

# **Analysis of Reliability and Warranty Claims in Products with Age and Usage Scales**

J.F. Lawless

University of Waterloo

M.J. Crowder

Imperial College London

K.-A. Lee

University of Waterloo

## **Abstract**

Failures or other adverse events in systems or products may depend on the age and usage history of the unit. This paper presents models that may be used to assess the dependence on age or usage in heterogeneous populations of products, and shows how to estimate model parameters based on different types of observational field data. The setting where the events in question are warranty claims is examined in some detail. Applications to the analysis of automobile warranty data are considered and used to illustrate the methodology.

Key words and phrases: accelerated time models; alternative time scales; missing data; random effects; recurrent events.

# 1. INTRODUCTION

When considering the reliability of systems, it is often important to consider both the age of the system (that is, the time since it was introduced into service) and its cumulative usage, measured according to some specified variable. For example, a motor vehicle's usage may be measured by distance driven in miles or kilometers; laser printers' usage may be measured by cumulative pages printed. A key question which we address here is whether certain events (e.g. warranty claims) are primarily dependent on the age of the system, on cumulative usage, or both. Warranty plans for certain types of products also specify limits of coverage in terms of both age and usage, and it may be of interest to assess the effect of changes in the limits. Notable in this context are North American automobile warranties, which have age and distance limits for specific systems on the vehicle, for example 3 years and 36,000 miles.

The purpose of this paper is to study the estimation of product reliability through data on warranty claims. To do this we require models that allow events of interest to depend on both the age and usage of a product. We also have to deal with the fact that warranty data typically have incomplete information about product usage, as we discuss below. Robinson and McDonald (1991), Lawless and Kalbfleisch (1992) and Lawless (1998) give background on warranty data and other information on the field reliability of products.

Various authors have considered the analysis of accumulating warranty claims as a function of the age (time in service) of a product; for example, see Kalbfleisch et al. (1991) and, for a review of methodology, Lawless (1998). The estimation of failure rates or time to failure distributions from such data is more problematic, however (e.g. Suzuki, 1985, 1993; Lawless and Kalbfleisch, 1992; Lawless, 1998). Aside from data quality issues related to the timing and diagnosis of failures, the main problem is that cumulative usage is often given only at the times of claims, and so is generally unavailable for units that do not have at least one warranty claim. Thus, in a followup study of units under warranty we also do not

know the end-of-study usage for many units. (We assume here that usage is automatically recorded, as with distance driven in vehicles or number of copies for copiers or printers, but that data on usage are not automatically available to the analyst.) Sometimes dates of purchase of a product are also largely unknown; for example, dates of purchase for electronic goods or appliances tend to be reported only if a warranty claim arises. We do not deal with this complication here, and assume that purchase dates of all units are known. Finally, for cases where warranty limits are two-dimensional (involving both age and usage), we may not know the age  $\tau_i$  at which a unit's coverage ceased. For example, with a 3 year/36,000 mile warranty for automobiles, many vehicles reach the distance limit in less than 3 years and for a vehicle with no claims we will not know when this occurred.

Several authors have estimated age-specific warranty claim rates or expected numbers of claims by using estimates of the probability a product is still under warranty at age  $t > 0$  (e.g. Hu and Lawless, 1996a; Chukova and Robinson, 2006). A similar approach leads to estimates of usage-based rates or expected numbers of claims per unit. However, these methods make the strong assumption that the duration of followup for a unit is independent of its claims process. This is very often false, for example, when claim rates are related to the usage rate for a unit. A primary objective of this paper is to deal with this problem.

Another issue of importance to manufacturers is whether claims and repair events are driven primarily by age, by usage or by both (e.g. Chukova and Robinson, 2006; Krivtsov, 2006). Lawless et al. (1995), Murthy et al. (1995), Ahn et al. (1998) and Jung and Bai (2007) have considered models involving age and usage as dual time scales, but estimation based on field data or warranty data has not been considered. In particular, the case where age or usage at the end of followup are missing for units with no claims has not been investigated, except for the analysis of time to the first claim (Lawless et al., 1995). Other authors have discussed bivariate age and usage models for failure times in reliability settings (e.g. Singpurwalla and Wilson, 1998; Yang and Nachlas, 2001) but they do not consider recurrent events or estimation as we do here.

This paper thus makes a number of new contributions. In Section 2 we discuss recurrent event processes where age and usage are both involved and formulate random effects models that reflect heterogeneity across product units. In Section 3 we develop estimation procedures based on warranty claims data for some models in which usage accumulates linearly with age. Such models are important for automobiles and other products, and we show how to deal with missing followup times. Estimates of Hu and Lawless (1996a) and Chukova and Robinson (2006) for event rates in terms of age or usage are considered in Section 4, and biases in such estimates are examined. Both real and simulated warranty claims data are examined in Section 5, and we demonstrate the limitations of data where only a small fraction of units have claims. In Section 6 we show how to investigate the effects of changes to a warranty plan's age or usage limits and in Section 7 we indicate extensions to the methodology and areas that deserve further investigation.

## 2. MODELS FOR REPEATED EVENTS WITH AGE AND USAGE SCALES

There can be multiple warranty claims on a unit and we adopt recurrent event terminology and notation (Cook and Lawless, 2007). Let  $t \geq 0$  denote age (time since sale) of a product unit and let  $U_i(t)$  denote the usage at age  $t \geq 0$  for the  $i$ 'th unit in some population or sample. Let  $N_i(t)$  be the number of events (claims) experienced by the unit up to age  $t$ , so that  $\{N_i(t), t \geq 0\}$  is the counting process for events on the age scale. In practice, events may be split into different types and there will be a counting process for each type; we consider  $N_i(t)$  as the number of events of some specified type. The joint analysis of different types is mentioned in Section 7. There may also be explanatory variables  $x_i$  associated with a unit but for most of this paper we do not consider them explicitly; product units will instead be subdivided into groups, if necessary. This is consistent with other papers on warranty data

analysis, but we briefly discuss the inclusion of covariates in Section 7.

For general discussion we assume that  $N_i(0) = 0$ , although in some settings we might wish to recognize repairs made under warranty before a unit is sold by allowing  $N_i(0)$  to be positive. The usage process or “curve”  $U_i(t)$  is non-decreasing for  $t \geq 0$ , and we assume that the usage process is external in the sense of Kalbfleisch and Prentice (2002, Section 6.3). That is, the usage curve is not influenced by the event process, and so we may condition on the usage history for a given unit, effectively treating it as a covariate. The assumption that  $U_i(t)$  is external is a reasonable approximation to reality while products are under warranty. Strictly speaking, units cannot be used while they are being repaired, but in the contexts we consider repair times are short. Non-external usage processes are more difficult to handle, and previous papers have not addressed this at all; we comment briefly on this topic in Section 7.

To specify models for repeated events, we introduce the additional notation  $\overline{N}_i(t) = \{N_i(s) : 0 \leq s \leq t\}$  and  $\overline{U}_i(t) = \{U_i(s) : 0 \leq s \leq t\}$ . We assume that there are specified warranty coverage limits  $T_0 > 0$  and  $U_0 > 0$  such that an event at age  $t$  is recorded only if  $t \leq T_0$  and  $U_i(t) \leq U_0$ . Either  $T_0$  or  $U_0$  may be infinite but typically at least  $T_0$  is finite. For convenience we will denote  $\overline{N}_i(T_0)$  as  $\overline{N}_i$  and  $\overline{U}_i(T_0)$  as  $\overline{U}_i$ , since we only consider  $t$  in the range  $(0, T_0)$ . For general discussion we treat events as being observed in continuous time but we will later also discuss discrete time models, since claims are generally recorded in terms of days.

An event process in continuous time can be modeled by specifying an intensity function  $\lambda(t; \overline{U}_i, \overline{N}_i(t-))$  that gives the instantaneous probability of an event at time  $t$ , conditional on  $\overline{U}_i$  and the history  $\overline{N}_i(t-)$  of previous events (e.g. see Cook and Lawless, 2007, Chapter 2). We make the reasonable assumption that  $\lambda(t; \overline{U}_i, \overline{N}_i(t-))$  depends on  $\overline{U}_i$  only through  $\overline{U}_i(t)$ . In addition, we will require a probability model for the usage paths  $\overline{U}_i$  in the population of units. An important feature of product usage is heterogeneity; usage paths tend to vary widely across units. In addition, unmeasured usage, environmental, or product quality

factors often create heterogeneity in the occurrence of failures or warranty claims across units. To accommodate such heterogeneity we introduce a random effect  $Z_i$  (perhaps a vector), and allow  $\bar{U}_i$  and the claim intensity to depend on  $Z_i$ . The  $Z_i$  are assumed to be independent and identically distributed across units, with distribution function  $G(\cdot)$ .

In this paper we assume that the event process is conditionally Poisson, with intensity function  $\lambda(t|\bar{U}, Z)$ , and that the intensity at time  $t$  depends on  $\bar{U}_i$  only through  $\bar{U}_i(t)$ . That is,

$$\lim_{\Delta t \downarrow 0} \frac{\Pr \{ \text{event in } [t, t + \Delta t) | \bar{U}_i, \bar{N}_i(t-), Z_i \}}{\Delta t} = \lambda(t|\bar{U}_i, Z_i) = \lambda(t|\bar{U}_i(t), Z_i),$$

where  $Z_i$  may also affect the distribution of  $\bar{U}_i$ . This model is very flexible, and allows for both heterogeneity in usage paths and in the event intensities given  $\bar{U}_i$ . If the process for unit  $i$  is observed from age 0 to some age  $\tau_i$  that is independent of the event process, it then follows that (e.g. Cook and Lawless, 2007, Ch. 2) the joint probability density for  $\bar{U}_i$  and the event history

$$\bar{N}_i(\tau_i) = \{n_i \text{ events over } 0 \leq t \leq \tau_i, \text{ at times } t_{ij}(j = 1, \dots, n_i)\} \quad (1)$$

is, conditional on  $Z_i$ ,

$$\left\{ \prod_{j=1}^{n_i} \lambda(t_{ij}|\bar{U}_i, Z_i) \right\} \exp \{ -\Lambda(\tau_i|\bar{U}_i, Z_i) \} p \{ \bar{U}_i(\tau_i) | Z_i \}, \quad (2)$$

where  $p(\bar{U}_i|Z_i)$  denotes the density of  $\bar{U}_i$  given  $Z_i$  and

$$\Lambda(\tau_i|\bar{U}_i, Z_i) = \int_0^{\tau_i} \lambda(t|\bar{U}_i, Z_i) dt \quad (3)$$

is the conditional expected number of events,  $E\{N_i(\tau_i)|\bar{U}_i, Z_i\}$ . The density function for the observable data  $\{\bar{N}_i(\tau_i), \bar{U}_i(\tau_i)\}$  is obtained by integrating (2) with respect to the distribution of  $Z_i$ .

There are some field reliability studies for which  $\tau_i$  is prespecified and essential data on  $\bar{U}_i$  are collected, and in such cases maximum likelihood estimation can be based on a likelihood function which is a product of terms (2) for a random sample of independent units. Indeed,

current technology allows usage curves  $\bar{U}_i$  to be recorded and stored for many products, and data of this type are expected to become increasingly common. However, in this paper we consider the more challenging setting exemplified by warranty data, and show how to estimate model parameters when  $\tau_i$  may depend on  $\bar{U}_i$ , or when data on  $\tau_i$  or  $\bar{U}_i$  are missing. In that case (2) is unavailable to us.

Even with the conditional Poisson assumption, there are many models for the conditional event intensity  $\lambda(t|\bar{U}, Z)$  and usage path processes  $\bar{U} = \{U(t), 0 \leq t \leq T_0\}$  that one might consider. There has been very little work on repeated events in this context; Crowder and Lawless (2007) provide a review, but we mention a few pertinent references. Lawless et al. (1995) mention an accelerated time model for which  $U_i(t) = Z_i t$  and  $\lambda(t|\bar{U}_i, z_i) = z_i^\beta \lambda_0(z_i^\beta t)$ , where  $\lambda_0(t)$  is a parametrically specified function. They investigate estimation for a corresponding failure time model, but not for the repeated events model. Murthy et al. (1995) consider a model with  $U_i(t) = Z_i t$  and  $\lambda(t|\bar{U}_i, z_i) = \alpha_0 + \alpha_1 z_i + \alpha_2 t + \alpha_3 z_i t$ , but do not address estimation. Other work has tended to be for the case of a single failure time  $T_i$  rather than recurrent events. Lawless et al. (1995) and Jung and Bai (2007) consider models with  $U_i(t) = Z_i t$  which assume  $T_i$  depends on  $Z_i$  but is independent of  $\bar{U}_i$ , given  $Z_i$ . Singpurwalla and Wilson (1998) take  $\bar{U}_i$  as the path of some stochastic process, and the hazard function for  $T_i$  given  $\bar{U}_i$  is of the form  $h(t|U_i(t))$ ; a random effect  $Z_i$  is not used.

Crowder and Lawless (2007) discuss pros and cons of different modeling strategies. In general, most models that involve stochastic processes for  $\bar{U}_i$  are difficult to fit to data, and for the types of incomplete data considered here, are very challenging. Consequently we consider in the next section models for which  $U_i(t)$  is a deterministic function of  $Z_i$  and  $t$ , with  $Z_i$  incorporating variability in usage paths across different units. They provide a good approximation to many warranty or field reliability settings.

### 3. A MODEL WITH RANDOM USAGE RATES

#### 3.1 A Family of Models

Families of models that have several advantages in the present setting are those where  $U_i(t) = U(t; Z_i, \psi)$  is a deterministic function of  $t$ , given the random effect  $Z_i$ ; the parameter  $\psi$  allows additional flexibility in the shapes of usage curves. In addition, it is assumed that  $\lambda(t|\bar{U}_i, Z_i) = \lambda(t|Z_i)$ ; this is especially convenient when data on  $\bar{U}_i$  is limited and, in some cases, entirely missing. In this paper we restrict attention to models where  $Z_i > 0$  is a random usage rate, with

$$U_i(t) = Z_i t, \quad t \geq 0 \quad (4)$$

and to start, we consider the Poisson model with conditional rate function

$$\lambda(t|\bar{U}_i, Z_i) = \lambda(t|Z_i) = Z_i^\beta \lambda_0(tZ_i^\beta; \alpha), \quad (5)$$

where  $\lambda_0(t; \alpha)$  is a baseline rate function specified in terms of a parameter vector  $\alpha$ . The  $Z_i$  are assumed independent with some distribution function  $G(Z; \gamma)$ . This model was mentioned by Lawless et al. (1995), but estimation was not pursued. The linear usage model (4) provides a good approximation for equipment such as motor vehicles over the early part of their time in service, and has often been used in dealing with warranty data (e.g. Lawless et al., 1995; Murthy et al., 1995; Chukova and Robinson, 2006; Jung and Bai, 2007).

The form of (5) is analogous to an accelerated failure time model. It is flexible and has an appealing property: the conditional expected number of events up to age  $t$  for unit  $i$  is

$$\begin{aligned} E\{N_i(t)|Z_i\} &= \int_0^t \lambda(s|Z_i) ds \\ &= \Lambda_0(tZ_i^\beta; \alpha), \end{aligned} \quad (6)$$

where

$$\Lambda_0(t; \alpha) = \int_0^t \lambda_0(s; \alpha) ds$$



is a baseline cumulative mean function. Thus, if  $\beta = 0$  the expected number of events is independent of the usage rate, whereas if  $\beta = 1$  then (6) equals  $\Lambda_0(U_i(t); \alpha)$  and depends only on the accumulated usage up to age  $t$ . For other values of  $\beta$ , event occurrence depends on both age and usage; note that  $\Lambda_0(tz^\beta)$  can be rewritten as  $\Lambda_0(t^{1-\beta}u(t)^\beta)$ . This model thus allows an examination of the important question as to whether failures are predominantly a function of age, usage, or both variables. Other models should, of course, be considered if (5) does not fit, and in the next subsection we extend (5) to allow for the frequently occurring phenomenon of extra – Poisson variation. Some other alternatives are mentioned in Section 7.

### 3.2 Maximum Likelihood Estimation

First, we give the contribution to the likelihood function for individual unit  $i$  in the “ideal” setting where it is observed up to a prespecified age  $\tau_i$ , and where  $Z_i$  is observed because the cumulative usage is observed at one or more times. The usage paths are not actually precisely linear in practice due to short term variations in the rates, so some pragmatic convention is usually needed to define the “observed”  $Z_i$ . We assume the additional variability introduced is small, and ignore it here; it is possible, but more complicated, to consider variation in  $Z_i$  measurements as a form of covariate measurement error. We adopt the convention of basing  $Z_i$  on the usage at the largest observation time.

The data on the  $i$ 'th unit is of the form (1), and by (2), (5) and (6) the likelihood function for  $\beta$  and  $\alpha$  based on  $M$  independent units is

$$L(\theta) = \prod_{i=1}^M \prod_{j=1}^{n_i} z_i^\beta \lambda_0 \left( t_{ij}; z_i^\beta; \alpha \right) e^{-\Lambda_0(\tau_i z_i^\beta; \alpha)}, \quad (7)$$

where  $\theta = (\alpha, \beta)$  is the vector of unknown parameters. In practice, parametric forms such as  $\lambda_0(t; \alpha) = \alpha_1 \alpha_2 (\alpha_1 t)^{\alpha_2 - 1}$  or  $\exp(\alpha_1 + \alpha_2 t)$  are useful. In this setting,  $Z_i$  functions as a covariate and its distribution can be estimated separately from the data  $z_1, \dots, z_M$  on the random sample of  $M$  units.

We now consider the common situation arising with warranty data and certain other

observational field data, where  $Z_i$  is observed only at the time of a repair or warranty claim; for units with no observed events the value of  $Z_i$  is consequently unknown. In addition,  $\tau_i$  may depend on  $Z_i$ , and so may also be missing. We consider the following set up. The elapsed time in service for the  $i$ 'th unit when the data are analyzed is denoted by  $T_i$  and we let  $\tau_i^* = \min(T_i, T_0)$ . If a usage limit  $U_0 < \infty$  is in effect, then  $U_0/Z_i$  is the time at which usage reaches  $U_0$ , so the total followup time for the unit under warranty claims reporting is  $\tau_i = \min(\tau_i^*, U_0/Z_i)$ . If  $N_i(\tau_i) > 0$  then  $Z_i$  is observed, because  $U_i(t) = Z_i t$  is recorded at the time of any warranty claims; in this case  $\tau_i$  also becomes known. If  $N_i(\tau_i) = 0$ , however, the value of  $Z_i$  is unknown and if  $U_0 < \infty$ , then  $\tau_i$  is also unknown.

For convenience, we label the units for which  $N_i(\tau_i) > 0$  as  $i = 1, \dots, m$  and those for which  $N_i(\tau_i) = 0$  as  $i = m + 1, \dots, M$ . The likelihood function is then (see (2) and also Lawless et al., 1995, Section 3.2)

$$L(\theta) = \prod_{i=1}^m \left\{ \prod_{j=1}^{n_i} z_i^\beta \lambda_0(t_{ij} z_i^\beta) \right\} e^{-\Lambda_0(\tau_i z_i^\beta)} g(z_i) \prod_{i=m+1}^M \int_0^\infty e^{-\Lambda_0(\tau_i z_i^\beta)} g(z_i) dz_i \quad (8)$$

where  $n_i = N_i(\tau_i)$ ,  $\lambda_0(t) = \lambda_0(t; \alpha)$ ,  $g(z) = g(z; \gamma)$  is the density function for  $Z_i$  and  $\theta = (\alpha, \beta, \gamma)$ . Note that in the terms for  $i = m + 1, \dots, M$  in (8),  $\tau_i = \min(\tau_i^*, U_0/z_i)$ , so that the integral can be written as

$$\int_0^{U_0/\tau_i^*} e^{-\Lambda_0(\tau_i^* z_i^\beta)} g(z_i) dz_i + \int_{U_0/\tau_i^*}^\infty e^{-\Lambda_0(U_0 z_i^{\beta-1})} g(z_i) dz_i. \quad (9)$$

In many contexts recurrent events exhibit extra-Poisson variation, even after conditioning on covariates and usage rates. This is generally due to heterogeneity in the users of different units and in the environment where they operate. To allow for this additional variability we consider a mixed Poisson process (Lawless, 1987) by extending (5) to

$$\lambda(t|\bar{U}_i, z_i, v_i) = v_i z_i^\beta \lambda_0(t z_i^\beta; \alpha), \quad (10)$$

where the  $v_i$  are independent and identically distributed random variables with mean 1 and variance  $\phi$ , and are independent of the  $Z_i$ . We assume for convenience, and with little loss

of flexibility, that the  $v_i$  have a gamma distribution. It is shown in Appendix 1 that in this case the likelihood function analogous to (8) is

$$L(\theta) = \prod_{i=1}^m \left\{ \prod_{j=1}^{n_i} z_i^\beta \lambda_0(t_{ij} z_i^\beta) \right\} \omega^\omega \frac{\Gamma(n_i + \omega)}{\Gamma(\omega)} \left\{ \omega + \Lambda_0(z_i^\beta \tau_i) \right\}^{-n_i - \omega} g(z_i) \\ \times \prod_{i=m+1}^M \int_0^\infty \left\{ 1 + \omega^{-1} \Lambda_0(z_i^\beta \tau_i) \right\}^{-\omega} g(z_i) dz_i, \quad (11)$$

where  $\omega = \phi^{-1}$  and  $\theta = (\alpha, \beta, \omega)$ .

The likelihoods (8) and (11) are most easily maximized by using general optimization software to maximize  $\log L(\theta)$ ; good software does not require formulas for derivatives (which can be messy) and will return a Hessian matrix  $H(\hat{\theta}) = (\partial^2 \log L(\theta) / \partial \theta \partial \theta')_{\hat{\theta}}$ , from which variance estimates for  $\hat{\theta}$  can be obtained. The possibility of estimating  $g(z)$  parametrically from (9) is discussed later, but frequently there is a good estimate, obtained from other sources (Lawless et al., 1995; Chukova and Robinson, 2006). As a result, it is often assumed that  $g(z)$  is known, although it will have been estimated.

## 4. MARGINAL RATE AND MEAN FUNCTIONS

### 4.1 Marginal Rate Functions and Bias in Naive Estimators

It is often of interest to estimate the marginal mean functions in terms of age and usage, which we denote respectively as

$$\Lambda_a(t) = E \{N_i(t)\}, \quad \Lambda_u(u) = E \{N_i^u(u)\},$$

where  $N_i(t)$  is the number of events on unit  $i$  up to age  $t$ , and  $N_i^u(u)$  is the number of events up to usage level  $u$ . These are related to the quantities in Section 2 by

$$\Lambda_a(t) = E_{Z_i, \bar{U}_i} \{E [N_i(t) | Z_i, \bar{U}_i]\} \quad (12)$$

$$\Lambda_u(u) = E_{Z_i, \bar{U}_i} \{E [N_i(U_i^{-1}(u)) | Z_i, \bar{U}_i]\}, \quad (13)$$

and can be estimated by working with models such as those in Section 3. The associated rate functions  $\lambda_a(t) = \Lambda'_a(t)$  and  $\lambda_u(u) = \Lambda'_u(u)$  can also be estimated.

Several authors have estimated  $\Lambda_a(t)$  and  $\Lambda_u(u)$  nonparametrically from warranty data by using estimates of the probability a product unit is still under warranty at age  $t$  or usage  $u$  to adjust raw warranty claim counts. See in particular Hu and Lawless (1996a) and Chukova and Robinson (2006). However, these authors assume that censoring or end-of-followup times are independent of the event processes. This assumption is violated when there is a warranty usage limit and events depend on both age and usage; this makes the proposed estimates biased. In this section we examine this bias in the setting of Section 3. We assume that data are of the form that leads to the likelihoods (8) and (11), that is, data on unit  $i$  are from age 0 to age  $\tau_i = \min(T_i, T_0, U_0/Z_i)$ . In addition, the usage rate  $Z_i$  is observed only if there is at least one claim over  $(0, \tau_i)$ . Let  $Y_i^a(t) = I(\tau_i \geq t) = I(T_i \geq t)I(U_0/Z_i \geq t)$ , where  $I(A)$  is the indicator function for an event  $A$  and we restrict attention to values  $t \leq T_0$ . Note that  $Y_i^a(t)$  indicates whether a unit is “at risk” of producing an observed event at age  $t$ .

Hu and Lawless (1996a) and Chukova and Robinson (2006) propose estimates of  $\lambda_a(t)$  in the discrete time case,

$$\hat{\lambda}_a(t) = \sum_{i=1}^M Y_i^a(t) n_i(t) / \sum_{i=1}^M \hat{P}_i^a(t), \quad t = 1, 2, \dots, T_0 \quad (14)$$

where  $n_i(t)$  is the number of claims at age  $t$  (typically measured in days) for unit  $i$ , and  $\hat{P}_i^a(t)$  is an estimate of  $E\{Y_i^a(t)\}$ , that is,

$$P_i^a(t) = \Pr(\tau_i \geq t) = \Pr(T_i \geq t) \Pr(Z_i \leq U_0/t). \quad (15)$$

It is assumed here that  $T_i$  is independent of  $Z_i$ . Often the  $T_i$  are known for all units and then

$$\hat{P}_i^a(t) = I(T_i \geq t) \hat{G}(U_0/t),$$

where  $\hat{G}(z)$  is an estimate of  $G(z) = \Pr(Z_i \leq z)$ . If  $Y_i^a(t)$  and  $n_i(t)$  are independent then  $E\{Y_i^a(t)n_i(t)\} = P_i^a(t)\lambda_a(t)$  and if  $M^{-1} \sum \hat{P}_i^a(t)$  is a consistent estimate of  $M^{-1} \sum P_i^a(t)$

then  $\hat{\lambda}_a(t)$  is a consistent estimate of  $\lambda_a(t)$ . However,  $Y_i^a(t)$  and  $n_i(t)$  are not in general independent, if events depend on both age and usage. Theorem 1 below gives the asymptotic bias of the estimates  $\hat{\lambda}_a(t)$  in that case.

Similarly, with usage expressed in discrete units, an analogous estimate of  $\lambda_u(u)$  is

$$\hat{\lambda}_u(u) = \sum_{i=1}^M Y_i^u(u) n_i^u(u) / \sum_{i=1}^M \hat{P}_i^u(u), \quad u = 1, 2, \dots, U_0 \quad (16)$$

where  $n_i^u(u)$  is the number of claims at usage  $u$  for unit  $i$ ,  $Y_i^u(u) = I(Z_i \tau_i^* \geq u)$  with  $\tau_i^* = \min(T_i, T_0)$ , and  $\hat{P}_i^u(u)$  is an estimate of

$$P_i^u(u) = \Pr(Z_i \geq u/\tau_i^*). \quad (17)$$

If  $Y_i^u(u)$  and  $n_i^u(u)$  are independent and  $M^{-1} \sum \hat{P}_i^u(u)$  is a consistent estimate of  $M^{-1} \sum P_i^u(u)$  then  $\hat{\lambda}_u(u)$  is a consistent estimate of  $\lambda_u(u)$ .

The following theorems, which are proved in Appendix 2, give the asymptotic biases for  $\hat{\lambda}_a(t)$  and  $\hat{\lambda}_u(u)$  when censoring times are not independent of the event processes. For convenience we revert to treating time and usage as continuous but the results also apply in the case of discrete time and usage scales.

**Theorem 1** Suppose that the event rate function  $\lambda(t|Z_i, \bar{U}_i)$  is of the form  $\lambda(t|Z_i)$  as in Section 3, that  $U_i(t)$  is given by (4), and that  $Z_i$  is independent of  $T_i$  with distribution function  $G(z)$ . Then the estimates  $\hat{\lambda}_a(t)$  and  $\hat{\lambda}_u(u)$  based on independent units  $i = 1, \dots, M$ , and given by (14) and (16), converge in probability to  $B_a(t)\lambda_a(t)$  and  $B_u(u)\lambda_u(u)$ , where

$$B_a(t) = \frac{E \left\{ \lambda(t|Z) \mid Z \leq U_0/t \right\}}{E \left\{ \lambda(t|Z) \right\}}, \quad (18)$$

$$B_u(u) = \frac{E \left\{ Z^{-1} \lambda(u/Z|Z) I(Z \geq u/\tau^*) \right\}}{E \left\{ Z^{-1} \lambda(u/Z|Z) \right\} \Pr(Z \geq u/\tau^*)}. \quad (19)$$

This result holds for both of models (5) and (10). In (19), the expectation in the numerator is with respect to both  $Z$  and  $T_i$  in  $\tau^* = \min(T_i, T_0)$ , and  $\Pr(Z \geq u/\tau^*)$  likewise treats the  $T_i$  as random. If all units are observed long enough that  $T_i \geq T_0$ , then  $\tau_i^* = T_0$  for all  $i = 1, \dots, M$  and (19) reduces to the simpler form (A5) in Appendix A. Note that if  $U_0 = \infty$ , (18) equals one and  $\hat{\lambda}_a(t)$  is consistent, as seems intuitively obvious when there is no usage limit. Even when  $T_0 = \infty$  (which rarely is the case in practice), however, (19) does not equal one because of finite limitations on the followup times  $T_i$ .

In many settings a model of the form (5), with  $\lambda_0(t)$  of the “power law” form  $\alpha_1 \alpha_2 t^{\alpha_2 - 1}$  has been found reasonable. The following theorem gives the asymptotic bias functions (18) and (19) for this model, and we illustrate its use in the next section.

**Theorem 2** Suppose that in addition to the conditions in Theorem 1, the conditional rate function  $\lambda(t|Z)$  is of the form  $Z^\beta \lambda_0(tZ^\beta)$  given in (5), with  $\lambda_0(t) = \alpha_1 \alpha_2 t^{\alpha_2 - 1}$ , where  $\alpha_1 > 0$ ,  $\alpha_2 > 0$ . Then (18) and (19) become

$$B_a(t) = \frac{E \{ Z^{\beta \alpha_2} | Z \leq U_0/t \}}{E \{ Z^{\beta \alpha_2} \}}, \quad (20)$$

$$B_u(u) = \frac{E \{ Z^{(\beta-1)\alpha_2} I(Z \geq u/\tau^*) \}}{E \{ Z^{(\beta-1)\alpha_2} \} \Pr \{ Z \geq u/\tau^* \}}. \quad (21)$$

## 4.2 An Illustration of Naive Estimator Bias

The expressions (20) and (21) take simple forms when the usage rates follow a lognormal distribution or a gamma distribution. To illustrate the gamma case, suppose that  $Z_i \sim \text{Gamma}(a, b)$ , which denotes the distribution with mean  $ab$ , variance  $a^2b$ , and density function

$$g(z; a, b) = \frac{1}{a^b \Gamma(b)} z^{b-1} e^{-z/a} \quad z > 0.$$

Straightforward calculations then show that (20) and (21) become

$$B_a(t) = \frac{G(U_0/t; a, b + \alpha_2 \beta)}{G(U_0/t; a, b)}, \quad (22)$$

$$B_u(u) = \frac{E \{ \bar{G}(u/\tau^*; a, b + \alpha_2(\beta - 1)) \}}{E \{ \bar{G}(u/\tau^*; a, b) \}}, \quad (23)$$

where  $G(z; a, b)$  is the distribution function corresponding to  $g(z; a, b)$ . In (23),  $\bar{G}(z; a, b) = 1 - G(z; a, b)$  and the expectation is with respect to the distribution of  $\tau^* = \min(T_i, T_0)$ . If  $T_i \geq T_0$  for all  $i$  then  $\tau^* = T_0$  and the expectation disappears.

For illustration we take a setting similar to one involving car data discussed in Section 5. We let  $t$  represent years and  $u$  thousands of miles driven, and consider warranty limits  $T_0 = 3$ ,  $U_0 = 36$ . The usage rates  $Z_i$  (in thousands of miles accumulated per year) are taken to be Gamma (2.5, 5.5), which gives  $E(Z_i) = 13.75$  and  $\text{Var}(Z_i) = 28.875$ . Finally, we assume that  $\lambda_0(t) = \alpha_1 \alpha_2 t^{\alpha_2 - 1}$ ; note that (22) and (23) do not depend on the value of  $\alpha_1$ .

Note that when  $\beta = 0$ ,  $B_a(t) = 1$ , that is, the naive marginal estimate (14) for  $\lambda_a(t)$  is unbiased asymptotically, or consistent. This is because when  $\beta = 0$ , the rate function (5) does not depend on  $Z_i$ , so that  $E\{n_i(t)|Z_i\} = E\{n_i(t)\}$ . When  $\beta = 1$ , (23) gives  $B_u(u) = 1$ , so that the naive estimate (16) for  $\lambda_u(u)$  is consistent. This is a consequence of the fact that when  $\beta = 1$  the rate function (5) gives  $\lambda(t|Z) = Z\lambda_0(tZ)$  and thus (see Appendix 2)

$$\lambda_u(u|Z) = Z^{-1}\lambda((uZ^{-1})Z|Z) = \lambda_0(u).$$

Hence  $E\{n_i^u(u)|Z_i\} = E\{n_i^u(u)\}$ , that is, the expected number of events as a function of usage is independent of the rate of usage  $Z_i$ .

Table 1 gives values of  $B_a(t)$  and  $B_u(u)$  for  $\beta = 0, 0.5, 1$  and  $\alpha_2 = 1, 1.1$ ; these represent a range of scenarios encountered with car warranty claims. We show results for the case where  $T_i = T_0 = 3$  for all units, so that  $\tau^* = 3$  in (23) and the expectations disappear. The bias functions  $B_a(t)$  with  $\beta > 0$  and  $B_u(u)$  with  $\beta < 1$  increase with  $t$  and  $u$  respectively; this is because at higher ages and usage (miles) there is a greater chance that a vehicle will no longer be covered by the warranty. This results in a selection effect against cars with high usage rates in the estimation of  $\lambda_a(t)$  and a selection effect against cars with low usage rates in the estimation of  $\lambda_u(u)$ . Each of these effects produces underestimation of the event rates in question. The key message is that if events are primarily a function of age ( $\beta$  close to 0)

then  $\lambda_u(u)$  is underestimated, and that if events are primarily a function of usage ( $\beta$  close to 1) then  $\lambda_a(t)$  is underestimated. In practice we should focus on  $\hat{\lambda}_a(t)$  in the former case and  $\hat{\lambda}_u(u)$  in the latter, and they are not substantially biased. However, this highlights the importance of knowing whether particular types of claims are mainly age- or usage-related. If one merely computes  $\hat{\lambda}_a(t)$  and  $\hat{\lambda}_u(u)$ , there is no indication as to the extent of bias in either estimate.

**Table 1.** Bias Functions (22) and (23) of Estimators (14) and (16)

$\beta$	$\alpha_2$	$B_a(t)$			$B_u(u)$		
		$t = 1$	$t = 2$	$t = 3$	$u = 12$	$u = 24$	$u = 36$
0	1.0	1	1	1	0.968	0.827	0.677
	1.1	1	1	1	0.962	0.806	0.645
0.5	1.0	0.998	0.918	0.805	0.988	0.923	0.840
	1.1	0.998	0.910	0.787	0.987	0.915	0.824
1.0	1.0	0.995	0.829	0.632	1	1	1
	1.1	0.995	0.811	0.600	1	1	1

## 5. APPLICATIONS TO CAR WARRANTY DATA

Automobile warranty data are a primary area of application for the methods here. In Section 5.2 we consider a set of warranty claims data for 44,890 cars of one type and model year, and manufactured in one plant during 2000-2001. The warranty limits were 3 years and 36,000 miles.

The fraction of vehicles that generate at least one warranty claim is small for most car systems. This can create difficulties in fitting models, even when the total number of vehicles  $M$  is large. Some parameter estimates may be highly correlated, and confidence limits for certain parameters may be wide. In addition, the information about the distribution of  $Z_i$  is limited and it may be important to have a good external estimate of  $g(z)$ . To illustrate the



effect of  $m/M$  on estimation, we consider simulated data sets, following which we examine the real warranty data.

### 5.1 Simulated Data Sets

To assess the amount of information about various model features, we simulated data on  $M = 1000$  vehicles under scenarios that give different expected values for  $m$ , the number of vehicles with one or more warranty claims within the age and mileage limits. To reflect real settings, we generated data under the following assumptions, with age  $t$  expressed in years and mileage  $u_i(t)$  in thousands of miles. The parameter values are realistic in car warranty settings where distance is the key factor (Lawless et al., 1995).

- (i) Mileage accumulation (usage) is given by (4), with  $\log Z_i \sim N(2.4, 0.58^2)$ .
- (ii) Each vehicle has a random effect  $v_i$  as in (10), which follows a gamma distribution with mean 1 and variance  $\phi = 1.0$ .
- (iii) Given  $v_i$  and  $Z_i$ , events (warranty claims) follow a nonhomogeneous Poisson process with intensity function of the form (10), with  $\beta = 0.95$  and

$$\lambda_0(t) = (\alpha_2/\alpha_1) (t/\alpha_1)^{\alpha_2-1}, \quad t > 0$$

with  $\alpha_2 = 1.10$  and  $\alpha_1$  taking different values which control the value of  $E(m)$ .

Estimation was based on maximization of the likelihood functions (11), and two cases were considered:

- (a)  $g(z)$  is assumed known, with  $\log Z \sim N(2.4, 0.58^2)$ , and  $\theta = (\alpha_1, \alpha_2, \beta, \phi)$  is estimated.
- (b)  $g(z)$  is assumed unknown, but it is assumed that  $\log Z_i \sim N(\mu, \sigma^2)$ . The parameter vector  $\theta = (\alpha_1, \alpha_2, \beta, \phi, \mu, \sigma)$  is estimated.

Estimation was carried out by maximizing  $\ell(\theta) = \log L(\theta)$  given by (11) using the general purpose optimization function `nlm` in *R*. Other general optimization functions, such as `NLP`

in SAS, work equally well. An advantage of these types of functions is that it is necessary only to provide code for  $\ell(\theta)$ , and not derivatives. The software can determine any needed derivatives numerically, and will also return a Hessian matrix  $H(\hat{\theta}) = (\partial^2 \ell(\theta) / \partial \theta \partial \theta')_{\theta=\hat{\theta}}$ , so that  $(-H(\hat{\theta}))^{-1}$  is an estimated asymptotic covariance matrix for  $\hat{\theta}$ . An adjustment that decreases the correlation between parameter estimates for  $\beta$  and  $\alpha_1$  is to re-parameterize  $\alpha_1$  as  $\alpha_{1c} = (\alpha_1 / 12^\beta)$  and scale  $Z_i$  as  $Z_{ci} = Z_i / 12$ ; we used  $(\alpha_{1c}, \alpha_2, \beta)$  in our procedures, and also  $\omega = \phi^{-1}$  in place of  $\phi$ .

The total number of vehicles ( $M = 1000$ ) is smaller than for many warranty data bases but large enough to indicate the relative amounts of information about different parameters. We show results for two individual samples, generated using values  $\alpha_1 = 500$  and  $\alpha_1 = 80$ . The first sample gave  $m = 38$  and the second  $m = 229$  vehicles with at least one claim. The ratios  $m/M$  of 0.038 and 0.229 for sample 1 and 2 are plausible, with the latter being at the high end of what is typically observed in practice for a vehicle system. Table 2 gives estimates and their standard errors, and a number of features stand out: with  $\mu, \sigma$  estimated and  $m/M$  small, the parameter  $\beta$  is imprecisely estimated, and sheds limited light on whether age or distance driven is more important; even when  $\mu$  and  $\sigma$  are known,  $\beta$  is rather imprecisely estimated when  $m/M$  is small;  $\mu$  is also imprecisely estimated when  $m/M$  is small;  $\omega$  (and  $\phi$ ) is imprecisely estimated, especially when  $m/M$  is small.

These patterns persist across other scenarios, including ones where  $\beta$  is close to zero and across a range of sample sizes. More detailed results from the simulated samples indicate that when  $m/M$  is small,  $\beta$  and  $\mu$  are rather highly confounded, giving high asymptotic correlation for  $\hat{\beta}$  and  $\hat{\mu}$ . The key problem is that when  $m/M$  is small (values of .10 or less are common for most car systems with three year warranties), and  $Z_i$  is unobserved unless  $n_i > 0$ , there is limited information about the effect (represented by  $\beta$ ) of the mileage rate on claims. A good external estimate of the mileage rate distribution  $g(z)$  improves the situation, but as we will see when considering the warranty data of the next section, the estimate of  $\beta$  depends heavily on  $g(z)$ , and so is susceptible to its misspecification. On the

other hand, if it is possible to observe  $Z_i$  for all vehicles, then estimation of  $g(z)$  and of  $\beta$  improves dramatically. This may not be feasible in general, but two other possibilities exist. One is to use the  $Z_i$  from all vehicles having any type of claims; this could be done by using models like those here to represent claims without reference to type. A second possibility is to obtain  $Z_i$  values for some of the vehicles with  $n_i = 0$  through a supplementary sample, and to extend methods of Hu and Lawless (1996b).

**Table 2.** Parameter estimates from two simulated samples generated with parameter values  $\beta = 0.95$ ,  $\alpha_2 = 1.1$ ,  $\omega = 1.0$ ,  $\mu = 2.4$ ,  $\sigma = 0.58$ , and (1)  $\alpha_1 = 500$  (2)  $\alpha_1 = 80$ . In each case  $M = 1000$  and the numbers of cars with one or more claims were (1)  $m = 38$ , (2)  $m = 229$ .

$m$	$M$	Parameter	True	(a) $\mu, \sigma$ known		(b) $\mu, \sigma$ estimated	
			Value	Est.	S.E.	Est.	S.E.
38	1000	$\beta$	0.95	0.86	0.28	1.06	3.15
		$\alpha_{1c}$	47.2	39.9	18.4	38.4	26.5
		$\alpha_2$	1.10	1.15	0.18	1.16	0.19
		$\omega$	1.00	3.92	18.0	6.21	109.0
		$\mu$	2.40	-	-	2.32	1.27
		$\sigma$	0.58	-	-	0.59	0.09
229	1000	$\beta$	0.95	0.93	0.13	1.08	0.28
		$\alpha_{1c}$	7.55	8.04	0.76	7.65	0.81
		$\alpha_2$	1.10	1.03	0.06	1.04	0.06
		$\omega$	1.00	1.66	0.63	2.09	1.13
		$\mu$	2.40	-	-	2.32	0.11
		$\sigma$	0.58	-	-	0.63	0.04

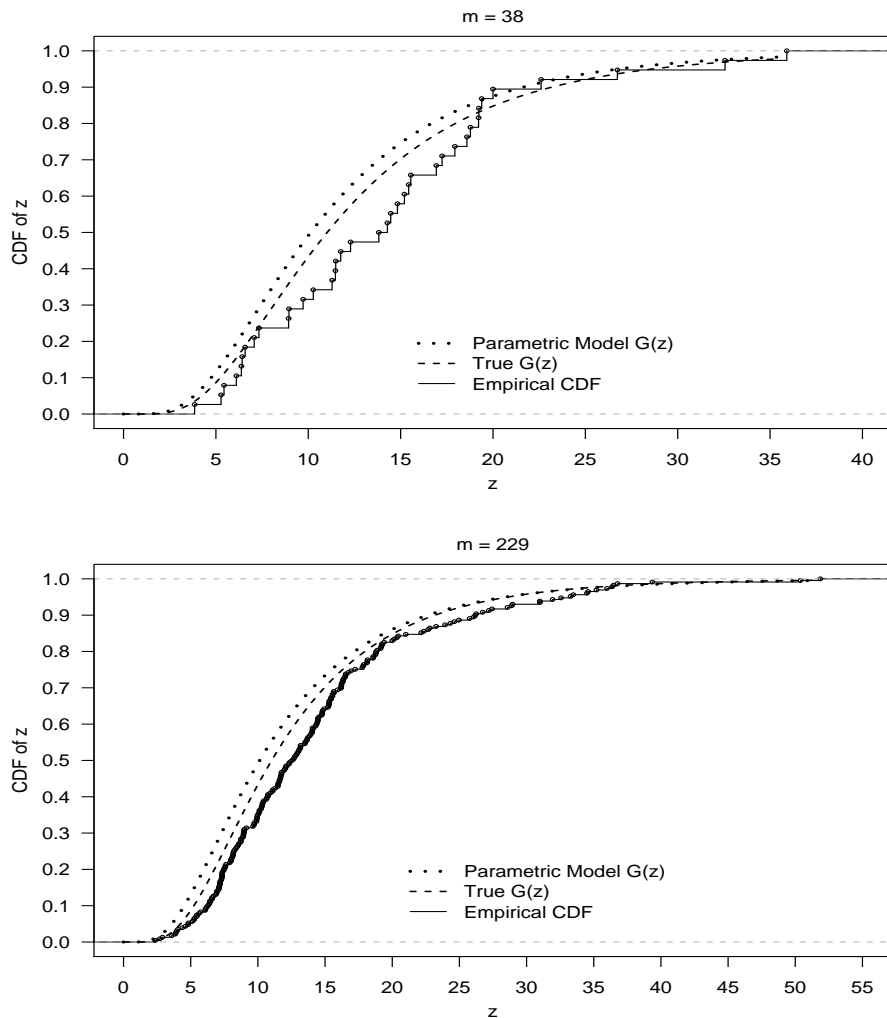


Figure 1. Estimated and true distribution functions  $G(z)$  for two simulated samples

Finally, it is possible to develop a nonparametric estimate of  $G(z)$ , as an alternative to the parametric estimates considered here; this is described in Appendix 3. It should be noted that a “naive” estimate obtained by taking the empirical cumulative distribution functions (ECDF) based on the  $Z_i$  observed for the  $m$  vehicles with claims can be highly biased if the mileage rate affects the occurrence of claims. Figure 1 shows (a) the true  $G(z)$ , (b) the parametric estimate  $G(z; \hat{\mu}, \hat{\sigma})$  and (c) the naive ECDF based on  $z_1, \dots, z_m$  for cars with claims, in each of samples 1 ( $m = 38$ ) and 2 ( $m = 229$ ). As we would expect since  $\beta = 0.95$ , the

naive estimate (c) appears biased, especially in sample 1, and over-estimates the quantiles of  $G(z)$ .

## 5.2 Car Warranty Data

We consider the car warranty data mentioned at the start of this section. The analysis of "Type P" claims will be discussed; they have been considered by Chukova and Robinson (2006), who calculated the age-based and mileage-based rate functions in (14) and (16). Claims occurring before the date of sale, and a few claims with inconsistent data (e.g. mileage at time of claim larger than  $U_0 = 36$  thousand miles) were dropped, leaving us with 1540 claims across the 44,890 vehicles. The data base we consider was closed before all cars had reached the age limit of 3 years from date of sale, but almost all cars had values of  $\tau_i$  (current age at data base closure) greater than 2 years. The 1540 claims were experienced by  $m = 1270$  cars; 1068 cars had 1 claim, 151 had 2 claims, and 40, 8, and 3 cars had 3, 4, and 5 or more claims, respectively.

We fit models of the type given in Section 3. The ratio  $m/M$  is small here, and the simulated scenarios in Section 5.1 suggest that the likelihood may not be highly informative about certain parameters, especially when  $G(z)$  is estimated. External information about  $G(z)$  is thus valuable, but even then estimation of  $\beta$  within the model

$$E \{N_i(t)|v_i, z_i\} = v_i \Lambda_0 \left( tz_i^\beta \right) \quad (24)$$

is sensitive to the assumed distribution  $G(z)$ .

We will consider models (10) with random effects, with two models for  $\Lambda_0(t)$  in (24);

$$(i) \quad \Lambda_0(t) = (t/\alpha_1)^{\alpha_2} \quad \text{and} \quad (ii) \quad \Lambda_0(t) = \frac{e^{\alpha_2 t} - 1}{\alpha_1 \alpha_2} \quad (25)$$

respectively, for  $t \geq 0$ . Model (i) has been found satisfactory in many reliability and warranty claim settings similar to the one here. Model checks described below show fairly small but systematic departures from the observed data and so model (ii), which corresponds to the exponential rate function  $\lambda_0(t) = \alpha_1^{-1} \exp(\alpha_2 t)$ , is considered as an alternative.

We describe below three model checks that can be applied in the present setting. The first is based on the fact that if  $N_i(\tau_i) = n_i > 0$ , then under the models (10) the quantities

$$r_{ij} = \frac{\Lambda(t_{ij}|z_i)}{\Lambda(\tau_i|z_i)} \quad j = 1, \dots, n_i \quad (26)$$

are distributed as uniform  $(0, 1)$  order statistics in a sample of size  $n_i$  (e.g. Cook and Lawless, 2007, Section 3.7.3). Therefore we treat the  $\hat{r}_{ij}$  obtained by inserting parameter estimates in (26) as “uniform” residuals; a probability plot of the  $n. = \sum_{i=1}^m n_i$  ordered values of the  $\hat{r}_{ij}$  against the uniform  $(0, 1)$  expected order statistics  $\ell/(n. + 1)$ ,  $\ell = 1, \dots, n.$  should be approximately linear with slope one if the model  $\Lambda_0(t; \alpha)$  is suitable.

A second type of model check is to compare observed and expected numbers of warranty claims at different ages  $0 \leq t \leq T_0$  (here  $T_0 = 3$  years). Under a model (10) with a warranty mileage limit  $U_0$  (here  $U_0 = 36$  thousand miles), the expected number of claims for a vehicle up to age  $t$  is given by

$$\begin{aligned} \Lambda^*(t) &= E\{N_i(t \wedge (U_0/Z))\}, \\ &= \int_0^{U_0/t} \Lambda_0(z^\beta t) g(z) dz + \int_{U_0/t}^\infty \Lambda_0(z^{\beta-1} U_0) g(z) dz. \end{aligned} \quad (27)$$

Inserting parameter estimates  $\hat{\alpha}$  in (27) gives the estimated expected claims curve  $\hat{\Lambda}^*(t)$ .

This may be plotted and compared with the Nelson-Aalen estimate

$$\hat{\Lambda}_{NA}^*(t) = \sum_{(i,j):t_{ij} \leq t} \left( \frac{1}{Y_{\cdot}(t_{ij})} \right), \quad (28)$$

where  $Y_{\cdot}(t_{ij}) = \sum_{\ell=1}^M I(\tau_\ell^* \geq t_{ij})$ . Note that (28) estimates the expected number of claims  $\Lambda^*(t)$  nonparametrically.

Figure 2 shows a plot of  $\hat{\Lambda}^*(t)$  based on model (i) in (25), along with  $\hat{\Lambda}_{NA}^*(t)$ . There is a fairly small but systematic difference between the two estimates of  $\Lambda^*(t)$ . Pointwise confidence limits for  $\Lambda^*(t)$  based on the Poisson variance estimate (e.g. Cook and Lawless, 2007, Section 3.4.1) for  $\hat{\Lambda}_{NA}^*(t)$ ,

$$\widehat{\text{Var}}_p \left\{ \hat{\Lambda}_{NA}^*(t) \right\} = \sum_{(i,j):t_{ij} \leq t} \left( \frac{1}{Y_{\cdot}(t_{ij})^2} \right),$$

are also shown in Figure 2. This estimate is certainly too small, given evidence shown below of extra-Poisson variation, and the confidence limits are pointwise and not simultaneous but even so we note that the model-based curve  $\hat{\Lambda}^*(t)$  falls only slightly outside the limits at values  $0 \leq t \leq 3$ . Nevertheless, in view of the systematic nature of the departure of the model-based estimate from the empirical (Nelson-Aalen) estimate, we consider alternative models for  $\Lambda_0(t)$ , and in particular, model (ii) in (25).

Figure 3 shows the analogous plot to Figure 2 for models (ii), and we see better agreement with the nonparametric estimate. We have shown parametric estimates both with  $\mu$  and  $\sigma$  assumed known, and estimated; the estimates are more or less indistinguishable. Uniform probability plots of residuals (26) are close to linear for model (ii), as well. Parameter estimates for the model are shown in Table 3, where the reparameterization  $\alpha_{1c} = \alpha_1/12^\beta$ ,  $\alpha_{2c} = \alpha_2(12^\beta)$  is used in conjunction with the rescaled usage rate  $Z_{ci} = Z_i/12$  in order to reduce correlation between  $\hat{\alpha}_1$  and  $\hat{\beta}$ . In this case we did not have a reliable external estimate of  $G(z)$  and so estimated it, assuming it to be log-normal. The maximum likelihood estimates for  $\alpha_{1c}$ ,  $\alpha_{2c}$ ,  $\beta$ ,  $\phi$ ,  $\mu$  and  $\sigma$  were obtained by maximizing the log-likelihood corresponding to (11). In order to illustrate the sensitivity of  $\hat{\beta}$  to assumptions about  $G(z)$ , we also show in Table 3 the estimates of  $\alpha_{1c}$ ,  $\alpha_{2c}$ ,  $\beta$  and  $\phi$  that are obtained when we assume respectively that  $\mu = 2.5$ ,  $\sigma = 0.7$  and  $\mu = 2.7$ ,  $\sigma = 0.7$ . We observe that  $\hat{\beta}$  changes substantially when  $\mu$  is changed from 2.5 (corresponding to a median usage rate of 12,182 miles per year) to 2.7 (corresponding to a median usage rate of 14,880 miles per year). Thus if an external estimate of  $G(z)$  is used, one should be confident that it applies to the population of vehicles represented in the warranty claims data base. On the other hand, estimation of the expected number of claims, as shown in Figure 3, is insensitive to the value of  $\mu$ , as we noted above.

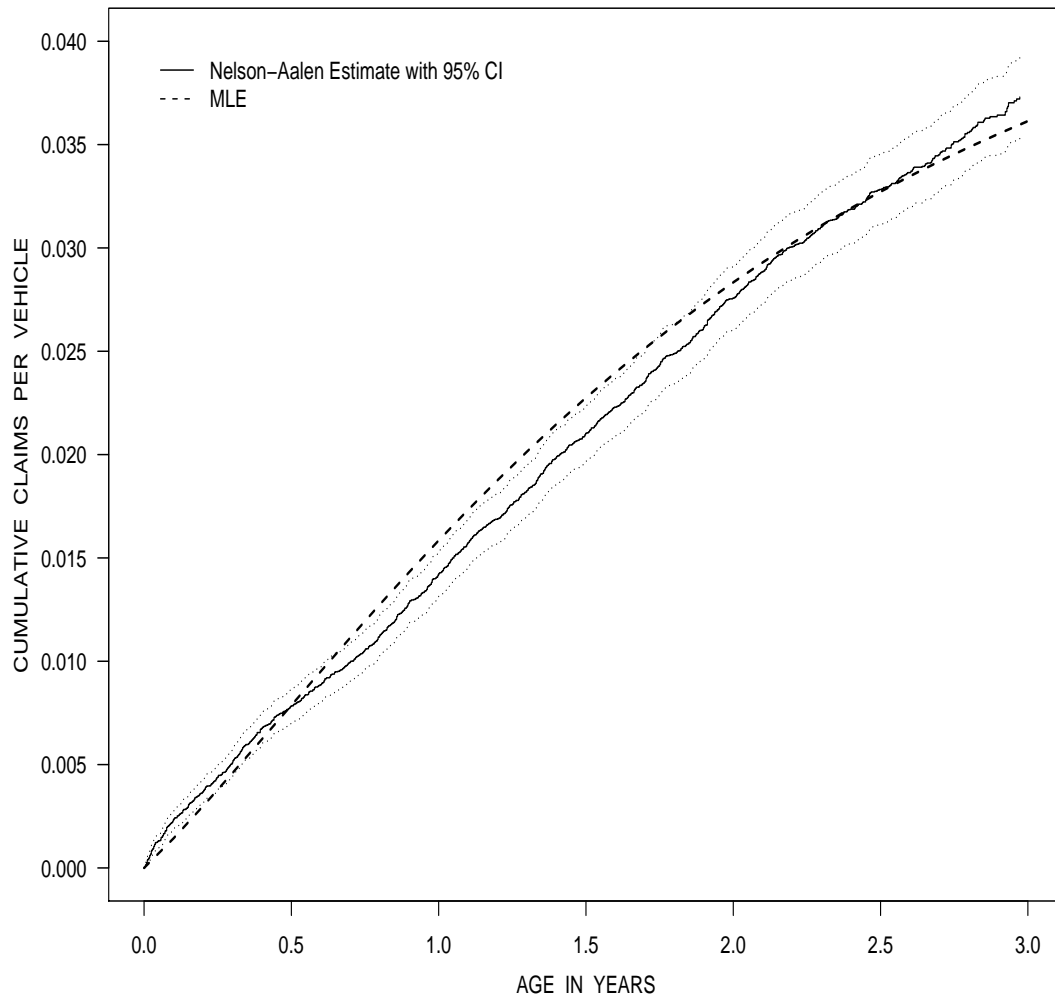


Figure 2. Parametric (model (25) (i)) and nonparametric estimates of  $\Lambda^*(t)$ , with pointwise Poisson confidence bands



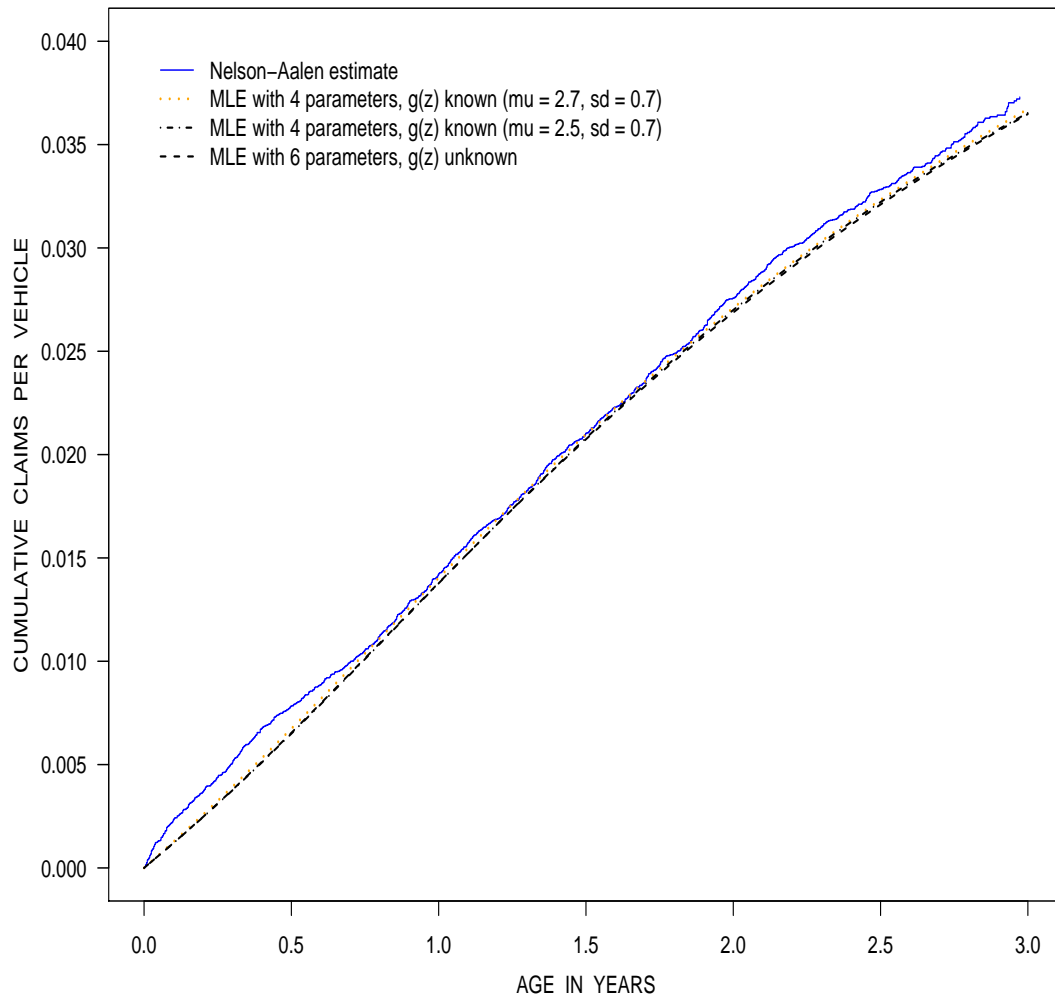


Figure 3. Parametric (model (25) (ii)) and nonparametric estimates of  $\Lambda^*(t)$

**Table 3.** Parameter Estimates for Model (10) with Mean Function (25) (ii)

Parameter	EST.	S.E.	EST.	S.E.	EST.	S.E.
$\alpha_{1c}$	94.7	5.3	94.2	5.2	87.8	5.0
$\alpha_{2c}$	0.320	0.037	0.303	0.036	0.288	0.037
$\beta$	0.576	0.096	0.616	0.038	0.307	0.040
$\phi$	10.1	0.8	10.3	0.8	10.4	0.8
$\mu$	2.52	0.065	2.5 <sup>1</sup>	-	2.7 <sup>1</sup>	-
$\sigma$	0.738	0.019	0.7 <sup>1</sup>	-	0.7 <sup>1</sup>	-

<sup>1</sup>  $\mu$  and  $\sigma$  are fixed at the values shown

Model (ii) in (25) gives a substantially larger maximum log-likelihood  $\ell(\hat{\theta})$  than model (i): -11885.26 versus -11917.27. It is possible to consider a three-parameter family  $\lambda_0(t; \alpha)$  that includes both model (i) and model (ii) (e.g. Lee, 1980) and within such a family model (ii) is thus better supported. We have not attempted to fit a three-parameter family in view of the limited information in these highly truncated claims data. Model (ii) is acceptable for estimation of the expected number of warranty claims, but we note one additional model check. In particular, we can for an arbitrary vehicle estimate the probability of  $r$  claims under model (10), without knowing the mileage rate  $Z_i$ . As shown in Appendix 1, this is

$$\hat{p}_i(r) = \frac{\Gamma(r + \hat{\omega})\hat{\omega}^{\hat{\omega}}}{\Gamma(\Lambda_0)r!} \int_0^\infty \frac{\Lambda_0 \left(\tau_i z_i^{\hat{\beta}}\right)^r}{\left\{\hat{\omega} + \Lambda_0 \left(\tau_i z_i^{\hat{\beta}}\right)\right\}^{r+\hat{\omega}}} \hat{g}(z) dz.$$

For each of  $r = 0, 1, 2, \dots$  we may then calculate expected frequencies  $e_r = \sum_{i=1}^M \hat{p}_i(r)$  and compare them with the observed frequencies  $f_r$  of cars with  $r$  claims. In fact both models considered here fit the data very well in this respect, which is not surprising because they have six parameters. For example, the model represented in Table 3 gives  $f_0 = 43620$ ,  $e_0 = 43624.9$ ;  $f_1 = 1068$ ,  $e_1 = 1073.1$ ;  $f_2 = 151$ ,  $e_2 = 154.3$ ;  $f_3 = 40$ ,  $e_3 = 29.4$ ;  $f_4 = 8$ ,  $e_4 = 6.4$ ;  $f_5 = 1$ ,  $e_5 = 1.5$ ;  $f_6 = 1$ ,  $e_6 = 0.4$ ;  $f_7 = 1$ ,  $e_7 = 0.1$ .

## 6. CHOICE OF WARRANTY LIMITS

In considering different possible limits on warranty coverage we need to investigate the distribution of

$$N(t, u) = \text{Number of claims with age } \leq t \text{ and usage } \leq u$$

for a given unit. No confusion should arise from using  $N(t, u)$  this way and  $N(t)$  to represent the number of claims occurring up to age  $t$ . As previously, we assume that  $U_i(t) = Z_i t$  and that conditional on the random variable  $v_i$  and usage rate  $Z_i$ , the events for unit  $i$  follow a Poisson process with mean function  $v_i \Lambda(t|Z_i)$ . To obtain the distribution of  $N(t, u)$  we observe that

$$N_i(t, u) = N_i(\min(t, u/Z_i))$$

and thus, denoting  $p_r(t|z) = \Pr\{N_i(t) = r|Z_i = z\}$ , we have

$$\begin{aligned} \Pr\{N(t, u) = r\} &= \int_0^\infty \Pr\{N_i(t, u) = r|Z_i = z\} g(z) dz \\ &= \int_0^{u/t} p_r(t|z) g(z) dz + \int_{u/t}^\infty p_r(u/z|z) g(z) dz \end{aligned}$$

For the case where  $\Lambda(t|z) = \Lambda_0(z^\beta t)$ , discussed earlier, and where  $v_i$  has a gamma distribution with mean 1 and variance  $\phi = \omega^{-1}$ , we have

$$p_r(t|z) = \frac{\omega^\omega \Gamma(r + \omega)}{\Gamma(\omega)} \frac{\Lambda_0(tz^\beta)^r}{\{\omega + \Lambda_0(tz^\beta)\}^{r+\omega}} \quad r = 0, 1, 2, \dots$$

and in the special Poisson process case where  $\phi = 0$ , we have

$$p_r(t|z) = e^{-\Lambda_0(tz^\beta)} \Lambda_0(tz^\beta)^r / r! \quad r = 0, 1, 2, \dots$$

For both models, we have

$$E\{N_i(t, u)\} = \int_0^{u/t} \Lambda_0(z^\beta t) g(z) dz + \int_{u/t}^\infty \Lambda_0(z^{\beta-1} u) g(z) dz. \quad (29)$$

As an example, we compute expected numbers of claims per vehicle using the parameter estimates for the model (25) (ii), given in Table 3. We consider three pairs of values for

$(t, u)$ : (3, 36), (4, 48) and (5,60). The first pair corresponds to the existing 3-year, 36000-mile warranty, and the other two represent longer warranty coverage. The estimates of (29) are respectively 0.0365, 0.0565 and 0.0824 so the indication is that extending the warranty even to 4 years, 48000 miles would increase expected claims substantially.

## 7. EXTENSIONS AND ADDITIONAL REMARKS

The fact that failures or other recurrent events may depend on both the age and usage history for a system leads to interesting modeling and inference problems. This paper has focussed on the rather difficult setting where the observed event data arise from warranty claims on a population of units. In this case the average number of claims per unit is typically very small and the usage histories are unobserved for units with no claims, making estimation much more difficult. We have adopted here a fairly simple model based on linear usage rates combined with a mixed Poisson process for events, conditional on the rates. This model appears satisfactory for the warranty data examined in the paper, but other models could of course be considered. The linear usage rate assumption is plausible for a substantial period following the purchase of many products, but investigation of the effects of mild variation around the linear function would be useful. As for the form of the event intensity  $\lambda(t|Z_i)$ , the accelerated time model (5) is convenient since it contains as special cases solely age-dependent and solely usage-dependent events. However, models of the multiplicative form  $\lambda(t|Z_i) = \lambda_0(t)r(Z_i; \beta)$  could easily be considered, say with  $r(Z_i; \beta) = \exp(\beta Z_i)$  or  $r(Z_i; \beta) = Z_i^\beta$ . Additive models, as in Murthy et al. (1995) can also be considered, but are more difficult to fit with warranty data.

We have seen in the case of warranty claims with low frequency that the amount of information about the relative effects of age and usage rate is limited, unless there is a fairly precise estimate of  $g(z)$ . An option not considered here is to do either some supplementary sampling to obtain  $z_i$  values for units with no claims, or to carry out random surveys of vehicles under warranty. Hu and Lawless (1996b) consider this in connection with failure

time models.

When there are multiple event types, the first priority is usually to fit models separately for each type. Poisson or mixed Poisson processes discussed here can be used for each event type. The same random variable  $Z_i$  applies across all event types, and this complicates estimation somewhat; this will be discussed in a separate article. Note also that the random variable  $Z_i$  induces association among event types. For example, if  $N_{i1}(t)$  and  $N_{i2}(t)$  count the number of type 1 and type 2 events, respectively, and if  $N_{i1}(t)$  and  $N_{i2}(t)$  are independent given  $Z_i$ , with  $E\{N_{i1}(t)|Z_i\} = \Lambda_1(t|Z_i)$  and  $E\{N_{i2}(t)|Z_i\} = \Lambda_2(t|Z_i)$ , then

$$\text{cov}(N_{i1}(t), N_{i2}(t)) = \text{cov}(\Lambda_1(t|Z_i), \Lambda_2(t|Z_i)) \quad (30)$$

where the covariance on the right side is with respect to the distribution of  $Z_i$ . For most models (30) will not have a simple form; an exception is for multiplicative models in which  $\Lambda_j(t|Z_i) = Z_i^{\beta_j} \Lambda_{0j}(t)$ . Association between event types, conditional on  $Z_i$ , can be modeled by adopting additional random effects along the lines of (10); see Cook and Lawless (2007, Chapter 6).

Comparison of different groups of units can be done by fitting separate models for each group. Regression models can also be considered. For example, if  $x_i$  is a  $p \times 1$  vector of covariates for unit  $i$  then we can extend the model (10) by taking

$$\lambda(t|\bar{U}_i, Z_i, v_i, x_i) = v_i Z_i^\beta e^{x_i' \gamma} \lambda_0\left(Z_i^\beta e^{x_i' \gamma} t\right), \quad (31)$$

where  $\gamma$  is a  $p \times 1$  vector of parameters. Fitting (31) with warranty data is likely to be somewhat challenging, but would be straightforward when  $Z_i$  is known for all units.

Finally, in many applications units tend to experience events rather frequently and there is also reasonably complete observation of the usage paths  $\bar{U}_i$  for most units. In this case models are much easier to fit, and a range of event intensity functions and processes for  $\bar{U}_i$  can be considered (Lawless and Crowder, 2007). Another problem of interest is when the occurrence of failures in a unit may affect their future usage. In this case  $\bar{U}_i(t)$  acts like an internal time-varying covariate with respect to event occurrence. Models that deal with

this and with settings where repairs involve substantial down-time for the system can be formulated, but are beyond the scope of this article.

## ACKNOWLEDGMENTS

The first author is supported by a grant from the Natural Sciences and Engineering Research Council of Canada. The authors thank Jeff Robinson for comments and for the data in Section 5.

## APPENDIX 1

Under the model (10) let  $g_v(v_i)$  denote the density function for  $v_i$ , which is assumed continuous. Since the  $v_i$  are unobservable, the likelihood function is obtained by integrating terms in (8) with respect to  $v_i$ , which is added to the intensity as in (10). That is,

$$\begin{aligned} L(\theta) &= \prod_{i=1}^m \int_0^\infty \left\{ \prod_{j=1}^{n_i} v_i z_i^\beta \lambda_0(t_{ij} z_i^\beta) \right\} e^{-v_i \Lambda_0(\tau_i z_i^\beta)} g_v(v_i) dv_i \\ &\quad \times \prod_{i=m+1}^M \int_0^\infty \int_0^\infty e^{-v_i \Lambda_0(\tau_i z_i^\beta)} g(z_i) g_v(v_i) dz_i dv_i. \end{aligned}$$

The second set of terms (for  $i = m + 1, \dots, M$ ) can be rewritten as

$$\prod_{i=m+1}^M \int_0^\infty \mathcal{L}_v(\Lambda_0(\tau_i z_i^\beta)) g(z_i) dz_i,$$

where  $\mathcal{L}_v(s) = E\{e^{-sv_i}\}$  is the Laplace transform of  $v_i$ .

If  $v_i$  has a gamma distribution with mean 1 and variance  $\phi = \omega^{-1}$ , then its Laplace transform is  $\mathcal{L}_v(s) = (1 + \omega^{-1}s)^{-\omega}$ . This produces the second term in (11). The first term in (11) is given by straightforward integration with

$$g_v(v) = \frac{\omega^\omega}{\Gamma(\omega)} v^{\omega-1} e^{-v\omega} \quad v > 0.$$

We can also derive the unconditional probability of  $r$  claims for a given vehicle, and the expected number of claims. Given  $\tau_i^* = \min(T_i, T_0)$ , but not  $Z_i$ , the probability of  $r$  claims is obtained by integrating the Poisson probability function with mean  $v_i Z_i^\beta \Lambda_0(\tau_i Z_i^\beta)$  with respect to  $v_i$  and  $Z_i$ . With  $v_i$  distributed according to a gamma distribution with mean 1 and

variance  $\phi = \omega^{-1}$ , and noting that  $\tau_i = \min(\tau_i^*, U_0/z_i)$ , we obtain  $p_i(r) = \Pr\{N_i(\tau_i^*) = r\}$  for  $r = 0, 1, 2, \dots$  as

$$p_i(r) = \frac{\Gamma(r + \omega)\omega^\omega}{\Gamma(\omega)r!} \int_0^\infty \frac{\Lambda_0 \left( \tau_i z_i^\beta \right)^r}{\left\{ \omega + \Lambda_0 \left( \tau_i z_i^\beta \right) \right\}^{r+\omega}} g(z_i) dz_i$$

## APPENDIX 2

**Proof of Theorem 1.** The estimate  $\hat{\lambda}_a(t)dt$  given by (14) converges in probability to

$$\begin{aligned} \lambda_a^*(t)dt &= E \{Y_i^a(t)n_i(t)\} / P_i^a(t) \\ &= \frac{E \{I(T_i \geq t)I(Z_i \leq U_0/t)n_i(t)\}}{\Pr(T_i \geq t) \Pr(Z_i \leq U_0/t)} \\ &= \frac{\int_0^{U_0/t} E \{n_i(t)|Z_i\} g(Z_i) dZ_i}{\Pr(Z_i \leq U_0/t)}, \end{aligned}$$

where we think in continuous time of  $n_i(t)$  as the number of events in the short interval  $(t, t + dt)$ , consider  $(n_i(t), T_i, Z_i)$  as i.i.d and use the fact that  $T_i$  and  $Z_i$  are independent. Since  $E\{n_i(t)|Z_i\} = \lambda(t|Z_i)dt$ , this gives

$$\lambda_a^*(t) = E \{ \lambda(t|Z) | Z \leq U_0/t \}. \quad (A1)$$

In addition,  $\lambda_a(t)dt = E\{n_i(t)\}$  so that

$$\lambda_a(t) = E \{ \lambda(t|Z) \} = \int_0^\infty \lambda(t|z)g(z)dz. \quad (A2)$$

This proves (18).

The conditional rate function in terms of usage satisfies

$$\lambda_u(u|Z)du = E \{n_i^u(u)|Z\} = Z^{-1}\lambda(uZ^{-1}|Z) du,$$

where  $n_i^u(u)$  is the number of events in the small interval  $(u, u + du)$  and we note that  $n_i^u(u) = n_i(u/z)$  with  $dt = du/z$ . The unconditional rate function is thus

$$\lambda_u(u) = E \{ Z^{-1}\lambda(uZ^{-1}|Z) \}. \quad (A3)$$



In addition, the estimate (16) converges in probability to

$$\begin{aligned}
\lambda_a^*(u)du &= E \{Y_i^u(u)n_i^u(u)\} / P_i^u(u) \\
&= \frac{E \{I(Z_i \geq u/\tau_i^*)n_i^u(u)\}}{\Pr(Z_i \geq u/\tau_i^*)} \\
&= \frac{E \int_{u/\tau_i^*}^{\infty} Z^{-1}\lambda(uZ^{-1}|Z)g(z)dzdu}{\Pr(Z_i \geq u/\tau_i^*)}, \tag{A4}
\end{aligned}$$

where the top and bottom of (A4) involve expectations with respect to the  $\tau_i$ . Together, (A3) and (A4) give (19).

If all units are observed for a time  $T_0$ , then  $\tau_i^* = T_0$  for  $i = 1, \dots, M$  and (A4) becomes

$$E \{Z^{-1}\lambda(uZ^{-1}|Z) | Z \geq u/T_0\}$$

and (19) becomes

$$B_u(u) = \frac{E \{Z^{-1}\lambda(uZ^{-1}|Z) | Z \geq u/T_0\}}{E \{Z^{-1}\lambda(uZ^{-1}|Z)\}}. \tag{A5}$$

**Proof of Theorem 2.** This follows immediately from replacing  $\lambda(t|Z)$  in (18) and (19) with  $Z^\beta[\alpha_1\alpha_2(tZ^\beta)^{\alpha_2-1}]$ .

### APPENDIX 3

Let the distinct values of  $z_i$  observed across units  $i = 1, \dots, m$  that have  $n_i > 0$  be denoted by  $z_1^*, \dots, z_R^*$ , and let  $m_r = \sum_{i=1}^m I(z_i = z_r^*)$  and  $g_r = \Pr(Z_i = z_r^*)$ , where we assume that the  $Z_i$  have support only on the set  $(z_1^*, \dots, z_R^*)$ . A nonparametric estimate of  $g = (g_1, \dots, g_R)$  can be obtained by maximizing the log likelihoods  $\ell(\theta) = \ell(\alpha, \beta, g)$  obtained from (8) or  $\ell(\theta) = \ell(\alpha, \beta, \phi, g)$  obtained from (11), replacing the integrals in (8) and (11) with sums over the  $z_r^*(r = 1, \dots, R)$ . A convenient approach is to alternatively update estimates of  $(\alpha, \beta)$  or  $\alpha, \beta, \phi$  and  $g$ , and to use an EM algorithm that treats the  $z_i(i = m + 1, \dots, M)$  as missing data.

Bearing in mind that  $\sum_{r=1}^R g_r = 1$ , we obtain an algorithm for updating the estimates  $g$  by treating the other parameters as fixed, and considering the estimating equations

$$\frac{m_r}{g_r} + \sum_{i=m+1}^M \frac{1}{g_r} \Pr(Z_i = z_r^* | n_i = 0) - \lambda = 0, \quad r = 1, \dots, R \quad (\text{A6})$$

where  $\lambda$  is a Lagrange multiplier. Since

$$\Pr(Z_i = z_r^* | n_i = 0) = \frac{g_r^{(0)} \Pr(n_i = 0 | z_r^*)}{\sum_{\ell=1}^R g_\ell^{(0)} \Pr(n_i = 0 | z_\ell^*)},$$

where  $g^{(0)} = (g_1^{(0)}, \dots, g_R^{(0)})$  is the current estimate of  $g$ , and using the fact that  $\lambda = M$  (multiply the terms of (A6) by  $g_r$  and sum over  $r = 1, \dots, R$ ), we find from (A6) that

$$g_r = M^{-1} \left[ m_r + \sum_{i=m+1}^M \left\{ \frac{g_r^{(0)} \Pr(n_i = 0 | z_r^*)}{\sum_{\ell=1}^R g_\ell^{(0)} \Pr(n_i = 0 | z_\ell^*)} \right\} \right] \quad (\text{A7})$$

where for the Poisson model giving (8) we have

$$\Pr(n_i = 0 | z_r^*) = e^{-\Lambda(\tau_i | z_r^*)} \quad (\text{A8})$$

and for the negative binomial model giving (11) we have

$$\Pr(n_i = 0 | z_r^*) = (1 + \omega^{-1} \Lambda(\tau_i | z_r^*))^{-\omega}. \quad (\text{A9})$$

In (A8) and (A9) we compute  $\Pr(n_i = 0 | z_r^*)$  using the most recent estimates of the other parameters.

A convenient algorithm for maximizing the semiparametric likelihoods  $\ell(\alpha, \beta, g)$  from (8) or  $\ell(\alpha, \beta, \phi, g)$  from (11) is as follows:

- (i) Take initial estimates  $g_r^{(0)} = m_r/m$  ( $r = 1, \dots, R$ ).
- (ii) Holding  $g = g^{(0)}$  fixed, maximize  $\ell(\alpha, \beta, g)$  or  $\ell(\alpha, \beta, \phi, g)$  using general purpose optimization software.
- (iii) Update  $g$  by computing (A7), using the current estimates of  $\alpha, \beta, \phi$ , for  $r = 1, \dots, R$  and then replace  $g^{(0)}$  with  $g$ .

Repeat steps (ii) to (iii) until convergence.

## REFERENCES

- Ahn, CW, Chae, KC and Clark, GM (1998). Estimating parameters of the power law process with two measures of failure time. *J. Qual. Tech.* **30**, 127-132.
- Chukova, S and Robinson, J (2006). Estimating mean cumulative functions from truncated automobile warranty data. *In Proceedings of 2004 Mathematical Methods in Reliability Conference*, eds.
- Cook, RJ and Lawless, JF (2007). *The Statistical Analysis of Recurrent Events*. Springer, New York.
- Crowder, MJ and Lawless, JF (2007). Failures and other events in systems with age and usage scales. Manuscript.
- Hu, XJ and Lawless, JF (1996a). Estimation of rate and mean functions from truncated recurrent event data. *J. Amer. Statist. Assoc.* **91**, 300-310.
- Hu, XJ and Lawless, JF (1996b). Estimation from truncated lifetime data with supplementary information on covariates and censoring times. *Biometrika* **83**, 747-761.
- Jung, M and Bai, DS (2007). Analysis of field data under two-dimensional warranty. *Rel. Eng. Syst. Safety* **92**, 135-143.
- Kalbfleisch, JD, Lawless, JF and Robinson, JA (1991). Methods for the analysis and prediction of warranty claims. *Technometrics* **33**, 273-285.
- Krivtsov, V. (2006). Personal communication.
- Lawless, JF (1987). Regression methods for Poisson process data. *J. Amer. Statist. Assoc.* **82**, 808-815.
- Lawless, JF (1998). Statistical analysis of product warranty data. *Int. Statist. Rev.* **66**, 41-60.
- Lawless, JF and Kalbfleisch, JD (1992). Some issues in the collection and analysis of field reliability data. *In Survival Analysis: State of the Art*, 141-151. Eds. J. Klein and P. Goel, Kluwer, Amsterdam.

- Lawless, J, Hu, J and Cao, J (1995). Methods for the estimation of failure distributions and rates from automobile warranty data. *Lifetime Data Anal.* **1**, 227-240.
- Lee, L. (1980). Testing adequacy of the Weibull and loglinear rate models for a Poisson process. *Techometrics*, **22**, 195-199.
- Murthy, DNP, Iskandar, BP and Wilson, RJ (1995). Two-dimensional failure free warranty policies: two-dimensional point process model. *Oper. Res.* **43**, 356-366.
- Robinson, J and McDonald, G (1991). Issues related to field reliability and warranty data. In *Data Quality Control: Theory and Pragmatics*. Eds. G. Liepins and V. Uppuluri. Dekker, New York.
- Singpurwalla, ND and Wilson, SP (1998). Failure models indexed by two time scales. *Adv. Appl. Prob.* **30**, 1058-1072.
- Suzuki, K (1985). Estimation of lifetime parameters from incomplete field data. *Technometrics* **27**, 263-272.
- Suzuki, K (1993). Estimation of lifetime distribution using the relationship of calendar time and age. *Rep. Statist. Appl. Res.* **40**, 10-22.
- Yang, SC and Nachlas, JA (2001). Bivariate reliability and availability modeling. *IEEE Trans. Rel.* **50**, 26-35.