1}^n \sum_{k=1}^{s_j} \tilde{N}_{j,k}(t; \boldsymbol{\theta}^{(i)})$, which
 14 are the conditional expectation of $N_{j,k}(t)$ and $N(t)$ given \mathbf{Y}_{obs} , evaluated at
 15 $\boldsymbol{\theta}^{(i)}$, respectively. Then $d\tilde{N}(t; \boldsymbol{\theta}^{(i)}) = \mathbb{E}\{dN(t) \mid \mathbf{Y}_{\text{obs}}, \boldsymbol{\theta}^{(i)}\}$ is the jump size of
 16 $\tilde{N}(t; \boldsymbol{\theta}^{(i)})$ at time t . Equation (8) can be rewritten to allow tied event times
 17 as follows:

$$(10) \quad \mathbb{E} \ell\{\boldsymbol{\beta}, \boldsymbol{\pi}, h_0(\cdot) \mid \mathbf{Y}_{\text{obs}}, \boldsymbol{\theta}^{(i)}\} \\ = \sum_{t \in \mathcal{T}} \left[-h_0(t) \sum_{j=1}^n \sum_{k=1}^{s_j} Y_{j,k}(t) p_{j,k}(\boldsymbol{\theta}^{(i)}) \exp(\mathbf{x}_j^\top \boldsymbol{\beta}) + d\tilde{N}(t; \boldsymbol{\theta}^{(i)}) \log h_0(t) \right] \\ + \sum_{j=1}^n \sum_{k=1}^{s_j} p_{j,k}(\boldsymbol{\theta}^{(i)}) \left[\delta_{j,k} \mathbf{x}_j^\top \boldsymbol{\beta} + \log \pi_{j,k} \right],$$

1 where $\mathcal{T} = \{t_{j,k} \mid \delta_{j,k} = 1, k \in \{1, \dots, s_j\}, j \in \{1, \dots, n\}\}$ is the collection
 2 of all observed possible event times.

3 Given $\boldsymbol{\beta}$ and $\boldsymbol{\pi}$, the baseline hazard h_0 only appears in the first term
 4 of (10), and the maximizer is

$$\hat{h}_0(t) = \frac{d\tilde{N}(t; \boldsymbol{\theta}^{(i)})}{\sum_{j=1}^n \sum_{k=1}^{s_j} Y_{j,k}(t) p_{j,k}(\boldsymbol{\theta}^{(i)}) \exp(\mathbf{x}_j^\top \boldsymbol{\beta})},$$

5 which is nonzero only for those $t \in \mathcal{T}$, similar to the ‘‘Breslow estima-
 6 tor’’ (Breslow, 1974). Further, for fixed $\boldsymbol{\beta}$, it is easy to check that $\pi_{j,k}^{(i+1)} =$
 7 $p_{j,k}(\boldsymbol{\theta}^{(i)})$ maximizes (10) by Lagrange multipliers method. Plugging these
 8 estimators back into (10), we get a profile likelihood in terms of $\boldsymbol{\beta}$

$$\begin{aligned} & \mathbb{E} \ell\{\boldsymbol{\beta}, \hat{\boldsymbol{\pi}}, \hat{h}_0 \mid \mathbf{Y}_{\text{obs}}, \boldsymbol{\theta}^{(i)}\} \\ &= \sum_{t \in \mathcal{T}} \left\{ -d\tilde{N}(t; \boldsymbol{\theta}^{(i)}) \left[1 - \log d\tilde{N}(t; \boldsymbol{\theta}^{(i)}) \right] \right\} + \sum_{j=1}^n \sum_{k=1}^{s_j} p_{j,k}(\boldsymbol{\theta}^{(i)}) \log p_{j,k}(\boldsymbol{\theta}^{(i)}) \\ & \quad + p\ell(\boldsymbol{\beta} \mid \boldsymbol{\theta}^{(i)}), \end{aligned}$$

9 where

$$(11) \quad p\ell(\boldsymbol{\beta} \mid \boldsymbol{\theta}^{(i)}) = \sum_{j=1}^n \sum_{k=1}^{s_j} \int_0^\infty I(\boldsymbol{\beta}, t \mid \boldsymbol{\theta}^{(i)}) d\tilde{N}_{j,k}(t; \boldsymbol{\theta}^{(i)}),$$

$$I(\boldsymbol{\beta}, t \mid \boldsymbol{\theta}^{(i)}) = \mathbf{x}_j^\top \boldsymbol{\beta} - \log \left(\sum_{l=1}^n \sum_{m=1}^{s_l} Y_{l,m}(t) p_{l,m}(\boldsymbol{\theta}^{(i)}) \exp(\mathbf{x}_l^\top \boldsymbol{\beta}) \right),$$

10 is the only part involving $\boldsymbol{\beta}$. This profiling approach is similar to the partial
 11 likelihood of Cox (1975) except that the distribution of the censoring time
 12 comes into play through $p_{j,k}$ ’s and $d\tilde{N}_{j,k}$ ’s. The estimator $\hat{\boldsymbol{\beta}}$ of $\boldsymbol{\beta}$ is obtained
 13 by maximizing (11). Once $\hat{\boldsymbol{\beta}}$ has converged, $\hat{h}_0(\cdot)$ and $\hat{\pi}_{j,k}$ ’s can be updated.

14 Maximizing the second part (9) involves nonparametric maximum likeli-
 15 hood estimator of the censoring distribution function $G(\cdot)$. We characterize
 16 the censoring time by its hazard function $h_c(\cdot)$. Similar to $h_0(t)$, the $h_c(\cdot)$
 17 that maximizes (9) is nonzero only at the observed censoring times. By the
 18 assumption we made, the only possible censoring time for each subject is its
 19 last record time. For $j \in \{1, \dots, n\}$, define $C_j(t; \boldsymbol{\theta}^{(i)}) = p_{j,s_j}(\boldsymbol{\theta}^{(i)}) \mathbf{1}(t_{j,s_j} \leq$
 20 $t, \delta_{j,s_j} = 0)$ and $C(t; \boldsymbol{\theta}^{(i)}) = \sum_{j=1}^n C_j(t; \boldsymbol{\theta}^{(i)})$. Let $dC(t; \boldsymbol{\theta}^{(i)})$ be the jump
 21 size of $C(t; \boldsymbol{\theta}^{(i)})$ at time t . Then we may rewrite (9) to allow tied censoring
 22 times as follows:

$$\mathbb{E} \ell_c(g(\cdot) \mid \mathbf{Y}_{\text{obs}}, \boldsymbol{\theta}^{(i)}) = \sum_{t \in \mathcal{C}} \left\{ dC(t; \boldsymbol{\theta}^{(i)}) \log h_c(t) - h_c(t) \sum_{j=1}^n \sum_{k=1}^{s_j} p_{j,k}(\boldsymbol{\theta}^{(i)}) Y_{j,k}(t) \right\},$$

1 where $\mathcal{C} = \{t_{j,s_j} \mid \delta_{j,s_j} = 0, j \in \{1, \dots, n\}\}$ is the collection of all observed
 2 censoring times. Maximizing it with respect to $h_c(t)$ gives

$$\hat{h}_c(t) = \frac{dC(t; \boldsymbol{\theta}^{(i)})}{\sum_{j=1}^n \sum_{k=1}^{s_j} p_{j,k}(\boldsymbol{\theta}^{(i)}) Y_{j,k}(t)}.$$

3 Therefore, for every record time $t_{j,k}$, we have

$$\hat{G}(t_{j,k}) = \exp \left\{ - \sum_{t \leq t_{j,k}} \hat{h}_c(t) \right\} = \exp \left\{ - \sum_{t \leq t_{j,k}} \frac{dC(t; \boldsymbol{\theta}^{(i)})}{\sum_{l=1}^n \sum_{m=1}^{s_l} p_{l,m}(\boldsymbol{\theta}^{(i)}) Y_{l,m}(t)} \right\}$$

4 and $\hat{g}(t_{j,k}) = \hat{h}_c(t) \hat{G}(t_{j,k})$.

5 We summarize the ECM estimation procedure in Algorithm 1. In our
 6 numerical studies, we stop the algorithm if $\|\boldsymbol{\beta}^{(i)} - \boldsymbol{\beta}^{(i-1)}\| / \|\boldsymbol{\beta}^{(i)} + \boldsymbol{\beta}^{(i-1)}\| <$
 7 10^{-6} and $\|\boldsymbol{\pi}^{(i)} - \boldsymbol{\pi}^{(i-1)}\| / \|\boldsymbol{\pi}^{(i)} + \boldsymbol{\pi}^{(i-1)}\| < 10^{-8}$.

8 **4.2. Initialization.** Since the maximum likelihood estimation problem
 9 here is non-convex, it may admit multiple local maxima. Therefore, we recom-
 10 mend setting multiple initial values of $\boldsymbol{\beta}$ and $\boldsymbol{\pi}$ to help identify a good
 11 solution, as allowed by the available computational resources. In particular,
 12 we propose two simple but pragmatic initialization procedures that work
 13 well even with limited resources.

14 The first procedure is as follows:

- 15 (i) Fit a regular Cox model on all the certain records (Case 1–2) and use
 16 the estimated coefficients to initialize $\boldsymbol{\beta}$; initialize $\hat{S}_{j,k}$ with the fitted
 17 survival function evaluated at $t_{j,k}$; initialize $\hat{h}_{j,k}$ with a nearest left
 18 neighbor interpolation of the fitted hazard function (if no left neighbor,
 19 use nearest right neighbor).
- 20 (ii) Switching event and censoring for all the certain records (Case 1–2),
 21 estimate the hazard function for censoring by the Nelson-Aalen estima-
 22 tor (without covariates) and obtain the corresponding survival function
 23 estimate; initialize $\hat{G}_{j,k}$ with the fitted survival function evaluated at
 24 $t_{j,k}$; initialize $\hat{h}_c(t_{j,k})$ with a nearest left neighbor interpolation of the
 25 fitted hazard function (if no left neighbor, use nearest right neighbor).
- 26 (iii) Plug $\hat{w}_{j,k} = h_{j,k}^* \hat{S}_{j,k} \hat{G}_{j,k}$, where $h_{j,k}^* = \delta_{j,k} \hat{h}_{j,k} + (1 - \delta_{j,k}) \hat{h}_c(t_{j,k})$,
 27 into (7) as $w_{j,k}$ and initialize $\pi_{j,k}$ as the resulting $p_{j,k}$.

Algorithm 1 Estimation procedure for integrative Cox model with uncertain event records. (The dependence of $\pi_{j,k}$'s, $\tilde{N}_{j,k}$'s, $d\tilde{N}(t)$, and $dC(t)$ on θ is dropped for ease of notation.)

initialize β and π ;
repeat
 for $j = 1, 2, \dots, n$ **do** ▷ Update $\tilde{N}_{j,k}(t)$'s
 for $k = 1, 2, \dots, s_j$ **do**

$$\tilde{N}_{j,k}(t) \leftarrow \pi_{j,k} \mathbf{1}(t_{j,k} \leq t, \delta_{j,k} = 1);$$

 end for
 end for
 for each $t \in \mathcal{T}$ **do** ▷ Update $\hat{h}_0(\cdot)$

$$h_0(t) \leftarrow \frac{d\tilde{N}(t)}{\sum_{j=1}^n \sum_{k=1}^{s_j} Y_{j,k}(t) \pi_{j,k} \exp(\mathbf{x}_j^\top \beta)}; H_0(t) \leftarrow \sum_{s \leq t} h_0(s);$$

 end for
 for each $t \in \mathcal{C}$ **do** ▷ Update $\hat{h}_c(\cdot)$

$$h_c(t) \leftarrow \frac{dC(t)}{\sum_{j=1}^n \sum_{k=1}^{s_j} Y_{j,k}(t) \pi_{j,k}}; H_c(t) \leftarrow \sum_{s \leq t} h_c(s),$$

 end for
 for $j = 1, 2, \dots, n$ **do** ▷ Update $\hat{\pi}_{j,k}$'s
 for $k = 1, 2, \dots, s_j$ **do**

$$S_{j,k} \leftarrow \exp \left\{ -H_0(t_{j,k}) \exp(\mathbf{x}_j^\top \beta) \right\}; \bar{G}_{j,k} \leftarrow \exp \left\{ -H_c(t_{j,k}) \right\};$$

$$w_{j,k} \leftarrow \pi_{j,k} [h_{j,k} S_{j,k} \bar{G}_{j,k}]^{\delta_{j,k}} [g_{j,k} S_{j,k}]^{1-\delta_{j,k}}; \pi_{j,k} \leftarrow \frac{w_{j,k}}{\sum_{k=1}^{s_j} w_{j,k}};$$

 end for
 end for

$$\beta \leftarrow \arg \max_{\beta} p\ell(\beta | \theta)$$
 ▷ Update $\hat{\beta}$
until Convergence

- 1 In the above procedure, letting $\hat{h}_{j,k}^* = 1$ in step (iii) leads to a simpler
2 alternative, which puts more weights to the uncertain event times before
3 the censoring time and thus may work better when Case 3a is estimated to
4 have a larger size than Case 3b. This gives a second initialization procedure.
5 The two initialization procedures were applied in the simulation studies
6 presented in Section 5 and the results were satisfactory in most scenarios.

- 7 **4.3. Inference.** In an EM or ECM algorithm, generally standard error
8 (SE) estimates for the parameter estimates cannot be easily produced along
9 with the estimation procedure. A few approaches have been proposed for

1 estimating the asymptotic covariance matrix for parameters of interest, in-
 2 cluding the supplemented EM (SEM) algorithm (Meng and Rubin, 1991),
 3 the profile likelihood approach (Murphy and van der Vaart, 2000), numerical
 4 differentiation methods based on forward difference and Richardson extrap-
 5 olation (Jamshidian and Jennrich, 2000), and their variants with profiling
 6 (Xu, Baines and Wang, 2014). Unfortunately, none of these methods is read-
 7 ily applicable to our case. In our work, we use the bootstrap (Efron, 1979,
 8 1981) method that performs resampling at the subject level for survival data
 9 for making inference. Efron (1981) proposed the SE be estimated as sample
 10 standard deviation of bootstrap estimates, or based on inter-quantile range
 11 and normal approximation. The p -values from the Wald test for testing the
 12 significance of each regression coefficient can then be computed.

13 4.4. *Connection with Trimmed Likelihood.* We show that the proposed
 14 method is closely connected to a trimmed likelihood approach, which offers
 15 an intuitive understanding of our method from the robust estimation per-
 16 spective. The trimmed likelihood (Rousseeuw, 1984; Hadi and Luceño, 1997;
 17 Neykov et al., 2007) is a general approach for conducting robust maximum
 18 likelihood estimation in the presence of outliers, in which the observations
 19 are trimmed according to their contributions to the likelihood function. Our
 20 probabilistic modeling approach via ECM provides an efficient way for tar-
 21 geting the computationally infeasible trimmed likelihood estimator.

22 Recall the observed-data likelihood formulation given in (3). Denote

$$r_{j,k}(\boldsymbol{\beta}) = [f_j(t_{j,k})\overline{G}(t_{j,k})]^{\delta_{j,k}} [g(t_{j,k})S_j(t_{j,k})]^{1-\delta_{j,k}},$$

23 where $j \in \{1, \dots, n\}$, $k \in \{1, \dots, s_j\}$. For ease of notation, here we do not
 24 explicitly write out the dependency of $r_{j,k}(\boldsymbol{\beta})$ on the observed data and
 25 assume other unknown quantities $h_0(\cdot)$ and $g(\cdot)$ have been profiled out. (In
 26 fact, the above can be regarded as a general survival modeling formulation
 27 in this section.) Then the proposed maximum likelihood estimator can be
 28 expressed as

$$(12) \quad (\hat{\boldsymbol{\beta}}, \hat{\boldsymbol{\pi}}) \in \arg \max_{\boldsymbol{\beta}, \boldsymbol{\pi}} \prod_{j=1}^n \left(\sum_{k=1}^{s_j} \pi_{j,k} r_{j,k}(\boldsymbol{\beta}) \right).$$

29 Here each $\boldsymbol{\pi}_j$ is a probability vector and there is no additional structural
 30 constraint on $\boldsymbol{\pi}$. Now, for each j , define $r_{j,(s_j)}(\boldsymbol{\beta})$ as the largest order statis-
 31 tic of $r_{j,k}(\boldsymbol{\beta}), k = 1, \dots, s_j$. Then, a trimmed likelihood estimator can be
 32 constructed as

$$(13) \quad \tilde{\boldsymbol{\beta}} \in \arg \max_{\boldsymbol{\beta}} \prod_{j=1}^n r_{j,(s_j)}(\boldsymbol{\beta}).$$

1 Intuitively, (13) shows that the optimal β is reached when for each patient
 2 with uncertain records, only the most plausible record (as judged by having
 3 the largest log-likelihood value among all the records) contributes to the
 4 overall log-likelihood function and the rest all get trimmed. Interestingly, it
 5 can be verified that the two methods in (12) and (13) share the same set of
 6 global solutions.

7 LEMMA 4.1. *The $\hat{\beta}$ from solving (12) is a solution of (13), and vice*
 8 *versa.*

9 To see this, note that for each $\pi_j = (\pi_{j,1}, \dots, \pi_{j,s_j})$, we have $\hat{\pi}_j =$
 10 $\arg \max_{\pi_j} \sum_{k=1}^{s_j} \pi_{j,k} r_{j,k}(\hat{\beta})$, because given $\beta = \hat{\beta}$ the problem in (12) is sep-
 11 arable in each set of π_j . Then the maximum can be attained at $\hat{\pi}_{j,k_j^0} = 1$ and
 12 $\hat{\pi}_{j,k} = 0$ for $k \neq k_j^0$, where $k_j^0 \in \arg \max_k r_{j,k}(\hat{\beta})$. It follows that the maxi-
 13 mum value of the objective function in (12) can be written as $\prod_{j=1}^n r_{j,(s_j)}(\hat{\beta})$,
 14 which, clearly reveals that $\hat{\beta}$ is a maximizer of the trimmed likelihood prob-
 15 lem in (13). On the other hand, let $\tilde{\pi}_j$ be that $\tilde{\pi}_{j,k_j^0} = 1$ and $\tilde{\pi}_{j,k} = 0$ for
 16 $k \neq k_j^0$, where $k_j^0 \in \arg \max_k r_{j,k}(\tilde{\beta})$ with some abuse of notation. Then it
 17 can be seen that $\{\tilde{\beta}, \tilde{\pi}\}$ is necessarily a solution of (12).

18 In practice, however, finding the global solution of the trimmed likelihood
 19 problem is infeasible via a naive exhaustive search approach. For example,
 20 in our suicide risk study, an exhaustive search amounts to fit $2^{584} \times 3^{39} \times 5^2$
 21 many Cox models. In contrast, our probabilistic modeling approach can be
 22 regarded as an efficient way for targeting the trimmed likelihood estimator
 23 via the ECM algorithm, with carefully constructed initial values.

24 5. Simulation Study.

25 5.1. *Simulation Settings.* Our simulation settings were designed to mimic
 26 the data integration process in the survival analysis of patients admitted to
 27 hospital due to unsuccessful suicide attempts in Connecticut. As shown in
 28 Figure 1, n_1 is the number of subjects with events observed for certain from
 29 dataset I (Case 1); n_2 is the number of subjects whose true event time is
 30 included in the matched event times (Case 3a); n_3 is the number of subjects
 31 whose true event time is censored but for whom some false event times are
 32 matched (Case 3b); n_4 is the number of subjects whose event times are cen-
 33 sored for certain since no match is found from dataset II (Case 2). As such,
 34 $n = \sum_{i=1}^4 n_i$ is the total sample size, and $n_2 + n_3 + n_4$ is the number of
 35 subjects that are censored before data matching.

TABLE 1
 Summary of different simulation settings. The number of subjects in Group 1 is fixed at $n_1 + n_2 = 200$.

Scenario #	CR1	MR	CMR	Group 1		Group 2		n	OCR
				n_1 (Case 1)	n_2 (Case 3a)	n_3 (Case 3b)	n_4 (Case 2)		
1	30	70	20	189	11	46	24	270	26
2	30	70	80	161	39	9	21	230	13
3	60	40	20	178	22	84	160	444	55
4	60	40	80	136	64	17	122	339	41
5	90	10	20	167	33	117	1350	1667	88
6	90	10	80	118	82	24	953	1177	83

CR1: Censoring rate before matching (%); MR: Matching rate (%); CMR: Correct matching rate (%); OCR: Oracle censoring rate (%).

1 We define a few quantities for designing the experiment: censoring rate
 2 of dataset I before matching (CR1) is $CR1 = 1 - n_1/n$; matching rate
 3 (MR) is $MR = (n_2 + n_3)/(n_2 + n_3 + n_4)$; correct matching rate (CMR)
 4 $CMR = n_2/(n_2 + n_3)$. MR is the proportion of subjects having matched
 5 records among subjects whose event times are censored from dataset I; CMR
 6 is the proportion of the subjects whose true event time is contained in the
 7 matched event times. In all the settings, we set $MR = 1 - CR1$, assuming
 8 that the lower the CR1, the more likely that dataset I misses true events
 9 among the censored records. The number of subjects who actually had events
 10 was fixed at $n_1 + n_2 = 200$ to keep an approximately same benchmark
 11 performance from oracle models under different settings.

12 Three levels of CR1 were considered, i.e., $CR1 \in \{30\%, 60\%, 90\%\}$, corre-
 13 sponding to moderate, heavy, and severe censoring, respectively. Two levels
 14 of CMR were considered, i.e., $CMR \in \{20\%, 80\%\}$; the larger the CMR,
 15 the more valuable information can be potentially recovered from dataset II.
 16 Given (CR1, MR, CMR) and with the condition $n_1 + n_2 = 200$, the values
 17 of n_i 's, $i = 1, \dots, 4$ were then completely determined. Table 1 summarizes
 18 the sample size and its decomposition into the four cases, for each of the 6
 19 simulation scenarios determined by the combinations of CR1 and CMR.

20 For ease of data generation, we divide the subjects into two groups:
 21 Group 1 contains those whose true event times are included in the observed
 22 data, not necessarily certain though (Case 1 and Case 3a); Group 2 con-
 23 tains those whose true event times are not in the observed data (Case 2
 24 and Case 3b). Define oracle censoring rate (OCR), $OCR = (n_3 + n_4)/n$, the
 25 proportion of Group 2 in the sample, which is unobserved but completely
 26 determined for each setting after the values of n_i 's are determined. Our

1 strategy was to generate true event time and censoring time for all subjects
2 for a given OCR first, identify subjects in Case 3a from Group 1, identify
3 subjects in Case 3b from Group 2, and then generate fake event times for
4 those in Case 3a and Case 3b, respectively.

5 The true event times were generated from Cox model (1) with a Weibull
6 baseline hazard function. Four independent covariates were included in the
7 model; the first three were from the standard normal distribution and the
8 fourth was from the Bernoulli distribution with rate 0.5. All four true regres-
9 sion coefficients were set to be 1. The censoring time was generated from the
10 uniform distribution over (0.5, 12.5). The Weibull-shape parameter was set
11 to be 2, 1, and 0.7 for the moderate, heavy, and severe censoring scenarios in
12 terms of CR1, respectively. The Weibull-scale parameter was tuned in each
13 setting so that the OCR determined in that setting is attained on average.

14 To identify Case 3a subjects from Group 1 and Case 3b subjects from
15 Group 2, we treated the data uncertainty as a missing-label problem: the
16 labels are observed for the $n_1 + n_4$ subjects in Cases 1–2, but are missing
17 for the $n_2 + n_3$ subjects in Case 3. Two missing mechanisms were considered
18 for the labels: missing completely at random (MCAR) and missing not at
19 random (MNAR). In the MCAR mechanism, the probability of a label be-
20 ing missing was completely random, regardless of the underlying true event
21 time. In the MNAR mechanism, the probability of a label being missing
22 was proportional to the true event time; the longer the true event time, the
23 more likely a subject was identified as Case 3a from Group 1 or Case 3b from
24 Group 2. Such decomposition ensures that the sample size decomposition of
25 each simulated data closely matches its corresponding setting in Table 1.

26 The last step was to generate fake event times for subjects in Case 3.
27 For subjects in Case 3a, their censoring times were observed and true event
28 times were included in the matches. The number of additional fake event
29 times was set to be zero or one with probability 0.9 and 0.1, respectively. In
30 other words, the possible records for each of them consisted of one observed
31 censoring time, one true event time, and one additional fake event time with
32 probability 0.1. For subjects in Case 3b, their true event times were censored
33 and the number of fake event times was set to be one or two with probabili-
34 ty 0.9 and 0.1, respectively. In other words, each of them had one observed
35 censoring time, one or two fake event times with probability 0.9 or 0.1, re-
36 spectively. Each fake event time was generated from Cox model (1) with
37 one extra covariate in addition to the existing four covariates, conditional
38 on that the fake event time was less than the censoring time (Nadarajah and
39 Kotz, 2006). This extra covariate took value -1 or 1 with equal probability,
40 and its coefficient was set to be 3.

1 5.2. *Competing Methods and Evaluation Metrics.* Three competing meth-
2 ods were considered, multiple imputation (MI) and two naive approaches.
3 MI was originally introduced by Rubin (1987) for non-response in surveys,
4 which imputes every missing value multiple times with draws from certain
5 distribution and summarized the results from the multiple versions of the
6 complete data. In our setup, the missing values are the truth indicators.
7 Given a simulated dataset, we imputed 200 times the truth indicators for
8 the subjects in Case 3 and took the average of the coefficient estimates from
9 fitting the regular Cox model with each imputed data as the final estimates.
10 Specifically, in each imputation, for each subject the truth indicator vec-
11 tor was generated from a multinomial distribution, where the probability of
12 censoring was set to be proportional to $n_4/(n_1+n_4)$, and the remaining prob-
13 ability was equally split among the uncertain event records. The two naive
14 approaches were based on the regular Cox model as well. The first (denoted
15 by C.Cox) fits the regular Cox model to dataset I, which treats all subjects
16 in Case 3 as censored, completely ignoring integration with dataset II. C.Cox
17 may give biased estimator for not considering the events missed by dataset I.
18 The second approach (denoted by U.Cox) excludes those subjects with mul-
19 tiple event times after matching with dataset II (Case 3) and fits the regular
20 Cox model with the remaining subjects with unique records (Case 1 and
21 Case 2). The data used by U.Cox is a subset of that used by C.Cox. By re-
22 moving subjects in dataset I whose event times were not uniquely recorded,
23 U.Cox may give less efficient but unbiased estimation under MCAR.

24 The proposed integrative Cox model is denoted by I.Cox. We also included
25 two oracle procedures where the true event indicators are known a priori:
26 the oracle Cox model (O.Cox) and the oracle Weibull model (O.Weibull).
27 They give the best achievable performances, infeasible in practice but can
28 be used as references in comparison.

29 We measured the estimation performance by the ℓ_2 -norm of $(\hat{\beta} - \beta_0)$, i.e.,
30 $\|\hat{\beta} - \beta_0\| = [(\hat{\beta} - \beta_0)^\top (\hat{\beta} - \beta_0)]^{1/2}$, where β_0 is the underlying true coefficient
31 vector, and $\hat{\beta}$ is its estimator. In addition, we estimated the baseline survival
32 functions from the purposed I.Cox model and two naive Cox methods, and
33 compared them with the true parametric curve over a tense time grid from
34 0 to 12 with step size of 0.1. For each subject with multiple records, the
35 estimated probabilities $\hat{\pi}_j$ from the proposed I.Cox model can be used to
36 identify the true record. We used the Bayes rule to select the record with
37 the largest estimated probability; by comparing to the underlying truth,
38 we computed the correct identification rate of the true records among the
39 subjects having uncertain records. The experiment was replicated 1,000 time
40 under each setting and the results were then averaged.

TABLE 2

Comparison on parameter estimation performance through mean of $100 \times \|\hat{\beta} - \beta_0\|$ (with the standard deviation given in parenthesis).

#	O.Weibull	O.Cox	I.Cox	U.Cox	C.Cox	MI
MCAR						
1	18.0 (7.3)	20.7 (8.6)	22.4 (9.7)	24.9 (9.9)	81.1 (22.0)	81.0 (9.9)
2	17.5 (7.7)	20.8 (8.9)	22.1 (9.6)	23.4 (10.0)	139.7 (14.7)	80.8 (11.0)
3	18.5 (7.9)	19.7 (8.7)	23.6 (10.2)	22.3 (9.2)	55.9 (17.2)	81.2 (9.9)
4	18.6 (7.6)	20.3 (8.7)	22.8 (10.1)	26.0 (11.4)	107.7 (15.9)	94.2 (12.0)
5	18.0 (8.0)	18.2 (8.1)	20.6 (9.0)	20.2 (9.0)	30.5 (12.3)	39.5 (11.7)
6	18.4 (8.2)	18.8 (8.3)	22.2 (9.8)	30.1 (12.6)	51.8 (12.3)	51.4 (11.3)
MNAR						
1	18.0 (7.3)	20.7 (8.6)	22.4 (9.4)	27.6 (10.6)	48.4 (21.0)	81.5 (10.0)
2	17.5 (7.7)	20.8 (8.9)	22.4 (9.7)	24.4 (10.4)	79.9 (20.9)	55.2 (11.5)
3	18.5 (7.9)	19.7 (8.7)	22.5 (10.2)	22.9 (9.7)	27.1 (11.8)	86.4 (8.9)
4	18.6 (7.6)	20.3 (8.7)	23.2 (10.2)	24.2 (10.7)	38.3 (15.0)	51.6 (12.0)
5	18.0 (8.0)	18.2 (8.1)	21.1 (9.2)	20.7 (9.3)	21.1 (9.0)	40.2 (9.3)
6	18.4 (8.2)	18.8 (8.3)	23.5 (10.3)	30.2 (13.0)	27.3 (11.3)	30.1 (10.8)

MCAR: Missing completely at random; MNAR: Missing not at random.

1 5.3. *Simulation Results.* Table 2 summarizes the simulation results on
2 parameter estimation. As expected, the two practically-infeasible oracle ap-
3 proaches perform the best, which provide benchmarks for comparison. The
4 I.Cox method and the U.Cox method appear to have a clear advantage over
5 the C.Cox method and the MI method under most settings. The disadvan-
6 tage of the C.Cox method is expected; subjects in Case 3a are mistakenly
7 treated as censored, which increases the variance in estimation due to less
8 events and introduces bias due to the mistakenly treated censoring. Its per-
9 formance is even worse in MCAR settings, because in the MNAR setting,
10 longer survival time is more likely to be uncertain such that the true event
11 time is more likely to be close to the censoring time than under MCAR. The
12 MI method performs worse than the C.Cox method in the setting with lower
13 CMR under MNAR, unlike in other settings where they are less different,
14 because the imputation does not account for the informative missingness
15 and lower CMR means higher noise in data integration.

16 Between I.Cox and U.Cox, it appears that the I.Cox method either sub-
17 stantially outperforms U.Cox, or has comparable performance comparing to
18 U.Cox. Specifically, when CR1 is moderate (30%) and MR is high (70%),
19 I.Cox outperforms U.Cox, with more advantage in the MNAR case than in
20 the MCAR case. When CR1 is heavy (60%) with 40% MR, I.Cox outper-
21 forms U.Cox in the cases where CMR is 80% and in the MNAR case with
22 20% CMR; otherwise, it has a close but slightly worse performance than

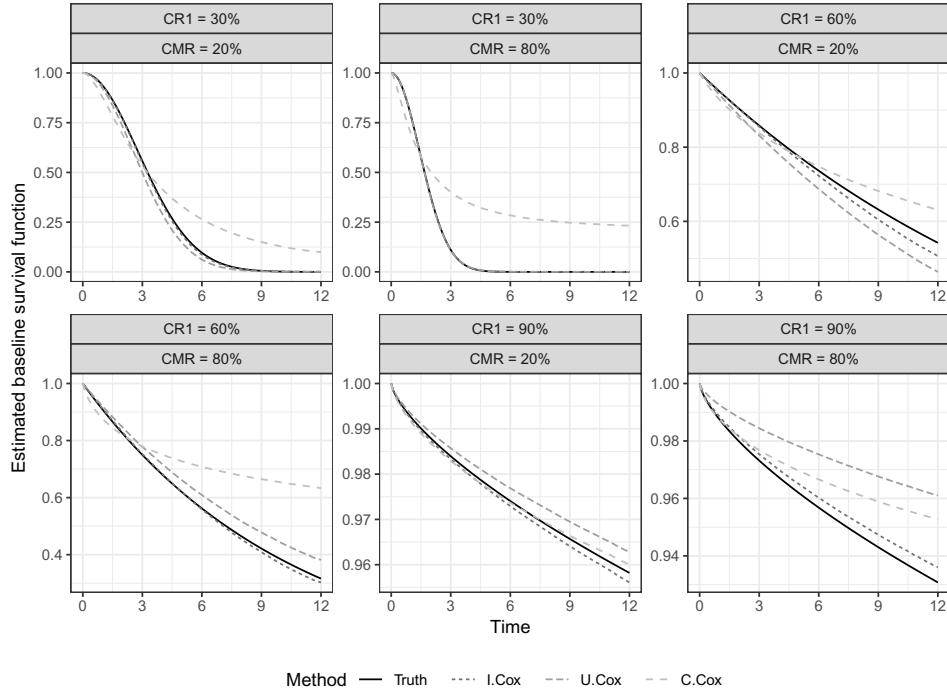


FIG 2. Mean of the estimated baseline survival function in various simulation settings when true record labels are missing at random (MCAR).

1 U.Cox. Lastly, under severe CR1 (90%) and low MR (10%), I.Cox outper-
 2 forms U.Cox when CMR is 80% and has a slightly worse performance when
 3 CMR decreased to 20%. It is not surprising that I.Cox does not always out-
 4 performs U.Cox, because the potential gain from data integration depends
 5 on the quality of both the original data (dataset I) and the matching data
 6 (dataset II). Indeed, I.Cox did not outperform U.Cox in scenario 3 and sce-
 7 nario 5 when CR1 is high and CMR is very low. In general, data integration
 8 is beneficial when the original data misses a substantial amount of true event
 9 records and thus may have inadequate or biased information for model esti-
 10 mation, and/or when the correct information that can be recovered by the
 11 matching data “exceeds” the accompanying noise/false information.

12 Figure 2 presents a visual comparison on the estimation of the baseline
 13 survival function from I.Cox and two naive Cox methods in different set-
 14 tings under MCAR. The true baseline survival curves are included. The
 15 I.Cox clearly performs the best overall, and in most cases the mean of its
 16 baseline survival function estimates over 1,000 replications is close to the

TABLE 3
Mean correct identification rates in percentage for subjects in Case 3 under different simulation settings.

	Scenario					
	1	2	3	4	5	6
MCAR	85.7	89.5	83.1	83.8	80.4	80.0
MNAR	87.5	88.6	90.5	85.5	94.2	86.0

1 corresponding true curve. In contrast, both U.Cox and C.Cox, especially
 2 the latter, may lead to substantial overestimation of the survival probabil-
 3 ities. We have also checked the variation of the estimated survival curves
 4 from these methods and I.Cox performs satisfactorily. See Section 2 of the
 5 Supplementary Materials (Wang et al., 2019) for an example plot of mean
 6 survival curves with pointwise 95% empirical confidence intervals and similar
 7 results under MNAR.

8 Table 3 reports the mean correct identification rate for all the subjects
 9 with uncertain records in Case 3 from the survival analysis with the I.Cox
 10 method. The rate ranges from 80.0% to 94.2%, which means that the true
 11 records can be correctly identifies by the I.Cox model for at least 80% of
 12 subjects having uncertain records in all cases. We remark that in practice
 13 the main focus of such integrative analysis is still on the estimation and
 14 inference of β ; one should be cautious on using the estimated probabilities
 15 to identify the true records, as the empirical evidence from our simulation
 16 study is certainly limited.

17 To check the performance of I.Cox in making inferences about the un-
 18 known covariate coefficients in comparison with U.Cox and O.Cox, we used
 19 bootstrap with 1,000 bootstrap samples. The confidence intervals based on
 20 sample standard deviation and inter-quantile produced estimates in good
 21 agreement, and we report those based on sample standard deviation. The
 22 results of point estimate, SE, and empirical coverage percentage for the coef-
 23 ficient of one continuous covariate and the binary covariate are summarized
 24 in Table 4. The bootstrap SE estimates appear to be close to the empirical
 25 SEs of the coefficient estimates in most of the settings. The coverage rate
 26 of 95% confidence intervals constructed from the SE estimates and normal
 27 approximation is close to the nominal level in most cases. The worst cases
 28 are for β_1 under MNAR when the censoring of dataset I is severe.

29 We have explored the asymptotic behaviors of the I.Cox estimator empir-
 30 ically. Following the original sample size decomposition given in Table 1, we
 31 increase the total sample size to 2, 4, 8, and 16 times under each original
 32 setting. The results show that the mean of $\|\hat{\beta} - \beta_0\|$ decreases as the sample

TABLE 4
Summaries of point estimate, standard error, and empirical coverage of 95% confidence intervals for two covariate coefficients.

#	Method	$\hat{\beta}_1$	SE($\hat{\beta}_1$)	ESE($\hat{\beta}_1$)	CP($\hat{\beta}_1$)	$\hat{\beta}_4$	SE($\hat{\beta}_4$)	ESE($\hat{\beta}_4$)	CP($\hat{\beta}_4$)
MCAR									
1	I.Cox	1.02	0.088	0.089	94.8	1.02	0.171	0.170	94.9
	U.Cox	0.94	0.086	0.088	87.0	0.93	0.162	0.168	92.2
	O.Cox	1.01	0.082	0.085	94.6	1.01	0.156	0.155	94.8
2	I.Cox	1.02	0.088	0.087	94.9	1.02	0.166	0.170	94.3
	U.Cox	1.01	0.095	0.094	95.7	1.01	0.176	0.184	94.0
	O.Cox	1.01	0.084	0.084	95.9	1.01	0.157	0.161	94.7
3	I.Cox	1.03	0.089	0.092	93.0	1.04	0.177	0.180	93.2
	U.Cox	0.96	0.086	0.087	91.4	0.96	0.164	0.163	94.1
	O.Cox	1.01	0.081	0.085	93.7	1.01	0.154	0.154	95.2
4	I.Cox	1.03	0.086	0.088	94.1	1.02	0.170	0.172	95.1
	U.Cox	1.05	0.099	0.096	94.5	1.04	0.192	0.189	95.2
	O.Cox	1.01	0.081	0.081	95.1	1.00	0.157	0.157	95.5
5	I.Cox	1.03	0.082	0.083	92.6	1.01	0.167	0.167	94.7
	U.Cox	1.02	0.084	0.085	94.3	1.03	0.164	0.161	95.3
	O.Cox	1.01	0.077	0.078	94.9	1.01	0.150	0.149	95.4
6	I.Cox	1.04	0.085	0.088	91.7	1.04	0.163	0.170	93.4
	U.Cox	1.09	0.106	0.105	87.7	1.10	0.199	0.196	93.0
	O.Cox	1.00	0.080	0.078	95.5	1.01	0.152	0.154	94.9
MNAR									
1	I.Cox	1.02	0.088	0.088	95.6	1.02	0.172	0.168	95.6
	U.Cox	0.91	0.086	0.086	80.5	0.91	0.162	0.165	90.3
	O.Cox	1.01	0.082	0.085	94.6	1.01	0.156	0.155	94.8
2	I.Cox	1.02	0.088	0.089	94.1	1.02	0.167	0.173	94.7
	U.Cox	0.96	0.092	0.095	92.3	0.95	0.175	0.183	92.7
	O.Cox	1.01	0.084	0.084	95.9	1.01	0.157	0.161	94.7
3	I.Cox	1.04	0.088	0.090	92.4	1.04	0.170	0.169	94.2
	U.Cox	0.95	0.087	0.090	89.2	0.95	0.164	0.163	93.3
	O.Cox	1.01	0.081	0.085	93.7	1.01	0.154	0.154	95.2
4	I.Cox	1.03	0.087	0.089	94.1	1.02	0.170	0.173	95.0
	U.Cox	1.00	0.094	0.095	95.6	1.00	0.192	0.188	95.3
	O.Cox	1.01	0.081	0.081	95.1	1.00	0.157	0.157	95.5
5	I.Cox	1.04	0.082	0.083	91.8	1.04	0.161	0.162	95.1
	U.Cox	1.03	0.084	0.085	93.3	1.04	0.165	0.165	94.8
	O.Cox	1.01	0.077	0.078	94.9	1.01	0.150	0.149	95.4
6	I.Cox	1.06	0.086	0.087	89.9	1.06	0.164	0.175	93.3
	U.Cox	1.08	0.103	0.104	88.2	1.11	0.202	0.209	92.1
	O.Cox	1.00	0.080	0.078	95.5	1.01	0.152	0.154	94.9

SE: Standard error estimate; ESE: Empirical standard error from point estimates;
 CP: Coverage probability (%) of 95% confidence intervals.

1 size increases, and the rate of convergence is approximately the square root
2 of the sample size. We have also done simulation studies with fixed total
3 sample size and the results are similar to what we have presented. More de-
4 tails are available in Section 3 and 4 of the Supplementary Materials ([Wang
5 et al., 2019](#)).

6 **6. Survival Analysis of the Connecticut Data.** We conducted a
7 marginal screening analysis using I.Cox over the aforementioned 58 indica-
8 tors of ICD-9 categories, with three demographic variables, age, male, (ver-
9 sus female) and White (versus non-White) always included in the model.
10 That is, each ICD-9 indicator was included as the fourth variable in the
11 screening process. The inference results were obtained based on 1,000 boot-
12 strap samples, following the procedure detailed in Section 4.3. After the p-
13 values of all the ICD-9 indicators were gathered from the marginal models,
14 the Benjamini–Hochberg procedure ([Benjamini and Hochberg, 1995](#)) was
15 applied to control the false discovery rate (FDR) at 5%. For comparison,
16 we repeated the same analysis using C.Cox, which ignored matching, and
17 U.Cox, which discarded all the uncertain events from matching.

18 The coefficient estimates for male and White from all the marginal models
19 were significant at 5% level. Males were at significantly higher risk of death
20 than females, and whites were at significantly higher risk than non-whites.
21 These findings of disparity in gender and race agree well with existing studies
22 (e.g., [Kung, Pearson and Wei, 2005](#); [Pena et al., 2012](#)). The age effect was less
23 significant compared with gender and race. Most estimates for the coefficient
24 of age from the marginal models were significantly greater than zero at 10%
25 level, providing mild evidence that the survival time after suicide attempt
26 tends to decrease with age for the patients in the study (age 15–30).

27 The screening analysis of ICD-9 codes revealed interesting and insightful
28 results. By controlling the FDR at 5% for the results from each method,
29 neither C.Cox nor U.Cox identified any significant ICD-9 category; in con-
30 trast, I.Cox identified four ICD-9 categories to be significantly associated
31 with the risk of death after unsuccessful suicide attempt. The p-values for
32 coefficient estimates of the four ICD-9 indicators are reported in the upper
33 part of Table 5. The coefficient for ICD-9 code 292 was significantly positive,
34 indicating that patients with drug-induced mental disorder had significantly
35 higher risk than others after controlling for age, gender, and race. Patients
36 with borderline personality disorders (ICD-9 code 301) were also found to
37 have a significantly higher risk of death. These results are supported by
38 several studies, e.g., [Harris and Barraclough \(1997\)](#), [Lieb et al. \(2004\)](#) and
39 [McGirr et al. \(2007\)](#), among others. The I.Cox model also suggests that pa-

TABLE 5

Selected ICD-9 categories by I.Cox and their brief descriptions. Columns 2–4 reports p-values (unadjusted) of coefficient estimates from I.Cox, C.Cox and U.Cox method, respectively, where the significance is indicated by asterisk and the sign of estimates is given in subscripts.

ICD-9	I.Cox	C.Cox	U.Cox	Description
Significant ICD-9 codes under 5% FDR control				
786	0.000* ₊	0.004 ₊	0.002 ₊	Dyspnea, respiratory abnormalities, and chest pain
V45	0.000* ₊	0.088 ₊	0.045 ₊	Postsurgical acquired absence of organ & other post-procedural status
292	0.001* ₊	0.007 ₊	0.007 ₊	Drug-induced mental disorders
301	0.002* ₊	0.069 ₊	0.066 ₊	Borderline personality disorder
Additional ICD-9 codes with individual p-value under 5%				
780	0.010* ₊	0.178 ₊	0.169 ₊	Alteration of consciousness, convulsions, and sleep disturbances
299	0.019* ₊	0.050* ₊	0.035* ₊	Pervasive developmental disorders
298	0.036* ₊	0.075 ₊	0.044* ₊	Other non-organic psychoses
304	0.041* ₊	0.014* ₊	0.011* ₊	Drug dependence (such as opioid type, cocaine, or cannabis)
966	0.041* ₊	0.140 ₊	0.129 ₊	Poisoning by anticonvulsants drugs
E98	0.043* ₋	0.046* ₋	0.065* ₋	Poisoning by analgesics, tranquilizers with undetermined reason
272	0.046* ₊	0.139 ₊	0.094 ₊	Disorders of lipid metabolism
070	0.053 ₊	0.008* ₊	0.008* ₊	Chronic viral hepatitis C
V65	0.143 ₊	0.027* ₊	0.047* ₊	Counseling on substance use and abuse
874	0.338 ₊	0.029* ₊	0.063 ₊	Open wound of neck without mention of complication
969	0.421* ₋	0.027* ₋	0.024* ₋	Poisoning by antidepressants, antipsychotics, and neuroleptics

I.Cox: Integrative Cox model; C.Cox: Regular Cox model fitted to dataset I before matching; U.Cox: Regular Cox model fitted to data with matched records removed.

1 tients with dyspnea respiratory abnormalities and chest pain (ICD-9 code
 2 786) had significantly higher risk. In the literature, chest pain was reported
 3 to have positive association between psychiatric illness and panic disorder
 4 by [Katon et al. \(1988\)](#) and [Fleet et al. \(1996\)](#), respectively, which provided a
 5 possible explanation. Patients having postsurgical acquired absence of organ
 6 and other postprocedural status (ICD-9 code V45) were also under higher
 7 risk of death, which may or may not be directly related to suicide.

8 We also checked the screening results without FDR control. The ad-
 9 ditional ICD-9 codes with unadjusted p-values under 5% are reported in
 10 the lower part of Table 5. For example, the effect of disorders of lipid
 11 metabolism indicated by ICD-9 code 272 was identified by I.Cox. The posi-
 12 tive association between suicidal behavior and lipid metabolism in depressive
 13 disorders was reported by [Koponen et al. \(2015\)](#). Overall, various mental

TABLE 6
Coefficient estimates from joint model including significant ICD-9 categories from marginal screening by I.Cox with FDR controlled at 5%.

Predictor	$\hat{\beta}$	$\exp(\hat{\beta})$	$SE(\hat{\beta})$	z	$\Pr(> z)$
<i>Demographics</i>					
Age	0.12	1.22	0.11	1.11	0.269
Male	1.81	6.11	0.32	5.63	0.000
White	2.18	8.86	0.38	5.78	0.000
<i>ICD-9 Code</i>					
786	1.54	4.67	0.36	4.23	0.000
V45	1.68	5.34	0.57	2.95	0.003
292	0.69	1.98	0.31	2.21	0.027
301	0.60	1.82	0.24	2.48	0.013

1 disorders, psychological issues, drug dependence and abuse appear to be as-
 2 sociated with shortened survival time after unsuccessful suicide attempts.
 3 Therefore, by taking the data uncertainty into consideration and utilizing
 4 information from the second data source, the proposed I.Cox method reveals
 5 much more insightful results than the naive approaches.

6 We then turned our attention to joint modeling, to check the estima-
 7 tion and predictive power of the joint model with all the identified ICD-9
 8 categories. Table 6 summarizes the refitted I.Cox model with the three de-
 9 mographic variables (age, gender, and race) and the four significant ICD-9
 10 indicators identified from marginal screening. The coefficient estimates of
 11 male, White and four ICD-9 indicators were all significantly positive at 5%
 12 level, consistent with the results from screening, while coefficient estimate of
 13 age was insignificant. Because neither of the naive Cox methods suggested
 14 any significant ICD-9 category with FDR controlled at 5% from marginal
 15 screening, their joint models only included the three demographic variables.
 16 We checked that the coefficient estimates were all significant at 5%.

17 For the three joint models resulting from I.Cox and two naive meth-
 18 ods, we performed an out-of-sample comparison analysis on their prediction
 19 performance. (We excluded age in the joint model of I.Cox since it was in-
 20 significant.) Specifically, we randomly split the patients into a training set
 21 and a test set. Patients having events and patients having censoring times
 22 were put in different strata so that the training set and the testing set had
 23 about the same censoring rate. For U.Cox, patients in Case 1 and Case 2
 24 were randomly selected into training set with probability 0.8; for C.Cox,
 25 patients having certain event times (Case 1) and the remaining patients
 26 having censoring times in dataset I were randomly selected into the training
 27 set, separately, with probability 0.8; for I.Cox, patients in Case 1, Case 2,

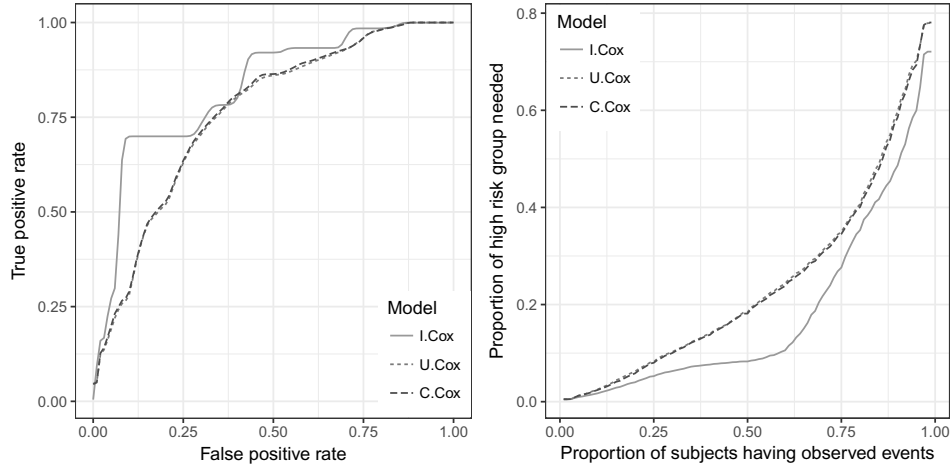


FIG 3. Out-of-sample comparison of the prediction performance on survival outcomes of I.Cox and the naive methods using random splitting.

1 and Case 3 were randomly selected into the training set, respectively, with
 2 probability 0.8, 0.8, and 1. As such, for each method, the testing set only
 3 consisted of patients whose records are certain (Case 1 and Case 2), which
 4 makes an objective evaluation of fitted models possible. In each split, a fitted
 5 model using the training set was used to predict the survival outcomes of
 6 patients in the testing set and classify them to a high risk group and a low
 7 risk group based on their risk scores. By comparing the group classification
 8 to the actual outcomes, we computed the receiver-operator characteristics
 9 (ROC) curve of the survival outcomes. The random split procedure was
 10 repeated 1,000 times and the results were then averaged.

11 Figure 3 presents the ROC curves (on the left panel), and the curves (on
 12 the right panel) showing the relationship between the size of the high risk
 13 group and the proportion of subjects having observed suicide death that
 14 were captured in the high risk group. Here the ROC curves are based on
 15 binary classification using the predicted risk scores; this is motivated by the
 16 clinical setting of a suicide prevention program, where a group of patients
 17 with high risk of suicidal death is identified and subsequently monitored for
 18 suicide prevention. We remark that one may also use a time-dependent ROC
 19 analysis (Heagerty and Zheng, 2005) to quantify the prediction performance
 20 of a survival model. On average, the area under curve (AUC) was 0.825 for
 21 I.Cox, 0.761 for C.Cox, and 0.757 for U.Cox. Therefore, the I.Cox model
 22 provided a better prediction on survival outcomes than both of the naive
 23 methods overall. The results on the right panel converted the ROC curves

1 based on the censoring rate and showed that in order to capture 60% of
2 the patients having observed events, the size of the high risk group needed
3 was 10.6% on average for I.Cox, much less than the sizes 23.8% and 24.3%
4 for C.Cox and U.Cox, respectively. Translating to the real clinical setting,
5 this means that in order to capture 60% of the patients that would die,
6 using I.Cox allows us to achieve this by monitoring only 10.6% of all the
7 patients, while using the native Cox methods will require 25%, a much larger
8 population.

9 **7. Discussion.** We studied a general survival modeling setup with inte-
10 grated data, in which the survival outcome, i.e., the time to certain event of
11 interest, needed to be captured from multiple datasets through record link-
12 age. Such problems are especially prevalent in medical research and health-
13 care analytics. Some commonly encountered events of interest include occur-
14 rence of disease, hospital readmission after discharge, and death following
15 certain diagnostics or treatment. However, patients' medical records are of-
16 ten scattered among many healthcare providers and government agencies.
17 These datasets are generally de-identified to protect patient privacy, but
18 due to limitations in the current healthcare system, the de-identification of
19 each dataset is often done separately before data integration, causing the
20 aforementioned record linkage issues. To the best of our knowledge, build-
21 ing healthcare information exchange system to connect healthcare providers
22 is still largely an ongoing effort. Moreover, analyzing uncertain survival or
23 time-to-event data is challenging due to censoring. When the censoring rate
24 is high, e.g., the event is rare, the information on event times can be quite
25 limited and the results could become sensitive to inaccuracies and anomalies
26 in event times. Therefore, properly handling the uncertainty in event times
27 holds the key to ensure the validity of statistical inference.

28 Data integration with partial identifier is a double-edged sword in inte-
29 grative statistical analysis. As a powerful tool to combine information from
30 multiple sources, integrative analysis with probabilistic uncertainty model-
31 ing needs to be applied with care depending on the degree of imperfectness
32 or noise. Imperfect data integration introduces noise and sometimes errors
33 into the integrated data, the consequence of which could outweigh the po-
34 tential gain in integrative data analysis. Although it is difficult to provide
35 a specific guideline on when to use integrative analysis, we suggest that
36 practitioner always perform out-of-sample analysis to evaluate and compare
37 different methods whenever possible. To ensure the evaluation is objective,
38 only the data without uncertainty should be used in testing.

39 Our case study has an additional distinguishing feature in that it is the

1 outcome variable (survival time) that is obtained from data integration. This
 2 is in contrast to other integrative data analysis settings where usually pre-
 3 dictors or features are obtained from multiple data sources. In our applica-
 4 tion, we obtained insightful results on potential risk factors associated with
 5 death following suicide attempt, which otherwise would have been missed
 6 by the naive approaches. Compared with the method of [Snapinn \(1998\)](#), our
 7 method is more attractive in that it does not require additional diagnostic
 8 variables or prior knowledge on the characterization of the truth indicators.

9 Several directions are worth pursuing for future research. The standard
 10 errors of the estimates cannot be easily produced along with the proposed
 11 estimation procedure. Although bootstrap is shown to perform well, the
 12 method would be more attractive in practice if a less computationally inten-
 13 sive inference approach were available. Under realistic settings of imperfect
 14 data linkage, the proposed method is shown to outperform several naive ap-
 15 proaches. A natural theoretical question of interest is to quantify how the
 16 potential gain from data integration is associated with the quality of the
 17 original data and the match data. Our model framework is flexible and can
 18 be further extended to other survival models such as parametric survival
 19 models and competing risk models. Other extensions include the modeling
 20 of censoring times with covariates and the incorporation of certain known
 21 missing mechanism of the label of true endpoint. In our application, we
 22 adopted a marginal screening approach to identify important predictors; it
 23 would be interesting to extend the proposed method to conduct variable
 24 selection with high-dimensional predictors through regularized estimation.
 25 The rareness of suicide attempt brings many challenges in its modeling and
 26 prediction, including the occurrence of quasi-complete separation; these is-
 27 sues will need to be carefully studied in the future.

28 It is promising to further explore the trimmed likelihood formulation to
 29 better understand the robustness of the proposed approach and design better
 30 algorithm to target its global optimal solution. This formulation also sheds
 31 light on the consistency of the resulting estimator of the proposed method
 32 through the perspective of robust estimation and outlier detection. It shows
 33 that at least two conditions, regarding the proportion and magnitude of the
 34 “outliers” — fake records — are required. First, the proportion of patients
 35 with uncertain records should be under control, e.g., $(n_2 + n_3)/n \rightarrow c$ for
 36 some $0 \leq c < 1$ as $n \rightarrow \infty$. Second, the fake records have to be distinguish-
 37 able from the true one; e.g., for patient j , we need $k^* = \arg \max_k r_{j,k}(\beta^*)$ for
 38 n sufficiently large, where the k^* th record is the truth and β^* denotes the
 39 true coefficient vector. A thorough investigation of the theoretical properties
 40 of the proposed method along this direction is of interest.

SUPPLEMENTARY MATERIAL

1 **Supplementary Materials: Integrative Survival Analysis with**
2 **Uncertain Event Times in Application to a Suicide Risk Study**
3 (doi: [COMPLETED BY THE TYPESETTER](#); .pdf). We provide detailed
4 derivations of the likelihood formulation, additional supporting tables/figures
5 from simulation studies, and discussions on the properties of the proposed
6 method.

7 **Acknowledgements.** Chen's research was partially supported by Na-
8 tional Science Foundation grants DMS-1613295 and IIS-1718798, and Na-
9 tional Institutes of Health grant R01-MH112148. Aseltine's research was
10 partially supported by National Institutes of Health grant R01-MH112148.

11 **References.**

- 12 ALBERT, A. and ANDERSON, J. A. (1984). On the Existence of Maximum Likelihood
13 Estimates in Logistic Regression Models. *Biometrika* **71** 1–10.
- 14 BENJAMINI, Y. and HOCHBERG, Y. (1995). Controlling the False Discovery Rate: A Prac-
15 tical and Powerful Approach to Multiple Testing. *Journal of the royal statistical society.*
16 *Series B (Methodological)* 289–300.
- 17 BOHENSKY, M. A., JOLLEY, D., SUNDARARAJAN, V., EVANS, S., PILCHER, D. V.,
18 SCOTT, I. and BRAND, C. A. (2010). Data Linkage: A powerful research tool with
19 potential problems. *BMC Health Services Research* **10** 346.
- 20 BOSTWICK, M. J., PABBATI, C., GESKE, J. R. and MCKEAN, A. J. (2015). Suicide
21 Attempt as a Risk Factor for Completed Suicide: Even More Lethal Than We Knew.
22 *The American Journal of Psychiatry* **173** 1094–1100.
- 23 BRESLOW, N. (1974). Covariance Analysis of Censored Survival Data. *Biometrics* **30** 89–
24 99.
- 25 CHEN, K. and ASELTINE, R. (2017). Using Hospitalization and Mortality Data to Target
26 Suicide Prevention Activities: A Demonstration from Connecticut. *Journal of Adoles-*
27 *cent Health* **61** 192–197.
- 28 COX, D. R. (1972). Regression Models and Life-Tables. *Journal of the Royal Statistical*
29 *Society. Series B (Methodological)* **34** 187–220.
- 30 COX, D. R. (1975). Partial Likelihood. *Biometrika* **62** 269–276.
- 31 DEMPSTER, A. P., LAIRD, N. M. and RUBIN, D. B. (1977). Maximum Likelihood from
32 Incomplete Data via the EM Algorithm. *Journal of the Royal Statistical Society. Series*
33 *B (Methodological)* **39** 1–38.
- 34 EFRON, B. (1979). Bootstrap Methods: Another Look at the Jackknife. *The Annals of*
35 *Statistics* **7** 1–26.
- 36 EFRON, B. (1981). Censored Data and the Bootstrap. *Journal of the American Statistical*
37 *Association* **76** 312–319.
- 38 FLEET, R. P., DUPUIS, G., MARCHAND, A., BURELLE, D., ARSENAULT, A. and BEIT-
39 MAN, B. D. (1996). Panic Disorder in Emergency Department Chest Pain Patients:
40 Prevalence, Comorbidity, Suicidal Ideation, and Physician Recognition. *The American*
41 *Journal of Medicine* **101** 371–380.
- 42 HADI, A. S. and LUCEÑO, A. (1997). Maximum Trimmed Likelihood Estimators: A Uni-
43 fied Approach, Examples, and Algorithms. *Computational Statistics & Data Analysis*
44 **25** 251–272.

- 1 HARRIS, E. C. and BARRACLOUGH, B. (1997). Suicide as An Outcome for Mental Disor-
 2 ders. A Meta-Analysis. *The British Journal of Psychiatry* **170** 205–228.
- 3 HARRON, K., GOLDSTEIN, H. and DIBBEN, C. (2015). *Methodological Developments in*
 4 *Data Linkage*. John Wiley & Sons.
- 5 HEAGERTY, P. J. and ZHENG, Y. (2005). Survival Model Predictive Accuracy and ROC
 6 Curves. *Biometrics* **61** 92–105.
- 7 HOF, M. H. P. and ZWINDERMAN, A. H. (2012). Methods for Analyzing Data from
 8 Probabilistic Linkage Strategies Based on Partially Identifying Variables. *Statistics in*
 9 *Medicine* **31** 4231–4242.
- 10 HOF, M. H. P. and ZWINDERMAN, A. H. (2015). A Mixture Model for the Analysis of
 11 Data Derived from Record Linkage. *Statistics in Medicine* **34** 74–92.
- 12 JAMSHIDIAN, M. and JENNRICH, R. I. (2000). Standard Errors for EM Estimation. *Journal*
 13 *of the Royal Statistical Society. Series B (Statistical Methodology)* **62** 257–270.
- 14 KALBFLEISCH, J. D. and PRENTICE, R. L. (2002). *The Statistical Analysis of Failure*
 15 *Time Data*. John Wiley & Sons.
- 16 KATON, W., HALL, M. L., RUSSO, J., CORMIER, L., HOLLIFIELD, M., VITALIANO, P. P.
 17 and BEITMAN, B. D. (1988). Chest Pain: Relationship of Psychiatric Illness to Coronary
 18 Arteriographic Results. *The American Journal of Medicine* **84** 1–9.
- 19 KOPONEN, H., KAUTIAINEN, H., LEPPÄNEN, E., MÄNTYSELKÄ, P. and VANHALA, M.
 20 (2015). Association Between Suicidal Behaviour and Impaired Glucose Metabolism in
 21 Depressive Disorders. *BMC Psychiatry* **15** 163.
- 22 KUNG, H.-C., PEARSON, J. L. and WEI, R. (2005). Substance Use, Firearm Availability,
 23 Depressive Symptoms, and Mental Health Service Utilization among White and African
 24 American Suicide Decedents Aged 15 to 64 Years. *Annals of Epidemiology* **15** 614–621.
- 25 LIEB, K., ZANARINI, M. C., SCHMAHL, C., LINEHAN, M. M. and BOHUS, M. (2004).
 26 Borderline Personality Disorder. *The Lancet* **364** 453–461.
- 27 MCGIRR, A., PARIS, J., LESAGE, A., RENAUD, J. and TURECKI, G. (2007). Risk Factors
 28 for Suicide Completion in Borderline Personality Disorder: A Case-Control Study of
 29 Cluster B Comorbidity and Impulsive Aggression. *The Journal of Clinical Psychiatry*
 30 **68** 721–729.
- 31 MEIER, A. S., RICHARDSON, B. A. and HUGHES, J. P. (2003). Discrete Proportional
 32 Hazards Models for Mismeasured Outcomes. *Biometrics* **59** 947–954.
- 33 MENG, X.-L. and RUBIN, D. B. (1991). Using EM to Obtain Asymptotic Variance-
 34 Covariance Matrices: The SEM Algorithm. *Journal of the American Statistical As-*
 35 *sociation* **86** 899–909.
- 36 MENG, X.-L. and RUBIN, D. B. (1993). Maximum Likelihood Estimation via the ECM
 37 Algorithm: A General Framework. *Biometrika* **80** 267–278.
- 38 MURPHY, S. A. and VAN DER VAART, A. W. (2000). On Profile Likelihood. *Journal of*
 39 *the American Statistical Association* **95** 449–465.
- 40 NADARAJAH, S. and KOTZ, S. (2006). R Programs for Truncated Distributions. *Journal*
 41 *of Statistical Software* **16** 1–8.
- 42 NEYKOV, N., FILZMOSER, P., DIMOVA, R. and NEYTCHEV, P. (2007). Robust Fitting of
 43 Mixtures Using the Trimmed Likelihood Estimator. *Computational Statistics & Data*
 44 *Analysis* **52** 299–308.
- 45 PATRICK, A. R., MILLER, M., BARBER, C. W., WANG, P. S., CANNING, C. F. and
 46 SCHNEEWEISS, S. (2010). Identification of Hospitalizations for Intentional Self-Harm
 47 When E-codes are Incompletely Recorded. *Pharmacoepidemiology and Drug Safety* **19**
 48 1263–1275.
- 49 PENA, J. B., MATTHIEU, M. M., ZAYAS, L. H., MASYN, K. E. and CAINE, E. D. (2012).
 50 Co-occurring Risk Behaviors Among White, Black, and Hispanic US High School Ado-

- 1 lescents with Suicide Attempts Requiring Medical Attention, 1999–2007: Implications
2 for Future Prevention Initiatives. *Social Psychiatry and Psychiatric Epidemiology* **47**
3 29–42.
- 4 PRITCHARD, C. and HANSEN, L. (2015). Examining Undetermined and Accidental Deaths
5 as Source of ‘Under-Reported-Suicide’ by Age and Sex in Twenty Western Countries.
6 *Community Mental Health Journal* **51** 365–376.
- 7 RICHARDSON, B. A. and HUGHES, J. P. (2000). Product Limit Estimation for Infectious
8 Disease Data When the Diagnostic Test for the Outcome is Measured with Uncertainty.
9 *Biostatistics* **1** 341–354.
- 10 ROUSSEEUW, P. J. (1984). Least Median of Squares Regression. *Journal of the American*
11 *Statistical Association* **79** 871–880.
- 12 RUBIN, D. B. (1987). *Multiple Imputation for Nonresponse in Surveys*. John Wiley &
13 Sons.
- 14 SNAPINN, S. M. (1998). Survival Analysis with Uncertain Endpoints. *Biometrics* **54** 209–
15 218.
- 16 SUOMINEN, K., ISOMETSA, E., SUOKAS, J., HAUKKA, J., ACHTE, K. and LÖNNQVIST, J.
17 (2004). Completed Suicide After a Suicide Attempt: A 37-Year Follow-Up Study. *Ameri-*
18 *cans Journal of Psychiatry* **161** 562–563.
- 19 TANCREDI, A. and LISEO, B. (2015). Regression Analysis with Linked Data: Problems
20 and Possible Solutions. *Statistica* **75** 19–35.
- 21 R DEVELOPMENT CORE TEAM (2017). R: A Language and Environment for Statistical
22 Computing R Foundation for Statistical Computing, Vienna, Austria ISBN 3-900051-
23 07-0.
- 24 TØLLEFSEN, I. M., THIBLIN, I., HELWEG-LARSEN, K., HEM, E., KASTRUP, M., NY-
25 BERG, U., ROGDE, S., ZAHL, P.-H., ØSTEVOLD, G. and EKEBERG, Ø. (2016). Accidents
26 and Undetermined Deaths: Re-evaluation of Nationwide Samples from the Scandinavian
27 Countries. *BMC Public Health* **16** 449.
- 28 WANG, W., ASELTINE, R., CHEN, K. and YAN, J. (2019). Supplementary Materials to
29 “Integrative Survival Analysis with Uncertain Event Times in Application to a Suicide
30 Risk Study”. *Annals of Applied Statistics*.
- 31 WINGLEE, M., VALLIANT, R. and SCHEUREN, F. (2005). A Case Study in Record Linkage.
32 *Survey Methodology* **31** 3–11.
- 33 XU, C., BAINES, P. D. and WANG, J. L. (2014). Standard Error Estimation Using the
34 EM Algorithm for the Joint Modeling of Survival and Longitudinal Data. *Biostatistics*
35 **15** 731–744.
- 36 ZHAO, Q., SHI, X., XIE, Y., HUANG, J., SHIA, B. and MA, S. (2015). Combining Multidi-
37 mensional Genomic Measurements for Predicting Cancer Prognosis: Observations from
38 TCGA. *Briefings in Bioinformatics* **16** 291–303.

39 WENJIE WANG
KUN CHEN
JUN YAN
DEPARTMENT OF STATISTICS
UNIVERSITY OF CONNECTICUT
STORRS, CT 06269
CENTER FOR POPULATION HEALTH
UConn Health
FARMINGTON, CT 06032
E-MAIL: wenjie.2.wang@uconn.edu
kun.chen@uconn.edu
jun.yan@uconn.edu

ROBERT ASELTINE
DIVISION OF BEHAVIORAL SCIENCE AND COMMUNITY HEALTH
CENTER FOR POPULATION HEALTH
UConn Health
FARMINGTON, CT 06032
E-MAIL: aseltine@uchc.edu