## Appendix: Mediation Analysis for Censored Survival Data under an Accelerated Failure Time Model

## October 26, 2016

For exposure A, mediator M and outcome T, let M(a) and T(a) = T(a, M(a)) define the counterfactual mediator and outcome had exposure taken value a. Likewise, let T(a, m) define the counterfactual outcome had exposure and mediator taken the value a and m, respectively. Finally let  $T(a, M(a^*))$  denote the counterfactual outcome had exposure taken value a and the mediator taken the value it would have under treatment  $a^*$ . The average pure or natural direct effect on the log-additive scale is then defined for  $a \neq a^*$ :

$$NDE(a, a^*) = E\{\log T(a, M(a^*))\} - E\{\log T(a^*)\}\}$$

and the natural indirect effect is defined as

$$NIE(a, a^*) = E\{\log T(a)\} - E\{\log T(a, M(a^*))\}\}$$

Equivalently, we could write the above expressions conditioning on a set of confounders, Z. Throughout, we make the assumption:

$$A \perp\!\!\!\perp \{T(a,m), M(a)\} \mid Z \tag{A1}$$

and we further suppose that we also have for all  $a, a^*$ :

$$T(a,m) \perp M(a^*)|A = a, Z \tag{A2}$$

Under these assumptions, it follows that  $NDE(a, a^*)$  and  $NIE(a, a^*)$  are identified empirically with<sup>2</sup>

$$E\{\log T(a, M(a^*))\} = \sum_{m, z} E\{\log T | a, m, z\} f(m | a^*, z) f(z)$$

**Derivation of the indirect effect under an AFT model:** Suppose that the following accelerated failure time model holds,

$$\log T = \beta_0 + \beta_a A + \beta_m M + \beta_z^T Z + \sigma \varepsilon \tag{A3}$$

where  $\varepsilon$  is an independent residual of arbitrary distribution and not necessarily mean zero.

Assume that M follows

$$M = \alpha_0 + \alpha_a A + \alpha_z^T Z + \xi \tag{A4}$$

where  $\xi$  is a mean zero error independent of A and Z. Then,

$$E \{ \log T(a, M(a^*)) \} = \sum_{m,z} E \{ \log T | a, m, z \} f(m | a^*, z)$$
$$= \beta_0 + \beta_a a + \beta_m E(M \mid a^*, z) + \beta_z^T z + \sigma \epsilon$$
$$= \beta_0 + \beta_a a + \beta_m \alpha_0 + \beta_m \alpha_a a^* + \alpha_z^T z + \beta_z^T z + \sigma \epsilon$$

which gives the following result,

$$NDE(a, a^{*}) = E \{\log T(a, M(a^{*}))\} - E \{\log T(a^{*}, M(a^{*}))\} \\ = E \{\log T(a, M(a^{*})) \mid Z\} - E \{\log T(a^{*}, M(a^{*})) \mid Z\} \\ = \beta_{a} (a - a^{*}) \\$$
$$NIE(a, a^{*}) = E \{\log T(a, M(a))\} - E \{\log T(a, M(a^{*}))\} \\ = E \{\log T(a, M(a)) \mid Z\} - E \{\log T(a, M(a^{*})) \mid Z\} \\ = \beta_{m} \alpha_{a} (a - a^{*})$$

Note that under the AFT model one has the stronger result that at the individual level,

$$NDE(a, a^*) = \log T(a, M(a^*)) - \log T(a^*, M(a^*))$$
  
=  $\beta_a (a - a^*)$ 

$$NIE(a, a^*) = \log T(a, M(a)) - \log T(a, M(a^*))$$
$$= \beta_m \alpha_a (a - a^*)$$

For binary A with a = 1 and  $a^* = 0$ , the indirect effect product method estimates in  $\beta_m \alpha_a$  and the natural direct effect is  $\beta_a$ . The expression for the difference method is obtained from (A3) and (A4):

$$\log T = \beta_0 + \beta_a A + \beta_m M + \beta_z^T Z + \sigma \varepsilon$$
  
=  $\beta_0 + \beta_a A + \beta_m (\alpha_0 + \alpha_a A + \alpha_z^T Z + \xi) + \sigma \varepsilon$   
=  $\beta_0 + \beta_m \alpha_0 + (\beta_a + \beta_m \alpha_a) A + (\beta_z^T + \alpha_z^T) Z + (\sigma \varepsilon + \beta_m \xi)$   
=  $\beta_0^* + \tau_a A + \beta_z^{*T} Z + \widetilde{\varepsilon}$  (A5)

where  $\tilde{\varepsilon}$  follows the distribution given by the convolution of the density of  $\sigma \varepsilon$  with that of  $\beta_m \xi$ , which is independent of A and Z. The total effect is given by  $\tau_a$  and the indirect effect from the difference method is:

$$\tau_a - \beta_a = \alpha_a \beta_m \tag{A6}$$

The difference method estimand is obtained by positing a second accelerated failure time model for T as a function of A and Z only, which shall be referred to as the reduced form model and would typically be specified as followed:

$$\log T = \beta_0^* + \tau_a A + \beta_z^{*T} Z + \sigma \nu \tag{A7}$$

where  $\sigma$  is some unknown scale parameter to be estimated. Therefore, when using the difference method, one must specify the correct distribution of  $\nu$  hoping to match that of  $\tilde{\varepsilon}$  in (A5) – failure to do so will result in model mis-specification.

Evaluating consistency of the maximum likelihood estimator for  $\tau_a$  under model misspecification and right censoring: Suppose that one mis-specifies the reduced form density of T given A and Z from model (A7) with the density  $f_T(t \mid X; \alpha, \beta, \sigma) = f_T(t \mid X)$  and survival function  $S_T(t \mid X; \alpha, \beta, \sigma) = S_T(t \mid X)$ . Let  $X = (A, Z^T)^T, \beta = (\tau_a, \beta_Z^{*T})$ , and  $\alpha$  is the intercept  $(\beta_0^* \text{ above})$ . We show below that the maximum likelihood estimator of  $\beta$ , and thus  $\tau_a$ , will be consistent in the absence of censoring. However, in the presence of censoring, the maximum likelihood estimator will not be consistent. We sketch the proof for the case of right censoring only.

The observed data is  $\min(T, C)$  and  $I(T \leq C)$  where T is event time and C is independent censoring time. The log likelihood for a single observation is:

$$\log \ell = I(T \le C) \log f_T(T \mid X) + I(T > C) \log S_T(C \mid X)$$

We can re-express this in terms of the rescaled residual error term,  $T_0 = Te^{-\alpha - \beta X} = \exp(\sigma \varepsilon)$ , which has density  $f_0(T_0 \mid X) = f_0(T_0)$  because the residual error is independent of X,

$$f_T(t \mid X) = f_0(te^{-\alpha - \beta X})e^{-\alpha - \beta X}$$

We can re-express the log likelihood:

$$\log \ell = I(T \le C) \log[f_0(te^{-\alpha - \beta X})e^{-\alpha - \beta X}] + I(T > C) \log S_0(Ce^{-\alpha - \beta X})$$
$$= I(T \le C) \log(f_0(te^{\alpha - \beta X})) - (\alpha + \beta X)I(T \le C) + I(T > C) \log(S_0(Ce^{-\beta X}))$$

The score function of  $\beta$ , can be expressed as:

$$\begin{split} U_{\beta}(\beta,\alpha) &= \frac{d}{d\beta} \bigg[ I(T \leq C) \log(f_0(te^{\alpha-\beta X})) - (\alpha+\beta X) I(T \leq C) + I(T > C) \log(S_0(Ce^{-\beta X}))) \bigg] \\ &= -XI(T \leq C) \frac{\dot{f}_0(te^{-\alpha-\beta X})te^{-\alpha-\beta X}}{f_0(te^{-\alpha-\beta X})} - XI(T \leq C) + I(T > C) \frac{\frac{d}{d\beta}S_0(Ce^{-\alpha-\beta X})}{S_0(Ce^{-\alpha-\beta X})} \\ &= -XI(T \leq C) \frac{\dot{f}_0(te^{-\alpha-\beta X})te^{-\alpha-\beta X}}{f_0(te^{-\alpha-\beta X})} - XI(T \leq C) \\ &+ \frac{I(T > C)}{S_0(Ce^{-\alpha-\beta X})} \frac{d}{d\beta} \bigg( 1 - \int_0^C f_0(te^{-\alpha-\beta X})e^{-\alpha-\beta X} dt \bigg) \\ &= -XI(T \leq C) \frac{\dot{f}_0(te^{-\alpha-\beta X})te^{-\alpha-\beta X}}{f_0(te^{-\alpha-\beta X})} - XI(T \leq C) \\ &+ \frac{I(T > C)}{S_0(Ce^{-\alpha-\beta X})} \bigg( - \int_0^C \frac{\frac{d}{d\beta}[f_0(te^{-\alpha-\beta X})e^{-\alpha-\beta X}]}{f_0(te^{-\alpha-\beta X})} f_0(te^{-\alpha-\beta X}) dt \bigg) \\ &= -XI(T \leq C) \frac{\dot{f}_0(te^{-\alpha-\beta X})te^{-\alpha-\beta X}}{f_0(te^{-\alpha-\beta X})} - XI(T \leq C) \\ &+ \frac{I(T > C)}{S_0(Ce^{-\alpha-\beta X})} \bigg( - \int_0^C \frac{\frac{d}{d\beta}[f_0(te^{-\alpha-\beta X})e^{-\alpha-\beta X}]}{f_0(te^{-\alpha-\beta X})} f_0(te^{-\alpha-\beta X}) dt \bigg) \\ &= -XI(T \leq C) \frac{\dot{f}_0(te^{-\alpha-\beta X})te^{-\alpha-\beta X}}{f_0(te^{-\alpha-\beta X})} - XI(T \leq C) \\ &+ \frac{I(T > C)}{S_0(Ce^{-\alpha-\beta X})} \bigg( X \int_0^C \frac{\dot{f}_0(te^{-\alpha-\beta X})te^{-\alpha-\beta X}}{f_0(te^{-\alpha-\beta X})} f_0(te^{-\alpha-\beta X})} f_0(te^{-\alpha-\beta X}) dt \bigg) \end{split}$$

where  $\dot{f}_0(\cdot)$  is the derivative of  $f_0(\cdot)$  with respect to its argument.

Let  $(\bar{\beta}, \bar{\alpha})$  denote the limiting value of the MLE, i.e.  $(\hat{\alpha}, \hat{\beta}) \xrightarrow{P} (\bar{\beta}, \bar{\alpha})$  where  $(\hat{\alpha}, \hat{\beta})$  is the MLE. Then,

$$E\left(\begin{bmatrix}U_{\beta}(\bar{\beta},\bar{\alpha})\\U_{\alpha}(\bar{\beta},\bar{\alpha})\end{bmatrix}\right)=0$$

Now, we can take the expectation of the score of  $\beta$  conditional on C and X. Note that  $f_0^*(\cdot)$  indicates the true law:

$$\begin{split} E[U_{\beta}(\bar{\beta},\bar{\alpha}) \mid C,X] &= E\left[-XI(T\leq C)\frac{\dot{f}_{0}(te^{-\bar{\alpha}-\bar{\beta}X})te^{-\bar{\alpha}-\bar{\beta}X}}{f_{0}(te^{-\bar{\alpha}-\bar{\beta}X})} - XI(T\leq C) \right. \\ &+ \frac{I(T>C)}{S_{0}(Ce^{-\bar{\alpha}-\bar{\beta}X})} \left(X\int_{0}^{C}\frac{\dot{f}_{0}(te^{-\bar{\alpha}-\bar{\beta}X})te^{-\bar{\alpha}-\bar{\beta}X}e^{-\bar{\alpha}-\bar{\beta}X} + f_{0}(e^{-\bar{\alpha}-\bar{\beta}X})e^{-\bar{\alpha}-\bar{\beta}X}}{f_{0}(te^{-\bar{\alpha}-\bar{\beta}X})} f_{0}(te^{-\bar{\alpha}-\bar{\beta}X}) \right] |C,X] \\ &= \int_{0}^{C} (-X)\frac{\dot{f}_{0}(te^{-\bar{\alpha}-\bar{\beta}X})te^{-\bar{\alpha}-\bar{\beta}X}}{f_{0}(te^{-\bar{\alpha}-\bar{\beta}X})} f_{0}^{*}(te^{-\alpha-\beta X})e^{-\alpha-\beta X}dt + \int_{0}^{C} (-X)f_{0}^{*}(te^{-\alpha-\beta X})e^{-\alpha-\beta X}dt \\ &+ \frac{\int_{C}^{C} f_{0}^{*}(te^{-\alpha-\beta X})e^{-\alpha-\beta X}dt}{S_{0}(Ce^{-\bar{\alpha}-\bar{\beta}X})} \left(X\int_{0}^{C}\frac{\dot{f}_{0}(te^{-\bar{\alpha}-\bar{\beta}X})te^{-\bar{\alpha}-\bar{\beta}X}}{f_{0}(te^{-\bar{\alpha}-\bar{\beta}X})} f_{0}^{*}(te^{-\alpha-\beta X})e^{-\alpha-\beta X}dt + \int_{0}^{C} (-X)f_{0}^{*}(te^{-\alpha-\beta X})e^{-\alpha-\beta X}dt \\ &+ \frac{S_{0}^{*}(Ce^{-\bar{\alpha}-\bar{\beta}X})}{f_{0}(te^{-\bar{\alpha}-\bar{\beta}X})} f_{0}^{*}(te^{-\alpha-\beta X})e^{-\alpha-\beta X}dt + \int_{0}^{C} (-X)f_{0}^{*}(te^{-\alpha-\beta X})e^{-\alpha-\beta X}dt \\ &+ \frac{S_{0}^{*}(Ce^{-\alpha-\beta X})}{f_{0}(te^{-\bar{\alpha}-\bar{\beta}X})}} f_{0}^{*}(te^{-\bar{\alpha}-\bar{\beta}X})e^{-\bar{\alpha}-\bar{\beta}X}dt + \int_{0}^{C} (-X)f_{0}^{*}(te^{-\bar{\alpha}-\bar{\beta}X})e^{-\bar{\alpha}-\bar{\beta}X}dt \\ &+ \frac{S_{0}^{*}(Ce^{-\bar{\alpha}-\bar{\beta}X})}{f_{0}(te^{-\bar{\alpha}-\bar{\beta}X})} f_{0}^{*}(te^{-\bar{\alpha}-\bar{\beta}X})} f_{0}^{*}(te^{-\bar{\alpha}-\bar{\beta}X})e^{-\bar{\alpha}-\bar{\beta}X}dt \\ &= \int_{0}^{C} (-X)\frac{\dot{f}_{0}(te^{-\bar{\alpha}-\bar{\beta}X})}{f_{0}(te^{-\bar{\alpha}-\bar{\beta}X})} f_{0}^{*}(te^{-\bar{\alpha}-\bar{\beta}X})} f_{0}^{*}(te^{-\bar{\alpha}-\bar{\beta}X})dt \\ &= \int_{0}^{C} (-X)\frac{\dot{f}_{0}(te^{-\bar{\alpha}-\bar{\beta}X})}{f_{0}(te^{-\bar{\alpha}-\bar{\beta}X})} f_{0}^{*}(te^{-\bar{\alpha}-\bar{\beta}X})}{f_{0}(te^{-\bar{\alpha}-\bar{\beta}X})}} f_{0}^{*}(te^{-\bar{\alpha}-\bar{\beta}X})dt \\ &+ \frac{S_{0}^{*}(Ce^{-\alpha-\bar{\beta}X})}{S_{0}(Ce^{-\bar{\alpha}-\bar{\beta}X})} \left(X\int_{0}^{C} \frac{\dot{f}_{0}(te^{-\bar{\alpha}-\bar{\beta}X})}{f_{0}(te^{-\bar{\alpha}-\bar{\beta}X})}} f_{0}^{*}(te^{-\bar{\alpha}-\bar{\beta}X})} f_{0}(te^{-\bar{\alpha}-\bar{\beta}X})dt \\ &+ \frac{S_{0}^{*}(Ce^{-\alpha-\bar{\beta}X})}{S_{0}(Ce^{-\bar{\alpha}-\bar{\beta}X})} \left(X\int_{0}^{C} \frac{\dot{f}_{0}(te^{-\bar{\alpha}-\bar{\beta}X})}{f_{0}(te^{-\bar{\alpha}-\bar{\beta}X})} f_{0}(te^{-\bar{\alpha}-\bar{\beta}X})} f_{0}(te^{-\bar{\alpha}-\bar{\beta}X})dt \right) \\ \end{bmatrix}$$

Note that the conditional mean for the score of  $\alpha$  is of similar form:

$$\begin{split} E[U_{\alpha}(\bar{\beta},\bar{\alpha}) \mid C,X] &= \int_{0}^{C} -\frac{\dot{f}_{0}(te^{-\bar{\alpha}-\bar{\beta}X})te^{-\bar{\alpha}-\bar{\beta}X} + f_{0}(te^{-\bar{\alpha}-\bar{\beta}X})}{f_{0}(te^{-\bar{\alpha}-\bar{\beta}X})} f_{0}^{*}(te^{-\alpha-\beta X})e^{-\alpha-\beta X}dt \\ &+ \frac{S_{0}^{*}(Ce^{-\alpha-\beta X})}{S_{0}(Ce^{-\bar{\alpha}-\bar{\beta}X})} \bigg( \int_{0}^{C} \frac{\dot{f}_{0}(te^{-\bar{\alpha}-\bar{\beta}X})te^{-\bar{\alpha}-\bar{\beta}X} + f_{0}(e^{-\bar{\alpha}-\bar{\beta}X})}{f_{0}(te^{-\bar{\alpha}-\bar{\beta}X})} f_{0}(te^{-\bar{\alpha}-\bar{\beta}X})e^{-\bar{\alpha}-\bar{\beta}X}dt \bigg) \end{split}$$

Let p be the mean vector for the vector X. Noting that  $E[U_{\beta}(\bar{\beta}, \bar{\alpha})] = 0$  and  $E[U_{\alpha}(\bar{\beta}, \bar{\alpha})] = 0$ , we can write,

$$\begin{split} E[U_{\beta}(\bar{\beta},\bar{\alpha})] &= E\left[U_{\beta}(\bar{\beta},\bar{\alpha}) - pU_{\alpha}(\bar{\beta},\bar{\alpha})\right] \\ &= E\left[-(X-p)\int_{0}^{C}\frac{\dot{f}_{0}(te^{-\bar{\alpha}-\bar{\beta}X})te^{-\bar{\alpha}-\bar{\beta}X} + f_{0}(te^{-\bar{\alpha}-\bar{\beta}X})}{f_{0}(te^{-\bar{\alpha}-\bar{\beta}X})}f_{0}^{*}(te^{-\alpha-\beta X})e^{-\alpha-\beta X}dt \\ &+ (X-p)\frac{S_{0}^{*}(Ce^{-\alpha-\beta X})}{S_{0}(Ce^{-\bar{\alpha}-\bar{\beta}X})}\left(\int_{0}^{C}\frac{\dot{f}_{0}(te^{-\bar{\alpha}-\bar{\beta}X})te^{-\bar{\alpha}-\bar{\beta}X} + f_{0}(e^{-\bar{\alpha}-\bar{\beta}X})}{f_{0}(te^{-\bar{\alpha}-\bar{\beta}X})}f_{0}(te^{-\bar{\alpha}-\bar{\beta}X})e^{-\bar{\alpha}-\bar{\beta}X}dt\right)\right] \end{split}$$

We will now plug in the true values to assess whether we get an unbiased score equation under

model mis-specification for  $\beta$ , i.e.  $E[U_{\beta}(\beta, \bar{\alpha})] = 0$ . Suppose that  $\bar{\beta} = \beta$ :

$$\begin{split} &= E\bigg[-(X-p)\int_{0}^{C}\frac{\dot{f}_{0}(te^{-\bar{\alpha}-\beta X})te^{-\bar{\alpha}-\beta X}+f_{0}(te^{-\bar{\alpha}-\beta X})}{f_{0}(te^{-\bar{\alpha}-\beta X})}f_{0}^{*}(te^{-\alpha-\beta X})e^{-\alpha-\beta X}dt \\ &\quad + (X-p)\frac{S_{0}^{*}(Ce^{-\alpha-\beta X})}{S_{0}(Ce^{-\bar{\alpha}-\beta X})}\bigg(\int_{0}^{C}\frac{\dot{f}_{0}(te^{-\bar{\alpha}-\beta X})te^{-\bar{\alpha}-\beta X}+f_{0}(e^{-\bar{\alpha}-\beta X})}{f_{0}(te^{-\bar{\alpha}-\beta X})}f_{0}(te^{-\bar{\alpha}-\beta X})e^{-\bar{\alpha}-\beta X}dt\bigg)\bigg] \\ &= E\bigg[-(X-p)\int_{0}^{Ce^{-\beta X}}\frac{\dot{f}_{0}(e^{-\bar{\alpha}}u)e^{-\bar{\alpha}}u+f_{0}(e^{-\bar{\alpha}}u)}{f_{0}(e^{-\bar{\alpha}}u)}f_{0}^{*}(e^{-\alpha}u)e^{-\alpha}du \\ &\quad + (X-p)\frac{S_{0}^{*}(Ce^{-\alpha-\beta X})}{S_{0}(Ce^{-\bar{\alpha}-\beta X})}\int_{0}^{Ce^{-\beta X}}\frac{\dot{f}_{0}(e^{-\bar{\alpha}}u)e^{-\bar{\alpha}}u+f_{0}(e^{-\bar{\alpha}}u)}{f_{0}(e^{-\bar{\alpha}}u)}f_{0}(e^{-\bar{\alpha}}u)}f_{0}(e^{-\bar{\alpha}}u)e^{-\alpha}du\bigg] \\ &= E\bigg[(p-X)\int_{0}^{Ce^{-\beta X}}\bigg[1-\frac{S_{0}^{*}(Ce^{-\alpha-\beta X})f_{0}(e^{-\bar{\alpha}}u)e^{-\alpha}}{S_{0}(Ce^{-\bar{\alpha}-\beta X})f_{0}^{*}(e^{-\alpha}u)e^{-\alpha}}}\bigg]\frac{\dot{f}_{0}(e^{-\bar{\alpha}}u)u+f_{0}(e^{-\bar{\alpha}}u)}{f_{0}(e^{-\bar{\alpha}}u)}f_{0}^{*}(e^{-\alpha}u)e^{-\alpha}}du\bigg] \\ &= \int_{0}^{\infty}E\bigg[(p-X)\bigg[1-\frac{S_{0}^{*}(Ce^{-\alpha-\beta X})f_{0}(e^{-\bar{\alpha}}u)e^{-\bar{\alpha}}}{S_{0}(Ce^{-\bar{\alpha}-\beta X})f_{0}^{*}(e^{-\alpha}u)e^{-\bar{\alpha}}}}\bigg]I(u < Ce^{-\beta X})\bigg]\frac{\dot{f}_{0}(e^{-\bar{\alpha}}u)}{f_{0}(e^{-\bar{\alpha}}u)}f_{0}^{*}(e^{-\alpha}u)e^{-\alpha}}du$$
(A8)

If there is no right censoring  $(C \to \infty)$ , and for every value of x:

$$\begin{split} &I\left(u < Ce^{-\beta x}\right) \to 1 \\ &\frac{S_T^*\left(Ce^{-\alpha - \beta x}\right)}{S_T\left(Ce^{-\bar{\alpha} - \beta x}\right)} \to 1 \end{split}$$

in which case the expectation evaluates to zero, and the score for  $\beta$  is unbiased. Additionally, if the model is not mis-specified, so that  $S_0^*(\cdot) = S_0(\cdot)$  and  $f_0^*(\cdot) = f_0(\cdot)$ , then the score for  $\beta$  will also be unbiased regardless of censoring. Thus, the association of X with T is consistent in the absence of censoring. However, in the presence of censoring, the above will not necessarily evaluate to zero. To show this, we consider a special case when X is binary:

$$\begin{split} &= \int_{0}^{\infty} E\bigg[ (p-X) X \bigg( \Big[ 1 - \frac{S_{0}^{*}(Ce^{-\alpha-\beta}) f_{0}(e^{-\bar{\alpha}}u)e^{-\bar{\alpha}}}{S_{0}(Ce^{-\bar{\alpha}-\beta}) f_{0}^{*}(e^{-\alpha}u)e^{-\alpha}} \Big] I(u < Ce^{-\beta}) - \Big[ 1 - \frac{S_{0}^{*}(Ce^{-\alpha}) f_{0}(e^{-\bar{\alpha}}u)e^{-\bar{\alpha}}}{S_{0}(Ce^{-\bar{\alpha}}) f_{0}^{*}(e^{-\alpha}u)e^{-\alpha}} \Big] I(u < C) \bigg) \\ &+ \Big[ 1 - \frac{S_{0}^{*}(Ce^{-\alpha}) f_{0}(e^{-\bar{\alpha}}u)e^{-\bar{\alpha}}}{S_{0}(Ce^{-\bar{\alpha}}) f_{0}^{*}(e^{-\alpha}u)e^{-\alpha}} \Big] I(u < C) \bigg] \frac{\dot{f}_{0}(e^{-\bar{\alpha}}u)u + f_{0}(e^{-\bar{\alpha}}u)}{f_{0}(e^{-\bar{\alpha}}u)} f_{0}^{*}(e^{-\alpha}u)e^{-\bar{\alpha}}} \bigg] I(u < C) \bigg) \\ &= \int_{0}^{\infty} E\bigg[ (X - p) X \bigg( \frac{S_{0}^{*}(Ce^{-\alpha-\beta})}{S_{0}(Ce^{-\bar{\alpha}-\beta})} I(u < Ce^{-\beta}) - \frac{S_{0}^{*}(Ce^{-\alpha})}{S_{0}(Ce^{-\bar{\alpha}})} I(u < C) \bigg) \frac{f_{0}(e^{-\bar{\alpha}}u)e^{-\bar{\alpha}}}{f_{0}^{*}(e^{-\alpha}u)e^{-\alpha}} \bigg] \\ &\times \frac{\dot{f}_{0}(e^{-\bar{\alpha}}u)u + f_{0}(e^{-\bar{\alpha}}u)}{f_{0}(e^{-\bar{\alpha}}u)} f_{0}^{*}(e^{-\alpha}u)e^{-\alpha} du \\ &= p(1 - p) \int_{0}^{\infty} E\bigg[ \frac{S_{0}^{*}(Ce^{-\alpha-\beta})}{S_{0}(Ce^{-\bar{\alpha}-\beta})} I(u < Ce^{-\beta}) - \frac{S_{0}^{*}(Ce^{-\alpha})}{S_{0}(Ce^{-\bar{\alpha}})} I(u < C) \bigg] e^{-\bar{\alpha}} [\dot{f}_{0}(e^{-\bar{\alpha}}u)u + f_{0}(e^{-\bar{\alpha}}u)] du \end{split}$$

The above expression will generally be nonzero except at exceptional laws, such as when  $\beta = 0$ . Therefore, in the presence of model mis-specification, censoring, and a non-null effect, the MLE of  $\tau_a$  will not be consistent.

Evaluating consistency of the maximum likelihood estimator for  $\tau_a$  under model misspecification and left truncation: Suppose that one mis-specifies the reduced form density of T given A and Z from model (A7) with the density  $f_T(t \mid X; \alpha, \beta, \sigma) = f_T(t \mid X)$  and survival function  $S_T(t \mid X; \alpha, \beta, \sigma) = S_T(t \mid X)$ . Let  $X = (A, Z^T)^T, \beta = (\tau_a, \beta_Z^{*T})$ , and  $\alpha$  is the intercept  $(\beta_0^* \text{ above})$ . We show below that the maximum likelihood estimator of  $\beta$ , and thus  $\tau_a$ , will be consistent in the absence of left truncation. However, in the presence of left truncation, the maximum likelihood estimate will not be consistent.

Let T be left truncated at V such that we consider  $T \mid T \geq V$  assuming that the truncation time is independent of T and X, but otherwise follows an unrestricted density. The log likelihood for a single observation subject to left truncation is:

$$\log \ell = \log f_T \left( T \mid X \right) - \log S_T \left( V \mid X \right)$$

We can re-express this in terms of the rescaled residual error term,  $T_0 = Te^{-\alpha-\beta X} = \exp(\sigma\varepsilon)$ , which has density  $f_0(T_0 \mid X) = f_0(T_0)$  because the residual error term is independent of X, the following way:

$$f_T(t \mid X) = f_0(te^{-\alpha - \beta X})e^{-\alpha - \beta X}$$

We can re-express the log likelihood:

$$\log \ell = \log[f_0(te^{-\alpha-\beta X})e^{-\alpha-\beta X}] - \log(S_0(Ve^{-\alpha-\beta X}))$$
$$= \log[f_0(te^{-\alpha-\beta X})] - (\alpha+\beta X) - \log(S_0(Ve^{-\alpha-\beta X}))$$

The score function of  $\beta$  can be expressed as:

$$\begin{split} U_{\beta}(\alpha,\beta) &= -X \frac{\dot{f}_{0}(te^{-\alpha-\beta X})te^{-\alpha-\beta X}}{f_{0}(te^{-\alpha-\beta X})} - X - \frac{\frac{d}{d\beta}S_{0}(Ve^{-\alpha-\beta X})}{S_{0}(V^{-\alpha-\beta X})} \\ &= -X \frac{\dot{f}_{0}(te^{-\alpha-\beta X})te^{-\alpha-\beta X}}{f_{0}(te^{-\alpha-\beta X})} - X - \frac{1}{S_{0}(Ve^{-\alpha-\beta X})} \frac{d}{d\beta} \left(1 - \int_{0}^{V} f_{0}(te^{-\alpha-\beta X})e^{-\alpha-\beta X}dt\right) \\ &= -X \frac{\dot{f}_{0}(te^{-\alpha-\beta X})te^{-\alpha-\beta X}}{f_{0}(te^{-\alpha-\beta X})} - X - \frac{1}{S_{0}(Ve^{-\alpha-\beta X})} \frac{d}{d\beta} \left(\int_{V}^{\infty} f_{0}(te^{-\alpha-\beta X})e^{-\alpha-\beta X}dt\right) \\ &= -X \frac{\dot{f}_{0}(te^{-\alpha-\beta X})te^{-\alpha-\beta X}}{f_{0}(te^{-\alpha-\beta X})} - X - \frac{1}{S_{0}(Ve^{-\alpha-\beta X})} \left(\int_{V}^{\infty} \frac{d}{d\beta}[f_{0}(te^{-\alpha-\beta X})e^{-\alpha-\beta X}]}{f_{0}(te^{-\alpha-\beta X})}f_{0}(te^{-\alpha-\beta X})} dt \right) \\ &= -X \frac{\dot{f}_{0}(te^{-\alpha-\beta X})te^{-\alpha-\beta X}}{f_{0}(te^{-\alpha-\beta X})} - X - \frac{1}{S_{0}(Ve^{-\alpha-\beta X})} \left(\int_{V}^{\infty} \frac{d}{d\beta}[f_{0}(te^{-\alpha-\beta X})e^{-\alpha-\beta X}]}{f_{0}(te^{-\alpha-\beta X})}f_{0}(te^{-\alpha-\beta X})} f_{0}(te^{-\alpha-\beta X})dt \right) \\ &= -X \frac{\dot{f}_{0}(te^{-\alpha-\beta X})te^{-\alpha-\beta X}}{f_{0}(te^{-\alpha-\beta X})} - X \\ &+ \frac{1}{S_{0}(Ve^{-\alpha-\beta X})} \left(X \int_{V}^{\infty} \frac{\dot{f}_{0}(te^{-\alpha-\beta X})te^{-\alpha-\beta X}}{f_{0}(te^{-\alpha-\beta X})}f_{0}(te^{-\alpha-\beta X})} f_{0}(te^{-\alpha-\beta X})dt \right) \\ &= -\frac{1}{S_{0}(Ve^{-\alpha-\beta X})} \left(X \int_{V}^{\infty} \frac{\dot{f}_{0}(te^{-\alpha-\beta X})te^{-\alpha-\beta X}}{f_{0}(te^{-\alpha-\beta X})}f_{0}(te^{-\alpha-\beta X})} f_{0}(te^{-\alpha-\beta X})dt \right) \\ &= -\frac{1}{S_{0}(Ve^{-\alpha-\beta X})} \left(X \int_{V}^{\infty} \frac{\dot{f}_{0}(te^{-\alpha-\beta X})te^{-\alpha-\beta X}}{f_{0}(te^{-\alpha-\beta X})}f_{0}(te^{-\alpha-\beta X})} f_{0}(te^{-\alpha-\beta X})dt \right) \\ &= -\frac{1}{S_{0}(Ve^{-\alpha-\beta X})}} \left(X \int_{V}^{\infty} \frac{\dot{f}_{0}(te^{-\alpha-\beta X})te^{-\alpha-\beta X}}{f_{0}(te^{-\alpha-\beta X})}f_{0}(te^{-\alpha-\beta X})} f_{0}(te^{-\alpha-\beta X})dt \right)$$

where  $f_0$  is the derivative of  $f_0$  with respect to its argument.

Let  $(\bar{\beta}, \bar{\alpha})$  denote the limiting value of the MLE, i.e.  $(\hat{\alpha}, \hat{\beta}) \xrightarrow{P} (\bar{\beta}, \bar{\alpha})$  where  $(\hat{\alpha}, \hat{\beta})$  is the MLE. Then,

$$E\left(\begin{bmatrix}U_{\beta}(\bar{\beta},\bar{\alpha})\\U_{\alpha}(\bar{\beta},\bar{\alpha})\end{bmatrix}\right)=0$$

Now, we can take the expectation of the score of  $\beta$  conditional on X with respect to the density of  $T \mid T > V$ . Note that  $f_0^*(\cdot)$  indicates the true law:

$$\begin{split} E[U_{\beta}(\bar{\beta},\bar{\alpha}) \mid X,V] &= \int_{V}^{\infty} (-X) \frac{\dot{f}_{0}(te^{-\bar{\alpha}-\bar{\beta}X})te^{-\bar{\alpha}-\bar{\beta}X}}{f_{0}(te^{-\alpha-\beta X})} \frac{f_{0}^{*}(te^{-\alpha-\beta X})}{S_{0}^{*}(Ve^{-\alpha-\beta X})} e^{-\alpha-\beta X} dt + \int_{V}^{\infty} (-X) \frac{f_{0}^{*}(te^{-\alpha-\beta X})}{S_{0}^{*}(Ve^{-\alpha-\beta X})} e^{-\alpha-\beta X} dt \\ &+ \frac{\int_{V}^{\infty} \frac{f_{0}^{*}(te^{-\alpha-\beta X})}{S_{0}^{*}(Ve^{-\alpha-\beta X})} e^{-\alpha-\beta X} dt}{S_{0}(Ve^{-\bar{\alpha}-\bar{\beta}X})} \left( X \int_{V}^{\infty} \frac{\dot{f}_{0}(te^{-\bar{\alpha}-\bar{\beta}X})te^{-\bar{\alpha}-\bar{\beta}X}}{f_{0}(te^{-\bar{\alpha}-\bar{\beta}X})} e^{-\alpha-\beta X} dt + \int_{V}^{\infty} (-X) \frac{f_{0}^{*}(te^{-\alpha-\beta X})}{S_{0}^{*}(Ve^{-\alpha-\beta X})} f_{0}(te^{-\bar{\alpha}-\bar{\beta}X}) dt \\ &= \int_{V}^{\infty} (-X) \frac{\dot{f}_{0}(te^{-\bar{\alpha}-\bar{\beta}X})te^{-\bar{\alpha}-\bar{\beta}X}}{f_{0}(te^{-\bar{\alpha}-\bar{\beta}X})} \frac{f_{0}^{*}(te^{-\alpha-\beta X})}{S_{0}^{*}(Ve^{-\alpha-\beta X})} e^{-\alpha-\beta X} dt + \int_{V}^{\infty} (-X) \frac{f_{0}^{*}(te^{-\alpha-\beta X})}{S_{0}^{*}(Ve^{-\alpha-\beta X})} e^{-\alpha-\beta X} dt \\ &+ \frac{1}{S_{0}(Ve^{-\bar{\alpha}-\bar{\beta}X})} \left( X \int_{V}^{\infty} \frac{\dot{f}_{0}(te^{-\bar{\alpha}-\bar{\beta}X})te^{-\bar{\alpha}-\bar{\beta}X}}{f_{0}(te^{-\bar{\alpha}-\bar{\beta}X})} f_{0}^{*}(te^{-\bar{\alpha}-\bar{\beta}X})} f_{0}^{*}(te^{-\bar{\alpha}-\bar{\beta}X}) dt \right) \\ &= \frac{1}{S_{0}^{*}(Ve^{-\alpha-\beta X})} \int_{V}^{\infty} (-X) \frac{\dot{f}_{0}(te^{-\bar{\alpha}-\bar{\beta}X})te^{-\bar{\alpha}-\bar{\beta}X}}{f_{0}(te^{-\bar{\alpha}-\bar{\beta}X})}} f_{0}^{*}(te^{-\bar{\alpha}-\bar{\beta}X}) e^{-\alpha-\beta X} dt \\ &+ \frac{1}{S_{0}(Ve^{-\bar{\alpha}-\bar{\beta}X})} \left( X \int_{V}^{\infty} \frac{\dot{f}_{0}(te^{-\bar{\alpha}-\bar{\beta}X})te^{-\bar{\alpha}-\bar{\beta}X}}{f_{0}(te^{-\bar{\alpha}-\bar{\beta}X})} f_{0}^{*}(te^{-\bar{\alpha}-\bar{\beta}X}) e^{-\alpha-\beta X} dt \\ &+ \frac{1}{S_{0}(Ve^{-\bar{\alpha}-\bar{\beta}X})} \left( X \int_{V}^{\infty} \frac{\dot{f}_{0}(te^{-\bar{\alpha}-\bar{\beta}X})te^{-\bar{\alpha}-\bar{\beta}X}}{f_{0}(te^{-\bar{\alpha}-\bar{\beta}X})} f_{0}^{*}(te^{-\bar{\alpha}-\bar{\beta}X}) e^{-\bar{\alpha}-\bar{\beta}X} dt \right) \end{split}$$

Note that the conditional mean for the score of  $\alpha$  is of similar form and satisfies:

$$\begin{split} E[U_{\alpha}(\bar{\beta},\bar{\alpha}) \mid X,V] &= \frac{1}{S_{0}^{*}(Ve^{-\alpha-\beta X})} \int_{V}^{\infty} (-X) \frac{\dot{f}_{0}(te^{-\bar{\alpha}-\bar{\beta}X})te^{-\bar{\alpha}-\bar{\beta}X} + f_{0}(te^{-\bar{\alpha}-\bar{\beta}X})}{f_{0}(te^{-\bar{\alpha}-\bar{\beta}X})} f_{0}^{*}(te^{-\alpha-\beta X})e^{-\alpha-\beta X}dt \\ &+ \frac{1}{S_{0}(Ve^{-\bar{\alpha}-\bar{\beta}X})} \bigg( X \int_{V}^{\infty} \frac{\dot{f}_{0}(te^{-\bar{\alpha}-\bar{\beta}X})te^{-\bar{\alpha}-\bar{\beta}X} + f_{0}(e^{-\bar{\alpha}-\bar{\beta}X})}{f_{0}(te^{-\bar{\alpha}-\bar{\beta}X})} f_{0}(te^{-\bar{\alpha}-\bar{\beta}X})e^{-\bar{\alpha}-\bar{\beta}X}dt \bigg) \end{split}$$

Let p be the mean vector for the vector X. Noting that  $E[U_{\beta}(\bar{\beta}, \bar{\alpha})] = 0$  and  $E[U_{\alpha}(\bar{\beta}, \bar{\alpha})] = 0$ , we can write,

$$\begin{split} E\left[U_{\beta}(\bar{\beta},\bar{\alpha})\right] &= E\left[U_{\beta}(\bar{\beta},\bar{\alpha}) - pU_{\alpha}(\bar{\beta},\bar{\alpha})\right] \\ &= E\left[-(X-p)\frac{1}{S_{0}^{*}(Ve^{-\alpha-\beta X})}\int_{V}^{\infty}\frac{\dot{f}_{0}(te^{-\bar{\alpha}-\bar{\beta}X})te^{-\bar{\alpha}-\bar{\beta}X} + f_{0}(te^{-\bar{\alpha}-\bar{\beta}X})}{f_{0}(te^{-\bar{\alpha}-\bar{\beta}X})}f_{0}^{*}(te^{-\alpha-\beta X})e^{-\alpha-\beta X}dt \\ &+ (X-p)\frac{1}{S_{0}(Ve^{-\bar{\alpha}-\bar{\beta}X})}\left(\int_{V}^{\infty}\frac{\dot{f}_{0}(te^{-\bar{\alpha}-\bar{\beta}X})te^{-\bar{\alpha}-\bar{\beta}X} + f_{0}(e^{-\bar{\alpha}-\bar{\beta}X})}{f_{0}(te^{-\bar{\alpha}-\bar{\beta}X})}f_{0}(te^{-\bar{\alpha}-\bar{\beta}X})e^{-\bar{\alpha}-\bar{\beta}X}dt\right)\right] \end{split}$$

We will now plug in the true values to assess whether we get an unbiased score equation under model mis-specification for  $\beta$ , i.e.  $E[U_{\beta}(\beta, \bar{\alpha})] = 0$ . Suppose that  $\bar{\beta} = \beta$ :

$$\begin{split} &= E\bigg[-(X-p)\frac{1}{S_{0}^{*}(Ve^{-\alpha-\beta X})}\int_{V}^{\infty}\frac{\dot{f}_{0}(te^{-\bar{\alpha}-\beta X})te^{-\bar{\alpha}-\beta X}+f_{0}(te^{-\bar{\alpha}-\beta X})}{f_{0}(te^{-\bar{\alpha}-\beta X})}f_{0}^{*}(te^{-\alpha-\beta X})e^{-\alpha-\beta X}dt \\ &\quad + (X-p)\frac{1}{S_{0}(Ve^{-\bar{\alpha}-\beta X})}\bigg(\int_{V}^{\infty}\frac{\dot{f}_{0}(te^{-\bar{\alpha}-\beta X})te^{-\bar{\alpha}-\beta X}+f_{0}(e^{-\bar{\alpha}-\beta X})}{f_{0}(te^{-\bar{\alpha}-\beta X})}f_{0}(te^{-\bar{\alpha}-\beta X})e^{-\bar{\alpha}-\beta X}dt\bigg)\bigg] \\ &= E\bigg[-(X-p)\frac{1}{S_{0}^{*}(Ve^{-\alpha-\beta X})}\int_{Ve^{-\beta X}}^{\infty}\frac{\dot{f}_{0}(e^{-\bar{\alpha}}u)e^{-\bar{\alpha}}u+f_{0}(e^{-\bar{\alpha}}u)}{f_{0}(e^{-\bar{\alpha}}u)}f_{0}^{*}(e^{-\alpha}u)e^{-\alpha}du \\ &\quad + (X-p)\frac{1}{S_{0}(Ve^{-\bar{\alpha}-\beta X})}\bigg(\int_{Ve^{-\beta X}}^{\infty}\frac{\dot{f}_{0}(e^{-\bar{\alpha}}u)e^{-\bar{\alpha}}u+f_{0}(e^{-\bar{\alpha}}u)}{f_{0}(e^{-\bar{\alpha}}u)}f_{0}(e^{-\bar{\alpha}}u)e^{-\bar{\alpha}}du\bigg)\bigg] \\ &= E\bigg[-(X-p)\int_{Ve^{-\beta X}}^{\infty}\bigg[\frac{1}{S_{0}^{*}(Ve^{-\alpha-\beta X})}-\frac{f_{0}(e^{-\bar{\alpha}}u)e^{-\bar{\alpha}}S_{0}(Ve^{-\bar{\alpha}-\beta X})}{f_{0}(e^{-\bar{\alpha}}u)}\bigg] \\ &\quad \times \frac{\dot{f}_{0}(e^{-\bar{\alpha}}u)e^{-\bar{\alpha}}}{f_{0}(e^{-\bar{\alpha}}u)}f_{0}^{*}(e^{-\alpha}u)e^{-\alpha}du\bigg] \\ &= \int_{0}^{\infty}E\bigg[-(X-p)\bigg[\frac{1}{S_{0}^{*}(Ve^{-\alpha-\beta X})}-\frac{f_{0}(e^{-\bar{\alpha}}u)e^{-\bar{\alpha}}S_{0}(Ve^{-\bar{\alpha}-\beta X})}{f_{0}^{*}(e^{-\alpha}u)}\bigg]I(u>Ve^{-\beta X})\bigg] \\ &\quad \times \frac{\dot{f}_{0}(e^{-\bar{\alpha}}u)e^{-\bar{\alpha}}S_{0}(Ve^{-\bar{\alpha}-\beta X})}{f_{0}(e^{-\bar{\alpha}}u)}}f_{0}^{*}(e^{-\alpha}u)e^{-\alpha}du\bigg]$$

If there is no left truncation (V = 0), and for every value of x:

$$I(u > Ve^{-\beta x}) = I(u > 0) = 1$$
  
$$\frac{1}{S_0^*(Ve^{-\alpha - \beta X})} - \frac{f_0(e^{-\bar{\alpha}}u)e^{-\bar{\alpha}}}{f_0^*(e^{-\alpha}u)e^{-\alpha}S_0(Ve^{-\bar{\alpha} - \beta X})} = 1 - \frac{f_0(e^{-\bar{\alpha}}u)e^{-\bar{\alpha}}}{f_0^*(e^{-\alpha}u)e^{-\alpha}}$$

in which case the expectation evaluates to zero, and the score for  $\beta$  is unbiased. Additionally, if the model is not mis-specified, so that  $S_0^*(\cdot) = S_0(\cdot)$  and  $f_0^*(\cdot) = f_0(\cdot)$ , then the score for  $\beta$  will also be unbiased. Thus, the association of X with T is consistent in the absence of censoring. However, in the presence of left truncation, the above will not necessarily evaluate to zero. To show this, we consider a special case when X is binary:

$$\begin{split} &= \int_{0}^{\infty} E \left[ -(X-p)X \left( \left[ \frac{1}{S_{0}^{*}(Ve^{-\alpha-\beta})} - \frac{f_{0}(e^{-\bar{\alpha}}u)e^{-\bar{\alpha}}}{f_{0}^{*}(e^{-\alpha}u)e^{-\alpha}S_{0}(Ve^{-\bar{\alpha}-\beta})} \right] I(u > Ve^{-\beta}) \right. \\ &\quad \left. - \left[ \frac{1}{S_{0}^{*}(Ve^{-\alpha})} - \frac{f_{0}(e^{-\bar{\alpha}}u)e^{-\bar{\alpha}}}{f_{0}^{*}(e^{-\alpha}u)e^{-\alpha}S_{0}(Ve^{-\bar{\alpha}})} \right] I(u > V) \right) \right. \\ &\quad \left. + \frac{f_{0}(e^{-\bar{\alpha}}u)e^{-\bar{\alpha}}}{f_{0}^{*}(e^{-\alpha}u)e^{-\alpha}S_{0}(Ve^{-\bar{\alpha}})} \right] I(u > V) \right] \frac{\dot{f}_{0}(e^{-\bar{\alpha}}u)e^{-\bar{\alpha}}u + f_{0}(e^{-\bar{\alpha}}u)e^{-\alpha}du \\ &\quad \left. + \frac{f_{0}(e^{-\bar{\alpha}}u)e^{-\bar{\alpha}}}{f_{0}^{*}(e^{-\alpha}u)e^{-\alpha}S_{0}(Ve^{-\bar{\alpha}})} \right] I(u > V) \right] \\ &\quad \left. + \frac{f_{0}(e^{-\bar{\alpha}}u)e^{-\alpha}}{f_{0}^{*}(e^{-\alpha}u)e^{-\alpha}S_{0}(Ve^{-\bar{\alpha}})} \right] I(u > V) \right] \\ &\quad \left. + \frac{f_{0}(e^{-\bar{\alpha}}u)e^{-\alpha}}{f_{0}^{*}(e^{-\alpha}u)e^{-\alpha}S_{0}(Ve^{-\bar{\alpha}-\beta})} \right] I(u > Ve^{-\beta}) \\ &\quad \left. - \left[ \frac{1}{S_{0}^{*}(Ve^{-\alpha})} - \frac{f_{0}(e^{-\bar{\alpha}}u)e^{-\bar{\alpha}}}{f_{0}^{*}(e^{-\alpha}u)e^{-\alpha}S_{0}(Ve^{-\bar{\alpha}-\beta})} \right] I(u > Ve^{-\beta}) \\ &\quad \left. - \left[ \frac{1}{S_{0}^{*}(Ve^{-\alpha})} - \frac{f_{0}(e^{-\bar{\alpha}}u)e^{-\bar{\alpha}}}{f_{0}^{*}(e^{-\alpha}u)e^{-\alpha}S_{0}(Ve^{-\bar{\alpha}-\beta})} \right] I(u > Ve^{-\beta}) \\ &\quad \left. - \left[ \frac{1}{S_{0}^{*}(Ve^{-\alpha})} - \frac{f_{0}(e^{-\bar{\alpha}}u)e^{-\bar{\alpha}}}{f_{0}^{*}(e^{-\alpha}u)e^{-\alpha}S_{0}(Ve^{-\bar{\alpha}-\beta})} \right] I(u > V) \right] \frac{\dot{f}_{0}(e^{-\bar{\alpha}}u)e^{-\bar{\alpha}}u}}{f_{0}(e^{-\bar{\alpha}}u)}f_{0}^{*}(e^{-\alpha}u)e^{-\alpha}du \\ &\quad \left. - \left[ \frac{1}{S_{0}^{*}(Ve^{-\alpha})} - \frac{f_{0}(e^{-\bar{\alpha}}u)e^{-\bar{\alpha}}}{f_{0}^{*}(e^{-\alpha}u)e^{-\bar{\alpha}}S_{0}(Ve^{-\bar{\alpha}-\beta})} \right] I(u > V) \right] \frac{\dot{f}_{0}(e^{-\bar{\alpha}}u)}{f_{0}(e^{-\bar{\alpha}}u)}f_{0}^{*}(e^{-\alpha}u)e^{-\alpha}du \\ &\quad \left. - \left[ \frac{1}{S_{0}^{*}(Ve^{-\alpha})} - \frac{f_{0}(e^{-\bar{\alpha}}u)e^{-\bar{\alpha}}}{f_{0}^{*}(e^{-\alpha}u)e^{-\bar{\alpha}}S_{0}(Ve^{-\bar{\alpha}-\beta})} \right] I(u > V) \right] \frac{\dot{f}_{0}(e^{-\bar{\alpha}}u)}{f_{0}(e^{-\bar{\alpha}}u)}f_{0}^{*}(e^{-\alpha}u)e^{-\alpha}du \\ &\quad \left. - \left[ \frac{1}{S_{0}^{*}(Ve^{-\alpha})} - \frac{f_{0}(e^{-\bar{\alpha}}u)e^{-\bar{\alpha}}}{f_{0}^{*}(e^{-\alpha}u)e^{-\bar{\alpha}}S_{0}(Ve^{-\bar{\alpha}-\beta})} \right] I(u > V) \right] \frac{\dot{f}_{0}(e^{-\bar{\alpha}}u)}{f_{0}(e^{-\bar{\alpha}}u)}f_{0}^{*}(e^{-\alpha}u)e^{-\alpha}du \\ &\quad \left. - \left[ \frac{1}{S_{0}^{*}(Ve^{-\alpha})} - \frac{f_{0}(e^{-\bar{\alpha}}u)e^{-\bar{\alpha}}}{f_{0}^{*}(e^{-\alpha}u)e^{-\bar{\alpha}}} \right] I(u > V) \right] \frac{\dot{f}_{0}(e^{-\bar{\alpha}}u)}{f_{0}(e^{-\bar{\alpha}}u)}f_{0}^{*}(e^{-\alpha}u)e^{-\alpha}du \\ &\quad \left. - \left[ \frac{1}{S_{0}^{*}(Ve^{-\alpha})}$$

The above expression will generally be nonzero except at exceptional laws, such as when  $\beta = 0$ . Therefore, in the presence of model mis-specification, censoring, and a non-null effect, the MLE of  $\tau_a$  will not be consistent.

An issue with equivalence of the product and difference method indirect effect under a AFT model with a Weibull outcome, no censoring: Consider model (A3) and (A4), where  $\varepsilon$  follows an extreme value distribution and  $\xi$  is normally distributed. Then the implied reduced form model is given by:

$$\log T = \beta_0 + \beta_a A + \beta_m M + \beta_z^T Z + \sigma \varepsilon$$
  
=  $\beta_0 + \beta_a A + \beta_m (\alpha_0 + \alpha_a A + \alpha_z^T Z \xi) + \beta_z^T Z + \sigma \varepsilon$   
=  $\beta_0 + \beta_m \alpha_0 + (\beta_a + \beta_m \alpha_a) A + (\alpha_z^T + \beta_z^T) Z + (\sigma \varepsilon + \beta_m \xi)$   
=  $\beta_0^* + \tau_a A + \beta_z^{*T} Z + \widetilde{\varepsilon}$  (A10)

where  $\beta_0^* = \beta_m \alpha_0 + \beta_0$ ,  $\beta_z^{*T} = \alpha_z^T + \beta_z^T$ ,  $\tilde{\varepsilon} = \beta_m \xi + \sigma \varepsilon$  and  $\tau_a = \alpha_a \beta_m + \beta_a$ . The above model is an AFT model since  $\tilde{\varepsilon}$  is independent of A and C which follows from  $(\xi, \varepsilon)$  independent of A and C. However, the reduced-form density of log T given A and C is of a complicated form given by the convolution of a normal density with an extreme value density:  $f_{\tilde{\varepsilon}}(\cdot) = \int_{\varepsilon} \frac{1}{\beta_m} f_{\xi}(\frac{\cdot -\sigma \epsilon}{\beta_m}) g(\varepsilon) d\varepsilon$ , where  $g(\varepsilon)$  is the extreme value density and  $\beta_m \neq 0$ . Thus,  $\tilde{\varepsilon}$  will not have an extreme value distribution, so that the reduced form model is mis-specified if an extreme value density is assumed for  $f_{\tilde{\varepsilon}}$ . As we showed in the previous section, in the presence of censoring, the estimator of  $\tau_a$  will therefore fail to be consistent; thus, the difference method indirect effect estimator will not be consistent for the indirect effect. However, according to our results, in the absence of censoring, the difference method estimator will be consistent for the indirect effect.

Equivalence of the product and difference method indirect effect under a AFT model with a log-normal outcome: In contrast, if  $\varepsilon$  and  $\xi$  are both normal, the reduced-form density of log T given A and C is of correct form because the convolution of two independent normal densities



Figure A1: Simulation Study, Product vs. Difference Method for the Indirect Effect

will also be a normal density. Due to this, the reduced form model (A7) will be correctly specified, so the estimator of  $\tau_a$  will be consistent. Thus, the difference method,  $\tau_a - \beta_a$ , will be a consistent estimator for the indirect effect.

## Monte Carlo variance for indirect effect estimates in the simulation study:

R Code for direct and indirect effect estimates from data application:

```
##Calculate indirect and direct effect estimates:
              #normally distributed time to event outcome
              #normally distributed mediator
              #no interaction between exposure and mediator
              #interval and right censoring
#exp is the exposure variable (A in the paper)
#med is the mediator variable (M in the paper)
#time1 is the left interval
#time2 is the right interval; NA for right censored data
#cov1,..,cov5 are the potential confounders
#choose the correct library in R
library(survival)
#full model
full.model <- survreg(Surv(time1,time2,type=c('interval2')) ~ exp + med + cov1 + cov2 +</pre>
     cov3 + cov4 + cov5, dist="gaussian")
#reduced model
exp.model <- survreg(Surv(time1,time2,type=c('interval2')) ~ exp + cov1 + cov2 + cov3</pre>
     cov4 + cov5, dist="gaussian")
#mediator model
med.model <- lm(med \sim exp + cov1 + cov2 + cov3 + cov4 + cov5)
#Calculating direct and indirect effects
nde <- full.model$coefficients[2]</pre>
nie.prod <- med.model$coefficients[2]*full.model$coefficients[3]</pre>
nie.diff <- exp.model$coefficients[2]-full.model$coefficients[2]</pre>
#Calculating standard errors for the indirect (product) and direct effect estimates
se_nde <- sqrt(full.model$var[2,2])</pre>
se_nie.prod <- sqrt((med.model$coefficients[2]^2)*full.model$var[3,3] +</pre>
     (full.model$coefficients[3]^2)*summary(med.model)$cov[2,2])
##Calculate indirect and direct effect estimates:
              #Weibull distributed time to event outcome
              #normally distributed mediator
              #no interaction between exposure and mediator
              #right censoring
#exp is the exposure variable (A in the paper)
#med is the mediator variable (M in the paper)
```

```
#outcome is the time of event or censoring
#censor is a binary variable indicating censoring
#cov1,..,cov3 are the potential confounders
#full model
full.model <- survreg(Surv(outcome, censor) ~ exp + med + cov1 + cov2 + cov3,</pre>
     dist="weibull")
#reduced model -- Recall the total effect is biased!
exp.model <- survreg(Surv(outcome, censor) ~ exp + cov1 + cov2 + cov3, dist="weibull")</pre>
#mediator model
med.model <- lm(med \sim exp + cov1 + cov2 + cov3)
#Calculating direct and indirect effects
nde <- full.model$coefficients[2]</pre>
nie.prod <- med.model$coefficients[2]*full.model$coefficients[3]</pre>
nie.diff <- exp.model$coefficients[2]-full.model$coefficients[2] #this is biased!</pre>
#Calculating standard errors for the indirect (product) and direct effect estimates
se_nde <- sqrt(full.model$var[2,2])</pre>
se_nie.prod <- sqrt((med.model$coefficients[2]^2)*full.model$var[3,3] +</pre>
```

```
(full.model$coefficients[3]^2)*summary(med.model)$cov[2,2])
```

## ##### NOTES #####

#The nie.diff estimator under the Weibull model will be biased in the presence of censoring #To calculate the standard errors for the indirect effect (difference), use the boostrap #Bootstrap code available upon request