

Attributable Risk Function with Clustered Survival Data

Changchun Xie^{1,2*}, Xuewen Lu³ and Janice Pogue²

¹Department of Environmental Health, University of Cincinnati, Ohio, USA

²Department of Clinical Epidemiology and Biostatistics, McMaster University, Ontario, Canada

³Department of Mathematics and Statistics, University of Calgary, Alberta, Canada

Abstract

The Attributable Fraction or risk function (ARF) is used to measure the impact of an exposure on occurrence of disease within a population. For any prospective cohort study, risk is likely to be estimated using time to event or survival data. Attributable risk function with right censored survival data has been discussed by Samuelsen and Eide. We propose a natural extension of the ARF to clustered survival data, which are common in medical research. We derive an estimator of the ARF. Simulation studies are conducted to evaluate the performance of our method and investigate the consequences of ignoring the cluster effect in analysis.

Keywords: Attributable risk function; Frailty models; Proportional hazards model

Introduction

The proportional hazards models are often used in prospective clinical and epidemiological studies to evaluate the association between time to a disease and exposures or risk factors [1]. These models allow the risk of outcomes over time to be estimated in the presence of censoring, and can incorporate time-dependent covariates and clustering of the individuals observed. The risk or hazard function given a risk factor covariate Z can be written as

$$\lambda(t | Z) = \lambda_0(t) \exp(\beta'Z), \quad (1)$$

where β is the regression parameter and $\lambda_0(t)$ is an unspecified baseline hazard function. This association, measured by the hazard ratio, does not take into account the prevalence of the risk factors in a given population. The attributable risk function has been used to measure the proportion of disease in a population associated with a given risk factor for binary outcomes [2]. When the outcomes are binary, the population attributable risk is usually defined as [3]

$$\phi = \frac{P(D = 1) - P(D = 1 | Z = 0)}{P(D = 1)}, \quad (2)$$

where D denotes a binary outcome and Z denotes the binary risk factor. For the time-to-event outcome T , A natural extension of ϕ for T is, for some $t > 0$, [4]

$$\tilde{\phi}(t) = \frac{P(T \leq t) - P(T \leq t | Z = 0)}{P(T \leq t)}. \quad (3)$$

Chen et al. [4] proposed an alternative measure of the attributable risk function for T :

$$\phi(t) = \frac{\lambda(t) - \lambda(t | z = 0)}{\lambda(t)}, \quad (4)$$

where $\lambda(t)$ is the population hazard function (see appendix for detail). Not like population attributable risk function for the binary outcomes, the attributable risk function for the time-to-event endpoints is not necessarily constant over time, even when the baseline hazard function itself and the exposure prevalence are constant [4]. However, Chen et al. only considered the case with one covariate and intervention at time 0. The adjusted attributable risk function cannot be obtained from their definition. In a very recent article, Samuelsen and Eide [5] proposed ARFs for studies with covariates and interventions that may vary over time:

$$\phi(t) = \frac{E[\lambda(t | Z)] - E[\lambda(t | Z^*)]}{E[\lambda(t | Z)]}, \quad (5)$$

where $Z = (Z_1, Z_2, Z_3, \dots, Z_p)$ and $Z^* = (Z_1^*, Z_2^*, \dots, Z_p^*)$ are the covariates without and with intervention respectively, $E[\lambda(t | Z)]$ and $E[\lambda(t | Z^*)]$ are their respective expected hazard function over the two populations. Since a vector of covariates is considered in the definition of ARF, the adjusted ARF can be calculated from this definition as Samuelsen and Eide showed in their examples. Chen et al. [4] have clearly shown the difference between the definition of ARF in (3) and that in (4). However, the relationship between the definition of ARF (4) proposed by Chen et al. [4] and that (5) provided by Samuelsen and Eide [5] is not that clear. We have shown that these two definitions are different if we assume the distribution of Z does not change with t , but they are the same when we consider distribution of Z among those at risk (i.e. condition on $T \geq t$), which typically changes over time (see appendix for detail). So far, these methods do not consider the case of ARF for clustered survival data, which are common in medical research with cluster randomized trials or community-based or family-based prospective cohort studies. When the cluster effects are sufficiently large, ignoring the clustering can lead to substantially biased estimators of regression coefficients [6], and leads to biased estimates of ARFs. In this paper, we extend the concept of ARF to clustered survival data. By using the gamma frailty model, which is a popular tool for addressing cluster effects in clustered survival data [7], our approach provides a practical method to calculate ARF in the presence of dependence in survival data due to cluster effects.

Attributable Risk Function for Clustered Survival Data

Assume n is the total number of subjects in a study with K clusters and n_k subjects in cluster k ($k = 1, \dots, K$), so that $n = \sum_{k=1}^K n_k$. Let T_{ki} be the

*Corresponding author: Changchun Xie, PhD, Department of Environmental Health, University of Cincinnati, Ohio, USA, E-mail: xiecn@UCMAIL.UC.EDU

Received January 04, 2012; Accepted March 23, 2012; Published March 23, 2012

Citation: Xie C, Lu X, Pogue J (2012) Attributable Risk Function with Clustered Survival Data. J Biomet Biostat S1:007. doi:10.4172/2155-6180.S1-007

Copyright: © 2012 Xie C, et al. This is an open-access article distributed under the terms of the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

survival times, L_{ki} be the censoring times and Z_{ki} be the P -dimensional vector of covariates for subject i ($i = 1, \dots, n_k$) in cluster ($k = 1, \dots, K$). Let T_k denote the vector $(T_{k1}, \dots, T_{kn_k})$ with L_k and Z_k defined similarly. Suppose that T_k, L_k and Z_k are independent across clusters and (T_k, L_k) ($k = 1, \dots, K$) are *i.i.d* with the components of T_k and L_k conditionally independent given covariates Z_k . Let $X_{ki} = \min(T_{ki}, L_{ki}), \Delta_{ki} = I(T_{ki} \leq L_{ki})$ and the at-risk indicator $Y_{ki}(t) = I(X_{ki} \geq t)$. We also suppose

$\max_k \sum_{i=1}^{n_k} \Delta_{ki} > 1$. The model we consider has the form

$$\lambda_{ki}(t) = u_k \lambda_0(t) \exp(\beta' Z_{ki}), \tag{6}$$

where the u_k is the common risk factor for all subjects in cluster k . We present two approaches to modelling cluster effects.

The fixed effects model

When the number of clusters, K , is small, the fixed effects model can be used by including indicator variables for clusters. Arbitrarily setting one cluster as the reference cluster, for example, cluster 1, we obtain

$$\lambda_{ki}(t) = \lambda_0(t) \exp(\alpha_k + \beta' Z_{ki}) \tag{7}$$

for $k = 1, \dots, K$, with $\alpha_1 = 0$. Using the same method suggested by Samuelsen and Eide [5], the attributable fractions can be calculated by treating the clusters as $(K - 1)$ dimensional vector of covariates. However, when K is large, compared to the sample size, and therefore there are too many parameters in the model, the asymptotics break down since $K \rightarrow \infty$ as $n \rightarrow \infty$ [7]. This approach is well-known to cause bias in parameter estimates.

The frailty model

The frailty model does not treat the cluster effects as parameters, but treats them as a sample from a frailty distribution. In this paper, we consider gamma frailty with mean 1 and variance θ . The density is given by

$$f(u) = \frac{u^{1/\theta-1} \exp(-u/\theta)}{\Gamma(1/\theta)\theta^{1/\theta}}. \tag{8}$$

Under the gamma frailty model, the marginal hazard function can be obtained from the expectation of the hazard function [8], conditional on being at risk at t and covariate Z :

$$\begin{aligned} \mu(t | Z) &= E(U | T \geq t, Z) \lambda_0(t) \exp(\beta' Z) \\ &= \{1 + \theta \Lambda_0(t) \exp(\beta' Z)\}^{-1} \lambda_0(t) \exp(\beta' Z), \end{aligned}$$

where $\Lambda_0(t) = \int_0^t \lambda_0(s) ds$. Note the average frailty value, $\{1 + \theta \Lambda_0(t) \exp(\beta' Z)\}^{-1}$ is a decreasing function of time, which is due to the fact that the subjects with high frailty values experience the event earlier on average and the population will contain more and more subjects with low frailty values. From this, the ARF is defined as

$$\phi(t) = \frac{E[\mu(t | Z)] - E[\mu(t | Z^*)]}{E[\mu(t | Z)]}. \tag{9}$$

Compared to the fixed effects model, the frailty model has the advantage of parsimony. The number of parameters to describe cluster effects does not increase with the number of clusters. In the gamma frailty, we have used only one parameter θ to describe the cluster effects (the heterogeneity of cluster effects). As θ increases, frailties become more dispersed and dependence increases [7].

Estimation Procedure

There are many methods available to fit the semiparametric gamma frailty model [8-11]. Assume we have the estimates $\hat{\theta}, \hat{\beta}$ and $\hat{\Lambda}_0(t)$ of θ, β and $\Lambda_0(t)$ respectively. We consider the following two situations:

The distributions of Z and Z^* do not depend on time t

The estimated ARF can be represented as

$$\hat{\phi}(t) = 1 - \frac{E\{[1 + \hat{\theta}\hat{\Lambda}_0(t)\exp(\hat{\beta}'Z^*)]^{-1}\exp(\hat{\beta}'Z^*)\}}{E\{[1 + \hat{\theta}\hat{\Lambda}_0(t)\exp(\hat{\beta}'Z)]^{-1}\exp(\hat{\beta}'Z)\}}. \tag{10}$$

When the covariate is one binary variable such that $P(Z = 1) = p, P(Z = 0) = 1 - p$ and $P(Z^* = 0) = 1$, the estimator simplifies to

$$\hat{\phi}(t) = 1 - \frac{\{1 + \hat{\theta}\hat{\Lambda}_0(t)\}^{-1}}{p\{1 + \hat{\theta}\hat{\Lambda}_0(t)\exp(\hat{\beta})\}^{-1}\exp(\hat{\beta}) + (1 - p)\{1 + \hat{\theta}\hat{\Lambda}_0(t)\}^{-1}}. \tag{11}$$

Although the distributions of Z and Z^* do not depend on time t , this function changes over time. When the effect of frailty does not exist, that is $\hat{\theta} = 0$, it becomes time-independent.

The distributions of Z and Z^* change over time t

In survival data, the distribution of some covariates for the subjects at risk usually changes over time. In order to get the estimate of $\phi(t)$, we need to estimate the population survival function conditional on covariate, Z :

$$\hat{S}(t | Z) = \{1 + \hat{\theta}\hat{\Lambda}_0(t)\exp(\hat{\beta}'Z)\}^{-1/\hat{\theta}}. \tag{12}$$

Following Samuelsen and Eide's [5] approach, the population hazard in a finite population with n individuals and covariates $z_i, i = 1, 2, \dots, n$ can be estimated by

$$\hat{\mu}(t) = \frac{\sum_{i=1}^n \{[1 + \hat{\theta}\hat{\Lambda}_0(t)\exp(\hat{\beta}'z_i)]^{-1}\hat{\lambda}_0(t)\exp(\hat{\beta}'z_i)\hat{S}(t | Z_i)\}}{\sum_{i=1}^n \hat{S}(t | Z_i)}. \tag{13}$$

The population hazard in a finite population with covariates $z_i^*, i = 1, 2, \dots, n$ can be estimated by

$$\hat{\mu}^*(t) = \frac{\sum_{i=1}^n \{[1 + \hat{\theta}\hat{\Lambda}_0(t)\exp(\hat{\beta}'z_i^*)]^{-1}\hat{\lambda}_0(t)\exp(\hat{\beta}'z_i^*)\hat{S}(t | Z_i^*)\}}{\sum_{i=1}^n \hat{S}(t | Z_i^*)}. \tag{14}$$

The estimate of $\phi(t)$ can be calculated as

$$\hat{\phi}(t) = 1 - \frac{\hat{\mu}^*(t)}{\hat{\mu}(t)} = 1 - \frac{\sum_{i=1}^n \{[1 + \hat{\theta}\hat{\Lambda}_0(t)\exp(\hat{\beta}'z_i^*)]^{-1}\exp(\hat{\beta}'z_i^*)\hat{S}(t | Z_i^*)\} \sum_{i=1}^n \hat{S}(t | Z_i)}{\sum_{i=1}^n \{[1 + \hat{\theta}\hat{\Lambda}_0(t)\exp(\hat{\beta}'z_i)]^{-1}\exp(\hat{\beta}'z_i)\hat{S}(t | Z_i)\} \sum_{i=1}^n \hat{S}(t | Z_i^*)}. \tag{15}$$

Simulations

In this section, we evaluate the performance of our method and investigate the consequences of ignoring the cluster effect in analysis. The time-to-event data were generated according to the model (6), where there are 200 clusters with cluster size 10; $\beta = 0, \log(2)$ respectively; u had gamma frailty with mean 1 and variance $\theta = 0.5, 1.64$, which corresponds to Kendall's $\tau = 0.2, 0.45$ respectively; baseline functions were constant of 0.01, 0.1 and 1 respectively. Each individual's binary exposure indicator was generated from Bernoulli distribution with $p = 0.25$. Censoring times were from Uniform distribution $U(0, L)$ where L was chosen to get about 10% and 30% of censored observations respectively. For each case described above, 1000 simulated data were generated. Following Chen et al. [4] we calculated the estimates and their associated variances at the 75 percentile and median of the marginal survival distribution, t_1 and t_2 , respectively. The results are

$\lambda_0(t) \equiv \lambda_0$	θ	Cens. %	$t_1 : S(t_1) = 0.75$				$t_2 : S(t_2) = 0.5$			
			Gamma Frailty		Ignoring Frailty		Gamma Frailty		Ignoring Frailty	
			Bias	SE	Bias	SE	Bias	SE	Bias	SE
0.01	1.64	10%	.0003	.0099	.0001	.0137	.0001	.0054	.0003	.0136
0.01	1.64	30%	.0002	.0112	.0007	.0157	.0001	.0062	.0001	.0155
0.01	0.5	10%	.0001	.0130	.0001	.0138	.0006	.0107	.0008	.0136
0.01	0.5	30%	.0001	.0151	.0000	.0155	.0002	.0125	.0004	.0159
0.1	1.64	10%	.0004	.0099	.0007	.0140	.0004	.0054	.0002	.0140
0.1	1.64	30%	.0003	.0114	.0005	.0156	.0004	.0061	.0012	.0156
0.1	0.5	10%	.0005	.0133	.0006	.0140	.0000	.0105	.0007	.0131
0.1	0.5	30%	.0004	.0157	.0001	.0162	.0002	.0127	.0000	.0159
1.0	1.64	10%	.0000	.0103	.0004	.0135	.0005	.0053	.0005	.0134
1.0	1.64	30%	.0004	.0116	.0003	.0157	.0002	.0062	.0002	.0147
1.0	0.5	10%	.0003	.0131	.0004	.0141	.0002	.0110	.0003	.0135
1.0	0.5	30%	.0009	.0148	.0007	.0159	.0002	.0134	.0000	.0167

Table 1: Estimation of attributable fractions in simulated clustered survival data: $\beta = 0$.

$\lambda_0(t) \equiv \lambda_0$	θ	Cens. %	$t_1 : S(t_1) = 0.75$				$t_2 : S(t_2) = 0.5$			
			Gamma Frailty		Ignoring Frailty		Gamma Frailty		Ignoring Frailty	
			Bias	SE	Bias	SE	Bias	SE	Bias	SE
0.01	1.64	10%	.0089	.0092	.0251	.0148	.0054	.0043	.0250	.0119
0.01	1.64	30%	.0087	.0094	.0131	.0150	.0055	.0044	.0339	.0116
0.01	0.5	10%	.0075	.0116	.0219	.0145	.0103	.0064	.0077	.0096
0.01	0.5	30%	.0075	.0126	.0140	.0153	.0099	.0071	.0125	.0102
0.1	1.64	10%	.0086	.0093	.0249	.0148	.0052	.0044	.0249	.0117
0.1	1.64	30%	.0086	.0097	.0119	.0153	.0054	.0045	.0346	.0123
0.1	0.5	10%	.0076	.0118	.0223	.0148	.0098	.0066	.0079	.0095
0.1	0.5	30%	.0075	.0129	.0141	.0158	.0094	.0067	.0134	.0097
1.0	1.64	10%	.0091	.0095	.0249	.0143	.0054	.0043	.0245	.0118
1.0	1.64	30%	.0092	.0098	.0122	.0158	.0052	.0047	.0354	.0122
1.0	0.5	10%	.0075	.0118	.0221	.0152	.0101	.0064	.0077	.0097
1.0	0.5	30%	.0076	.0129	.0144	.0152	.0099	.0067	.0130	.0098

Table 2: Estimation of attributable fractions in simulated clustered survival data: $\beta = \log(2)$.

shown in Table 1 and 2. Here the bias is the absolute difference between the average of the 1000 estimates and the true attributable fraction and SE is the sample standard error. As shown in Table 1, where $\beta=0$, SE increases when the frailty is ignored. However, in Table 2, where $\beta = \log(2)$, both bias and SE increase when the frailty is ignored, especially for large θ .

Discussion

Clustered survival data often occur in many practical areas, where a group of related subjects constitutes a cluster, such as a group of patients from the same hospital, a group of students from the same school, a group of people from the same community or a group of genetically related members from the same family. The attributable fraction with clustered survival data is discussed in this article as a measure of the proportion of disease over time with associated risk factors within given populations. Our simulations show that ignoring cluster effects can cause the increase of both bias and the sample standard error of the attributable fraction, especially for large cluster effects.

In this paper, we only consider gamma frailty. However, gamma frailty can be easily replaced with other frailty distributions which have a simple Laplace transform representation, such as the inverse Gaussian distribution, the positive stable distribution, or others. We focus on ARF when an intervention takes place at time 0. However, we can also consider ARF when an intervention takes place at time t and use $\varphi_A(t)$ to denote it. Under the proportional hazard model conditional

on frailty U the estimate of $\varphi_A(t)$ can be expressed as

$$\hat{\varphi}_A(t) = 1 - \frac{\sum_{i \in R(t)} \{1 + \hat{\theta} \hat{\lambda}_0(t) \exp(\hat{\beta}' z_i^*)\}^{-1} \exp(\hat{\beta}' z_i^*)}{\sum_{i \in R(t)} \{1 + \hat{\theta} \hat{\lambda}_0(t) \exp(\hat{\beta}' z_i)\}^{-1} \exp(\hat{\beta}' z_i)}, \quad (16)$$

where $R(t)$ is the risk set at time t i.e. the set of subjects with $X_i \geq t$.

In this paper, our interest is in the modelling cluster effect on ARF. To make inference for ARF, one may use the bootstrap method [12] to find the standard errors of the estimates and construct confidence intervals. Applications of the proposed methods to a real data analysis and associated inference issues will be investigated in our future research.

References

1. Cox DR (1972) Regression models and life-tables. J R Statist Soc B 34: 187-220.
2. Walter SD (1976) The estimation and interpretation of attributable risk in health research. Biometrics 32: 829-849.
3. Levin ML (1953) The occurrence of lung cancer in man. Acta Unio Int Contra Cancrum 9: 531-541.
4. Chen YQ, Hu C, Wang Y (2006) Attributable risk function in the proportional hazards model for censored time-to-event. Biostatistics 7: 515-529.
5. Samuelsen SO, Eide GE (2008) Attributable fractions with survival data. Stat Med 27: 1447-1467.
6. Henderson R, Oman P (1999) Effect of frailty on marginal regression estimates in survival analysis. J R Statist Soc B 61: 367-379.

7. Glidden DV, Vittinghoff E (2004) Modelling clustered survival data from multicentre clinical trials. *Stat Med* 23: 369-388.
8. Nielsen GG, Gill RD, Andersen PK, Sorensen TIA (1992) A counting process approach to maximum likelihood estimation in frailty models. *Scand J Statist* 19: 25-43.
9. Klein JP (1992) Semiparametric estimation of random effects using the Cox model based on the EM algorithm. *Biometrics* 48: 795-806.
10. Therneau TM, Grambsch PM (2000) Modelling survival data: extending the Cox model. Springer, New York.
11. Ha ID, Lee Y, Song JK (2001) Hierarchical likelihood approach for frailty models. *Biometrika* 88: 233.
12. Efron B, Tibshirani RJ (1993) An introduction to the bootstrap. Chapman and Hall.

Appendix

Assume the covariate Z at time 0 is a binary variable such that $P(Z=1) = p$, $P(Z=0) = 1-p$. What is the relationship between definition of the attributable risk function, $\phi(t)$ in (4) provided by Chen et al. (2006) and the definition of $\phi(t)$ in (5) provided by Samuelsen et al. (2008)? The survival function given the risk factor covariate Z can be written as

$$S(t | Z=1) = e^{-\int_0^t \lambda_0(u) \exp(\beta) du}, S(t | Z=0) = e^{-\int_0^t \lambda_0(u) du}. \quad (17)$$

The population (or marginal) survival function is

$$S(t) = pe^{-\int_0^t \lambda_0(u) \exp(\beta) du} + (1-p)e^{-\int_0^t \lambda_0(u) du}. \quad (18)$$

The population density function can be obtained from the population survival function as follows:

$$f(t) = -S'(t) = \lambda_0(t) \exp(\beta) pe^{-\int_0^t \lambda_0(u) \exp(\beta) du} + \lambda_0(t)(1-p)e^{-\int_0^t \lambda_0(u) du}. \quad (19)$$

The population hazard function is

$$\lambda(t) = \frac{f(t)}{S(t)} = \frac{\lambda_0(t) \exp(\beta) pe^{-\int_0^t \lambda_0(u) \exp(\beta) du} + \lambda_0(t)(1-p)e^{-\int_0^t \lambda_0(u) du}}{pe^{-\int_0^t \lambda_0(u) \exp(\beta) du} + (1-p)e^{-\int_0^t \lambda_0(u) du}}. \quad (20)$$

Based on the definition provided by Chen et al. (2006),

$$\phi(t) = \frac{\lambda(t) - \lambda(t | Z=0)}{\lambda(t)} = 1 - \frac{\lambda_0(t)}{\lambda(t)}, \quad (21)$$

where $\lambda(t)$ is give in (20). If we assume the distribution of Z does not change with t , the expectations of $\lambda(t | Z)$ and $\lambda(t | Z=0)$ are given by

$$E(\lambda(t | Z)) = \lambda_0(t) \exp(\beta) p + \lambda_0(t)(1-p), E(\lambda(t | Z=0)) = \lambda_0(t). \quad (22)$$

Then, based on the definition provided by Samuelsen et al. (2008)

$$\phi(t) = \frac{E[\lambda(t | Z)] - E[\lambda(t | Z=0)]}{E[\lambda(t | Z)]} = \frac{\exp(\beta)p - p}{\exp(\beta)p + (1-p)}, \quad (23)$$

which does not equal the value in (21). However, distribution of Z among those at risk typically changes over time (Samuelsen et al., 2008). In fact, the conditional distribution of Z given $T \geq t$ can be obtained by

$$P(Z=1 | T \geq t) = \frac{S(t | Z=1)p}{S(t)}, P(Z=0 | T \geq t) = \frac{S(t | Z=0)(1-p)}{S(t)}. \quad (24)$$

Then, from (17), (20) and (24), we have

$$\begin{aligned} E(\lambda(t | Z)) &= \lambda(t | Z=1)P(Z=1 | T \geq t) + \lambda(t | Z=0)P(Z=0 | T \geq t) \\ &= \{\lambda_0(t) \exp(\beta) S(t | Z=1)p + \lambda_0(t) S(t | Z=0)(1-p)\} / S(t) \\ &= \{\lambda_0(t) \exp(\beta) pe^{-\int_0^t \lambda_0(u) \exp(\beta) du} + \lambda_0(t)(1-p)e^{-\int_0^t \lambda_0(u) du}\} / S(t) \\ &= \lambda(t), \end{aligned}$$

and

$$E(\lambda(t | Z=0)) = \lambda(t | Z=0) = \lambda_0(t).$$

So,

$$\phi(t) = \frac{E[\lambda(t | Z)] - E[\lambda(t | Z=0)]}{E[\lambda(t | Z)]} = \frac{\lambda(t) - \lambda(t | Z=0)}{\lambda(t)}.$$

The two definitions are the same.

This article was originally published in a special issue, **Advances in Markov Chain Monte Carlo Methods and Survival Analysis** handled by Editor(s). Dr. Faming Liang, Texas A&M University, USA; Dr. Nengjun Yi, University of Alabama at Birmingham, USA; Dr. Wenqing He, University of Western Ontario, Canada; Dr. Liuquan Sun, Institute of Applied Mathematics, Academy of Mathematics and Systems Science, China