



A Model for Informative Censoring

William A. Link

To cite this article: William A. Link (1989) A Model for Informative Censoring, Journal of the American Statistical Association, 84:407, 749-752, DOI: [10.1080/01621459.1989.10478829](https://doi.org/10.1080/01621459.1989.10478829)

To link to this article: <https://doi.org/10.1080/01621459.1989.10478829>



Published online: 12 Mar 2012.



Submit your article to this journal [↗](#)



Article views: 35



Citing articles: 10 View citing articles [↗](#)

A Model for Informative Censoring

WILLIAM A. LINK*

In the usual model for censored survival analysis, observations are of the form $X = \min(T, C)$, where T and C are nonnegative random variables representing *lifetime* and *censoring time*, respectively. Under the usual assumption of independence of T and C , the Kaplan–Meier estimator (KME) is the appropriate estimator of $S(t) = \Pr(T > t)$. The KME can lead to substantial overestimates of survival probabilities if the event of censoring indicates an unfavorable prognosis for future survival. We consider an alternative model in which censoring only occurs in a subpopulation defined by the frailty distribution. A self-consistent estimator of the survival function appropriate to the model is obtained.

KEY WORDS: Frailty; Kaplan–Meier estimator; Self-consistency.

1. INTRODUCTION

Suppose that T_1, T_2, \dots, T_n is a random sample of *lifetimes* (nonnegative continuous random variables) with common survival function $S(t) = \Pr(T > t)$. We consider the problem of estimating $S(\cdot)$ when the T 's are not directly observable; rather, one is able to observe $(X_1, \delta_1), (X_2, \delta_2), \dots, (X_n, \delta_n)$, where $X_i \leq T_i$ and δ_i is a binary random variable equaling 1 if $X_i = T_i$ and 0 otherwise.

The problem of estimating a survival function in the presence of random right censoring has been extensively studied. The majority of research has centered on the independent censoring model, in which C_1, C_2, \dots, C_n are *censoring times*, independent of T_1, T_2, \dots, T_n , and $X_i = \min(T_i, C_i)$. Under this model, the Kaplan–Meier estimator (KME) (Kaplan and Meier 1958) is the appropriate estimator of $S(\cdot)$.

A great deal is known about the KME. Földes and Rejtő (1981) established strong uniform consistency of the KME, and Gill (1983) provided weak convergence results on the entire positive half-line. Furthermore, it is well known that the KME is the generalized maximum likelihood estimator of the survival function (e.g., see Miller 1981, pp. 57–59). Bivariate versions of the KME were considered by Campbell (1981) and Korwar and Dahiya (1982).

It is not difficult to envision situations in which the assumption of independent censoring is inappropriate. If the only observations available are the pairs (X, δ) , however, the independence assumption is completely untestable. It has been shown by Cox (1959) and Tsiatis (1975) that “there always exist independent censoring models consistent with any probability distribution for the observable pair (X, δ) ” (Lagakos 1979, p. 152). The consequence of this is that, if it is believed that the independence assumption is unwarranted, an equally untestable assumption about the joint distribution of (T, C) must be made.

As unpalatable as this may be, it is a practical necessity. Lagakos (1979) mentioned the following three situations in which the independence assumption is of questionable validity:

1. a clinical trial in which some patients remove themselves from study

for reasons possibly related to therapy and thereby censor their survival time under test conditions;

2. a clinical trial in which those patients experiencing a specific critical event such as metastatic spread of disease are, by design, removed from study and no longer followed for survival time;

3. a clinical trial or animal experiment in which failure times from causes of secondary interest are recorded as censored observations of the failure times from the causes of primary interest. (p. 151)

In each of these cases, censoring indicates an unfavorable prognosis for future survival. The KME, which does not take this into account, will tend to overestimate the true survival probabilities. It is apparent, therefore, that when censoring carries an unfavorable prognosis for future survival, reasonable estimators should be bounded above by the KME and below by the empirical survival function of the observed random variable X , which we refer to as $\hat{H}(\cdot)$.

Williams and Lagakos (1977) and Lagakos and Williams (1978) defined a “cone-class” of censoring models that are indexed by a parameter $\theta \in [0, 1]$. When $\theta = 0$, censoring immediately precedes death; \hat{H} is the appropriate estimator of $S(\cdot)$. When $\theta = 1$, the KME is the appropriate estimator. Their procedure for estimating $S(\cdot)$ under the cone-class model assumption involves specifying that $S(\cdot)$ is a member of some parametric family of distributions, the parameters of which, along with various model parameters, are estimated by maximum likelihood.

Another procedure for obtaining alternative estimators to the KME was proposed by Robertson and Uppuluri (1984). Their procedure is based on a modification of the well-known “redistribute to the right” algorithm due to Efron (1967). There do not appear to be easily constructive models justifying most redistribution schemes.

There would appear to be a need for some simple models and corresponding survival function estimators, applicable when censoring carries an unfavorable prognosis for future survival. In this article a model in which censoring can occur only in a “high-risk” (low-risk) subpopulation is proposed. This model suggests a modification of Efron’s self-consistency algorithm that leads to a modified Kaplan–Meier estimator (MKME).

2. THE MODEL

The model for prognostic censoring considered here postulates heterogeneity in survival probabilities as de-

* William A. Link is Mathematical Statistician, Patuxent Wildlife Research Center, U.S. Fish and Wildlife Service, Laurel, MD 20708. This research was conducted while Link was a doctoral candidate at the University of Massachusetts at Amherst. The author thanks Ramesh Korwar (his doctoral advisor), the referee, and the associate editor, whose helpful comments improved the quality of the manuscript.

scribed by the "frailty model" of Vaupel, Manton, and Stallard (1979). Associated with each lifetime T is a random variable Z , called the frailty. The frailty model specifies that "the death rate or hazard at age t for a person with frailty z is assumed to be of the form $\mu(t; z) = z\mu(t)$, where $\mu(t)$ is independent of z and describes the age effect" (Hougaard 1984, p. 75). Thus conditional on frailty $Z = z$, the survival function $\Pr(T > t | Z = z)$ is given by

$$S(t | Z = z) = \exp \left\{ - \int_0^t z\mu(s) ds \right\}. \quad (2.1)$$

The objective is to estimate the population survival function $S(t) = E\{S(t | Z)\}$. The data set to be used in estimation consists of pairs (X, δ) , where $X \leq T$ and δ is a binary random variable equaling 1 if $X = T$ and 0 otherwise.

The independent censoring model describes the relationship between T and X by introducing a censoring time C , independent of T , and defining X as the minimum of T and C . The alternative proposed here is to assume that censoring is only possible for individuals with high (or low) frailty values. That is, let C be a *potential* censoring time: a censored observation is recorded iff $T \leq C$ and Z is large (small). Letting A denote the set of values of Z for which censoring is possible, the relationship between X and T is given by

$$X = (1 - \gamma_A)T + \gamma_A \min\{T, C\},$$

where γ_A is the indicator function for the event $\{Z \in A\}$.

From (2.1) it is seen that $S(t | Z = z)$ is decreasing in z , so individuals with high frailties tend to have smaller lifetimes. It follows that for $A = \{z | z \geq a\}$, censoring will be heavier on small observations and lighter on large observations than under the independent censoring model. Under these conditions, use of the KME leads to overestimates of $S(t)$. For $A = \{z | z \leq a\}$, the opposite effect occurs.

The model can be used to assess the robustness of the independent censoring model. In Sections 3 and 4, this is illustrated with examples in which the frailty is taken to be a unit exponential random variable. The unit exponential frailty is attractive because for any survival function $S(t)$ there exists an age effect function $\mu(t)$ that combined with a unit exponential frailty distribution yields the desired population survival function [Take $M(t)$, the integrated age effect, to be $(1 - S(t))/S(t)$.] Note that the unit exponential frailty implies that the distribution of frailty among survivors at time t is exponential with parameter $1/S(t)$, so $E(Z | T > t) = S(t)$.

3. SURVIVAL FUNCTION ESTIMATOR

Letting $0 = x_0 < x_1 < x_2 < \dots < x_n$ represent the ordered times of observation and $\delta_{(1)}, \delta_{(2)}, \dots, \delta_{(n)}$ represent the corresponding values of δ , the KME is the unique limit (as $K \rightarrow \infty$) of the sequence of functions obtained by

$$\tilde{S}^{(K+1)}(t) = \frac{1}{n} \left\{ \sum_{i=1}^n I(x_i > t) + \sum_{i=1}^n (1 - \delta_{(i)}) \frac{\tilde{S}^{(K)}(t)}{\tilde{S}^{(K)}(x_i)} \right\}.$$

Consequently, the KME satisfies

$$\hat{S}(t) = \frac{1}{n} \left\{ \sum_{i=1}^n I(x_i > t) + \sum_{i=1}^n (1 - \delta_{(i)}) \frac{\hat{S}(t)}{\hat{S}(x_i)} \right\},$$

which is to say that the estimated probability of survival beyond time t is the percentage of observations (censored or uncensored) beyond time t plus the estimated percentage that would have survived beyond time t but were censored before t . The KME is said to be "self-consistent" because, in the independent censoring model, for $x_i < t$,

$$\Pr(T > t | X = x_i, \delta = 0) = S(t)/S(x_i).$$

Under the model discussed in Section 2,

$$\Pr(T > t | X = x_i, \delta = 0) = \frac{S(t | Z \in A)}{S(x_i | Z \in A)},$$

suggesting that the self-consistency algorithm be replaced by

$$\tilde{S}^{(K+1)}(t) = \frac{1}{n} \left\{ \sum_{i=1}^n I(x_i > t) + \sum_{i=1}^n (1 - \delta_{(i)}) \left(\frac{\tilde{S}^{(K)}(t | Z \in A)}{\tilde{S}^{(K)}(x_i | Z \in A)} \right) \right\}, \quad (3.1)$$

where $\tilde{S}^{(K)}(t | Z \in A)$ is the estimate of $S(t | Z \in A)$ based on $\tilde{S}^{(K)}(t)$. The limit as $k \rightarrow \infty$ of this sequence of survival function estimators will be referred to as the modified Kaplan-Meier estimator (MKME).

Assuming a unit exponential frailty, $S(t | Z \in A)$ is given by

$$S(t) \exp \left(-a \frac{1 - S(t)}{S(t)} \right),$$

for $A = \{z | z \geq a\}$ ($a > 0$), and by

$$S(t) \left(\frac{1 - \exp(-a/S(t))}{1 - e^{-a}} \right),$$

for $A = \{z | z \leq a\}$ ($a > 0$). Substitution of these values in the algorithm described by (3.1) yields the MKME for high (low) exponential frailty censoring. Note that for $t > x$, $\Pr(T > t | X = x, \delta = 0)$ is less than $S(t)/S(x)$ in the case of high exponential frailty censoring (indicating that the KME will overestimate survival probabilities) and greater than $S(t)/S(x)$ in the case of low exponential frailty censoring (indicating that the KME will underestimate survival probabilities).

4. SIMULATION RESULTS

Samples of 4,000 pairs (U_T, γ_A) were generated by letting $U_T = E/(E + Z)$ and $\gamma_A = I(Z \in A)$, where E and Z are independent unit exponential variates and $A \subseteq \mathbf{R}^+$. Conditional on Z , the hazard function for U_T is $\mu(t; z) = z\mu(t)$, $t \in (0, 1)$, where the age effect is given by $\mu(t) = (1 - t)^{-2}$. Unconditionally, U_T is uniformly distributed on $(0, 1)$, so the U_T 's can be thought of as the quantiles of a random sample from an arbitrary continuous survival function $S(\cdot)$. Observations of the form (U_X, δ) were then

Table 1. Survival Rates From KME (MKME) Under High Exponential Frailty Censoring ($n = 4,000$)

	Proportion subject to censoring								
	.10	.20	.30	.40	.50	.60	.70	.80	.90
True	Estimated								
.900	.901 (.892)	.902 (.902)	.903 (.901)	.896 (.897)	.904 (.900)	.907 (.900)	.902 (.891)	.906 (.901)	.903 (.900)
.800	.795 (.796)	.819 (.800)	.805 (.799)	.805 (.785)	.804 (.802)	.815 (.790)	.809 (.781)	.809 (.808)	.806 (.810)
.700	.694 (.696)	.722 (.693)	.716 (.699)	.718 (.683)	.712 (.707)	.720 (.682)	.723 (.671)	.713 (.711)	.711 (.707)
.600	.597 (.591)	.628 (.596)	.620 (.596)	.622 (.594)	.626 (.616)	.634 (.570)	.630 (.575)	.626 (.617)	.613 (.609)
.500	.498 (.495)	.533 (.494)	.527 (.496)	.526 (.501)	.532 (.511)	.554 (.476)	.546 (.484)	.535 (.512)	.523 (.516)
.400	.402 (.397)	.435 (.398)	.440 (.395)	.434 (.390)	.441 (.413)	.469 (.389)	.463 (.385)	.451 (.406)	.446 (.399)
.300	.299 (.302)	.326 (.295)	.338 (.295)	.328 (.289)	.354 (.306)	.372 (.301)	.374 (.286)	.377 (.307)	.361 (.302)
.200	.194 (.202)	.212 (.200)	.230 (.205)	.219 (.193)	.254 (.197)	.268 (.201)	.286 (.188)	.296 (.209)	.281 (.207)
.100	.102 (.103)	.109 (.094)	.112 (.103)	.114 (.093)	.141 (.092)	.142 (.102)	.161 (.092)	.206 (.102)	.199 (.119)
	Proportion censored								
	.021	.053	.086	.130	.166	.222	.265	.330	.413

obtained, where

$$U_X = (1 - \gamma_A)U_T + \gamma_A \min\{U_T, U_C\}$$

and

$$\delta = 1 - \gamma_A I(U_T > U_C),$$

and the potential censoring variable U_C , uniformly distributed on $(0, 1)$, was generated independently of E and Z . Sets A of the form $A = (a, \infty)$ and $[0, a)$ were used (for high and low exponential frailty censoring, respectively), with a chosen so that $\Pr(Z \in A) = .1(.1).9$. Results for these cases are summarized in Tables 1 and 2, respec-

tively. As anticipated, the KME tends to overestimate survival probabilities in the former case and to underestimate survival probabilities in the latter case.

For both sets of simulations the MKME, obtained as the limit as $k \rightarrow \infty$ of the sequence given by (3.1), was computed. These results are also summarized in Tables 1 and 2.

5. DISCUSSION

The model considered in this article offers an alternative to the usual independent censoring model. The event of censoring is related to the frailty of the individual. Simulations involving an exponential frailty distribution show

Table 2. Survival Rates From KME (MKME) Under Low Exponential Frailty Censoring ($n = 4,000$)

	Proportion subject to censoring								
	.10	.20	.30	.40	.50	.60	.70	.80	.90
True	Estimated								
.90	.901 (.904)	.904 (.897)	.898 (.900)	.901 (.892)	.905 (.896)	.894 (.897)	.900 (.900)	.902 (.905)	.892 (.895)
.80	.806 (.807)	.804 (.798)	.800 (.801)	.792 (.796)	.799 (.793)	.791 (.803)	.801 (.792)	.808 (.810)	.784 (.797)
.70	.705 (.701)	.709 (.701)	.694 (.704)	.687 (.698)	.693 (.689)	.685 (.706)	.688 (.698)	.713 (.707)	.682 (.695)
.60	.604 (.595)	.599 (.594)	.574 (.595)	.574 (.586)	.577 (.589)	.570 (.605)	.582 (.596)	.613 (.600)	.582 (.599)
.50	.502 (.499)	.488 (.496)	.462 (.499)	.454 (.486)	.473 (.490)	.462 (.502)	.476 (.493)	.501 (.492)	.485 (.504)
.40	.390 (.399)	.369 (.393)	.341 (.405)	.335 (.381)	.359 (.390)	.341 (.398)	.365 (.398)	.389 (.383)	.383 (.401)
.30	.273 (.292)	.253 (.292)	.229 (.297)	.213 (.285)	.241 (.297)	.226 (.292)	.257 (.305)	.290 (.281)	.285 (.301)
.20	.164 (.196)	.140 (.194)	.117 (.202)	.115 (.187)	.138 (.182)	.127 (.188)	.168 (.186)	.183 (.171)	.185 (.203)
.10	.056 (.101)	.041 (.095)	.037 (.088)	.044 (.082)	.066 (.078)	.066 (.095)	.090 (.099)	.104 (.071)	.074 (.118)
	Proportion censored								
	.802	.163	.224	.273	.343	.373	.413	.462	.478

that the KME can produce substantial errors in estimation of survival probabilities.

As pointed out by Hougaard (1984), the frailty model is overparameterized unless, perhaps, the frailty is common to several individuals. In the present setting, the value of the frailty model may lie in producing a class of survival function estimates so that the potential effect of incorrectly assuming the independent censoring model can be assessed.

[Received December 1987. Revised February 1989.]

REFERENCES

- Campbell, G. (1981), "Nonparametric Bivariate Estimation With Randomly Censored Data," *Biometrika*, 68, 417-422.
- Cox, D. R. (1959), "The Analysis of the Exponentially Distributed Lifetimes With Two Types of Failure," *Journal of the Royal Statistical Society, Ser. B*, 59, 411-421.
- Efron, B. (1967), "The Two Sample Problem With Censored Data," in *Proceedings of the 5th Berkeley Symposium on Mathematical Statistics and Probability* (Vol. 4), Berkeley: University of California Press, pp. 831-853.
- Földes, A., and Rejtő, L. (1981), "Strong Uniform Consistency of the Product-Limit Estimator Under Variable Censoring," *Zeitschrift für Wahrscheinlichkeitstheorie und Verwandte Gebiete*, 58, 95-108.
- Gill, R. (1983), "Large Sample Behavior of the Product-Limit Estimator on the Whole Line," *The Annals of Statistics*, 11, 49-58.
- Hougaard, P. (1984), "Life Table Methods for Heterogeneous Populations: Distributions Describing the Heterogeneity," *Biometrika*, 71, 75-83.
- Kaplan, E. L., and Meier, P. (1958), "Nonparametric Estimation From Incomplete Observations," *Journal of the American Statistical Association*, 53, 457-481.
- Korwar, R. M., and Dahiya, R. C. (1982), "Estimation of a Bivariate Distribution Function From Incomplete Observations," *Communications in Statistics, Part A—Theory and Methods*, 12, 887-897.
- Lagakos, S. W. (1979), "General Right Censoring and Its Impact on the Analysis of Survival Data," *Biometrics*, 35, 139-156.
- Lagakos, S. W., and Williams, J. S. (1978), "Models for Censored Survival Analysis: A Cone Class of Variable-Sum Models," *Biometrika*, 65, 181-189.
- Miller, R. G. (1981), *Survival Analysis*, New York: John Wiley.
- Robertson, J. B., and Uppuluri, V. R. R. (1984), "A Generalized Kaplan-Meier Estimator," *The Annals of Statistics*, 12, 366-371.
- Tsiatis, A. (1975), "A Nonidentifiability Aspect of the Problem of Competing Risks," *Proceedings of the National Academy of Science*, 72, 20-22.
- Vaupel, J. W., Manton, K. G., and Stallard, E. (1979), "The Impact of Heterogeneity in Individual Frailty on the Dynamics of Mortality," *Demography*, 16, 439-454.
- Williams, J. S., and Lagakos, S. W. (1977), "Models for Censored Survival Analysis: Constant-Sum and Variable-Sum Models," *Biometrika*, 64, 215-224.