

**Nonparametric Estimation of the Number of Components  
of A Superposition of Renewal Processes**

ANUP DEWANJI

*Applied Statistics Unit, Indian Statistical Institute, 203 B. T. Road,  
Kolkata 700108, India*

SUBRATA KUNDU AND TAPAN K. NAYAK

*Department of Statistics, George Washington University, Washington, D.C., USA*

**Abstract**

Suppose all events occurring in an unknown number ( $\nu$ ) of iid renewal processes, with a common renewal distribution  $F$ , are observed for a fixed time  $\tau$ , where both  $\nu$  and  $F$  are unknown. The individual processes are not known apriori, but for each event, the process that generated it is identified. For example, in software reliability application, the errors (or bugs) in a piece of software are not known apriori, but whenever the software fails, the error causing the failure is identified. We present a nonparametric method for estimating  $\nu$  and investigate its properties. Our results show that the proposed estimator performs well in terms of bias and asymptotic normality, while the MLE of  $\nu$  derived assuming that the common renewal distribution is exponential may be seriously biased if that assumption does not hold.

**Keywords:** Asymptotic normality; bias; Kaplan-Meier estimator; profile likelihood; software reliability.

# 1. Introduction

Suppose signals are received, in continuous time, from an unknown number ( $\nu$ ) of sources, where for each signal, its source is also observed. Thus, any specific source is detected when its first signal is received. Obviously, any source that does not send a signal during the observation period remains undetected. While each source induces a point process, the observed event (signal receiving) times come from the superposition of all the processes. The situation considered here may arise in continuous time capture-recapture experiments, where each animal or species acts as one source. In software testing, each error (or bug) is a source for causing software failures and the overall failure process is a superposition of an unknown number of point processes.

We shall consider the case where the overall process is observed for a fixed time  $\tau$ . Formally, we observe the values of the following random variables:  $R$  = number of detected processes (or sources),  $M_i$  = the number of events in the  $i$ th detected process,  $i = 1, \dots, R$ , and  $T_{ij}$  = the  $j$ th inter-event (renewal) time in the  $i$ th detected process,  $i = 1, \dots, R, j = 1, \dots, M_i$ . Then,  $S_i = \tau - \sum_{j=1}^{M_i} T_{ij}$  is the last censored event time for the  $i$ th detected process. We shall let  $M = \sum_{i=1}^R M_i$  denote the total number of events observed. Following the standard convention we shall use upper case letters to denote random variables and lower case letters to denote their observed values.

It is most convenient to assume that all component processes are iid Poisson processes with a common but unknown rate  $\lambda$  (e.g., Becker, 1984; Nayak, 1988, 1991; Chao and Lee, 1993). Under the Poisson assumption, the maximum likelihood estimators (MLE) of  $\nu$  and  $\lambda$  as well as their distributions have been obtained in Nayak (1988, 1991), which are reviewed briefly in Section 2. More generally, Dewanji et al. (1995) assumed that the

component processes are iid renewal processes with a common renewal density  $f_\theta(\cdot)$  and cdf  $F_\theta(\cdot)$ , where  $\theta$  is an unknown parameter vector. They discussed maximum likelihood estimation of  $\nu$  and  $\theta$  and derived asymptotic normality of a class of estimators, which includes the MLE. We shall review some of their results in the next section. Several authors have discussed the estimation of a population size based on continuous time capture-recapture data under various settings; we refer to Chao and Lee (1993), Lloyd (1994), Wilson and Anderson (1995), Xi et al. (2007), and Yoshida et al. (1996) for further reading and additional references.

In this paper, we focus on estimating  $\nu$ , assuming that all components are iid renewal processes. In contrast to Dewanji et al. (1995), which assumes a parametric family for the renewal density, we present a nonparametric estimator. Naturally, any estimator of  $\nu$  derived for a parametric family of renewal distributions is likely to be sensitive to the assumptions about the renewal distribution. Because of mathematical and computational simplicity one may be tempted to use the MLE ( $\hat{\nu}_{exp}$ ) of  $\nu$  derived assuming that the renewal distribution is exponential. A natural question is: Can the simple estimate based on exponential distribution be used regardless of the true renewal distribution? Thus, we also explore the bias of  $\hat{\nu}_{exp}$  when the true distribution is not exponential.

In Section 2, we describe briefly some results of Nayak (1988) and Dewanji et al. (1995) for Poisson and parametric renewal processes. In Section 3, we introduce our nonparametric method along with two computational methods for obtaining the nonparametric MLE ( $\hat{\nu}_{np}$ ) of  $\nu$ . In section 4, we discuss some asymptotic properties of  $\hat{\nu}_{np}$ , as  $\nu \rightarrow \infty$ . In Section 5, we investigate, by means of simulation, the performance of both of  $\hat{\nu}_{np}$  and  $\hat{\nu}_{exp}$  under different renewal distributions. We find that the nonparametric estimator performs well in terms of bias and asymptotic normality, while  $\hat{\nu}_{exp}$  can be severely biased when

the true renewal distribution is not exponential.

## 2. Preliminaries

In a parametric setup, the likelihood function for the observed data described in the previous section can be written as (see Dewanji et al., 1995)

$$L(\nu, \theta) = \frac{\nu!}{(\nu - r)!} [\bar{F}(\tau; \theta)]^{\nu - r} \prod_{i=1}^r \left[ \left\{ \prod_{j=1}^{m_i} f(t_{ij}; \theta) \right\} \bar{F}(s_i; \theta) \right], \quad (1)$$

where  $f(\cdot; \theta)$  is the renewal density function with  $\bar{F}(\cdot; \theta)$  being the corresponding survival function. For any given  $\theta$  (with  $F(\tau; \theta) < 1$ ), the likelihood in (1) is maximized (with respect to  $\nu$ ) by

$$\hat{\nu}(\theta) = \lfloor \frac{r}{1 - \bar{F}(\tau; \theta)} \rfloor \quad (2)$$

where  $\lfloor x \rfloor$  is the largest integer not exceeding  $x$ . Thus, if  $\theta$  is known,  $\hat{\nu}(\theta)$  in (2) is the MLE of  $\nu$ . For most distribution families, the maximization of (1) with respect to  $\theta$ , for given  $\nu$ , can be done only numerically. For the joint estimation of  $\nu$  and  $\theta$ , Dewanji et al. (1995) suggested the following iterative procedure. Starting with an initial estimate  $\nu^{(0)}$  of  $\nu$ , calculate an estimate  $\theta^{(0)}$  of  $\theta$  by maximizing  $L(\nu^{(0)}, \theta)$ . Then, use  $\theta^{(0)}$  in (2) to obtain  $\nu^{(1)} = \hat{\nu}(\theta^{(0)})$ , and then calculate  $\theta^{(1)}$  by maximizing  $L(\nu^{(1)}, \theta)$  with respect to  $\theta$ . Repeat this process until convergence is achieved. We may also note that the MLE of  $\nu$  and  $\theta$  may not exist. For example, if the renewal distribution is exponential, the MLE exists if and only if  $m > r$ , where  $m = \sum_{i=1}^r m_i$ .

As a simpler alternative, Dewanji et al. (1995) suggested a conditional maximum likelihood (CML) method, where  $\theta$  is estimated by maximizing (with respect to  $\theta$ ) the

conditional likelihood, given  $R = r$ , which is given by

$$r! \prod_{i=1}^r \left[ \frac{f(t_{i1}; \theta)}{F(\tau; \theta)} \right] \prod_{i=1}^r \left[ \left\{ \prod_{j=2}^{m_i} f(t_{ij}; \theta) \right\} \bar{F}(s_i; \theta) \right]. \quad (3)$$

The maximizer of (3), denoted by  $\hat{\theta}_c$ , is a conditional maximum likelihood estimator (CMLE) of  $\theta$ . The CMLE of  $\nu$  is then obtained by using  $\theta = \hat{\theta}_c$  in (2).

Dewanji et al. (1995) also derived certain asymptotic properties (as  $\nu \rightarrow \infty$ ) of the MLE and CMLE. In particular, (i) if (3) is maximized uniquely, then under certain mild conditions,  $\hat{\theta}_c$  is a consistent estimator of  $\theta$ , (ii) the CMLE and MLE are asymptotically equivalent and (iii) for a broad class of estimators  $(\hat{\nu}, \hat{\theta})$ , which includes the MLE and the CMLE, as  $\nu \rightarrow \infty$ ,

$$[\nu^{1/2}(\hat{\theta} - \theta), \nu^{-1/2}(\hat{\nu} - \nu)] \xrightarrow{L} \mathcal{N}(0, \Sigma), \quad (4)$$

where

$$\Sigma = \begin{pmatrix} I(\theta) & -\delta \\ -\delta' & F(\tau; \theta)/\bar{F}(\tau; \theta) \end{pmatrix}^{-1},$$

$\delta = \frac{\partial}{\partial \theta} \log \bar{F}(\tau; \theta)$  and  $I(\theta)$  is the information matrix based on observing a single component process for time  $\tau$ . One can estimate  $I(\theta)$  consistently using the observed information matrix, given by

$$-\frac{1}{\hat{\nu}} \frac{\partial^2 \log L(\hat{\nu}, \hat{\theta})}{\partial \theta \partial \theta'}.$$

In practice,  $\Sigma$  may be consistently estimated by replacing  $\theta$  by  $\hat{\theta}$  and  $I(\theta)$  by its consistent estimate as given above.

For exponential renewal distribution with density  $f(t) = \lambda e^{-\lambda t}$ ,  $\lambda > 0$ ,  $t > 0$ , the likelihood in (1) reduces to

$$L(\nu, \lambda) = \frac{\nu!}{(\nu - r)!} \lambda^m \exp[-\nu \lambda \tau]. \quad (5)$$

It follows that the maximum of (5) does not exist if  $m = r$ . Otherwise, the MLE of  $\nu$ , to be denoted by  $\hat{\nu}_{exp}$ , is determined (see Nayak, 1988) by the function

$$g(\nu) = \left(1 - \frac{1}{\nu}\right)^m + \frac{r}{\nu}, \quad \nu > r.$$

If  $g(r+1) < 1$ , then  $\hat{\nu}_{exp} = r$ . If  $g(r+1) > 1$ , then  $\hat{\nu}_{exp} = \max\{\nu : g(\nu) > 1\}$ . The MLE  $\hat{\lambda}$  of  $\lambda$  is then obtained as  $\hat{\lambda} = m/(\tau\hat{\nu}_{exp})$ . From the general result discussed above, the asymptotic distribution of  $[\nu^{1/2}(\hat{\lambda} - \lambda), \nu^{-1/2}(\hat{\nu}_{exp} - \nu)]$ , as  $\nu \rightarrow \infty$ , is normal, whose covariance matrix  $\Sigma$  can be consistently estimated by

$$\hat{\Sigma} = \begin{pmatrix} \frac{m}{\hat{\nu}_{exp}\hat{\lambda}^2} & \tau \\ \tau & e^{\hat{\lambda}\tau} - 1 \end{pmatrix}^{-1}.$$

In particular, the asymptotic variance of  $\nu^{-1/2}(\hat{\nu}_{exp} - \nu)$  can be estimated by

$$m[(e^{\hat{\lambda}\tau} - 1)m - \hat{\nu}_{exp}\hat{\lambda}^2\tau^2]^{-1}. \quad (6)$$

### 3. Nonparametric Estimation

As discussed in Section 1, we shall take  $\nu$  as the parameter of primary interest and to develop a nonparametric method for estimating  $\nu$ , we shall let the renewal distribution be arbitrary with pdf  $f(\cdot)$  and cdf  $F(\cdot)$ . Here, the renewal distribution acts as a nuisance parameter, and the likelihood function in (1) takes the form

$$L(\nu, f) = \frac{\nu!}{(\nu - r)!} [\bar{F}(\tau)]^{\nu-r} \prod_{i=1}^r \left[ \left\{ \prod_{j=1}^{m_i} f(t_{ij}) \right\} \bar{F}(s_i) \right]. \quad (7)$$

To find the nonparametric MLE of  $\nu$  we need to maximize (7) with respect to  $\nu$  and  $f$ .

In (7), the factor  $[\nu!/(\nu - r)!]$  does not involve the renewal distribution. So, if  $\nu$  is known, the product of all terms, excluding  $[\nu!/(\nu - r)!]$ , should be considered for

estimating  $f$ . Note that the product of those terms has the form of the Kaplan-Meier likelihood (Kaplan and Meier, 1958) with censored data. In the spirit of the Kaplan-Meier estimator, it is enough to consider all renewal distributions with mass concentrated at the observed renewal times  $\{t_{ij}, i = 1, \dots, r, j = 1, \dots, m_i, \}$  and another point larger than  $\tau$ . Thus, letting  $x_1 < \dots < x_n$  denote the distinct ordered values of the  $t_{ij}$ 's, we shall consider all the probability distributions with sample space  $\{x_1, x_2, \dots, x_{n+1}\}$ , where  $x_{n+1}$  is some value greater than  $\tau$ . Let  $\pi_l = P(X = x_l), l = 1, \dots, n+1$ , with  $\sum_{l=1}^{n+1} \pi_l = 1$ , and  $f_l$  denote the frequency of  $x_l$ , for  $l = 1, \dots, n$ . Now, the likelihood function in (7) can be written in terms  $\pi_1, \dots, \pi_{n+1}$  as

$$L(\nu, \tilde{\pi}) = \frac{\nu!}{(\nu - r)!} (\pi_{n+1})^{\nu-r} \prod_{l=1}^n (\pi_l)^{f_l} \prod_{i=1}^r \left( \sum_{x_h > s_i} \pi_h \right), \quad (8)$$

Then, (8) can be maximized to obtain the MLE of  $\pi_1, \dots, \pi_{n+1}$  and  $\nu$ .

However, it is more convenient to work with the discrete hazard components

$$\lambda_l = \frac{\pi_l}{\sum_{j=1}^n \pi_j}, \quad l = 1, \dots, n.$$

The transformation from  $(\pi_1, \dots, \pi_{n+1})$  to  $(\lambda_1, \dots, \lambda_n)$  is one-to-one, with  $\pi_1 = \lambda_1$ ,

$$\pi_l = \lambda_l \prod_{j=1}^{l-1} (1 - \lambda_j), \quad l = 2, \dots, n \quad \text{and} \quad \pi_{n+1} = \prod_{j=1}^n (1 - \lambda_j).$$

Note that, in our context, the survival function  $\bar{F}(t)$  at time  $t$  can be written, in terms of

$\tilde{\lambda} = (\lambda_1, \dots, \lambda_n)$ , as

$$\bar{F}(t; \tilde{\lambda}) = \prod_{l: x_l \leq t} (1 - \lambda_l), \quad (9)$$

which implies that  $\bar{F}(\tau; \tilde{\lambda}) = \prod_{l=1}^n (1 - \lambda_l)$ . In terms of  $\tilde{\lambda}$ , the likelihood function is

$$L(\nu, \tilde{\lambda}) = \frac{\nu!}{(\nu - r)!} \prod_{l=1}^n (\lambda_l)^{f_l} \prod_{l=1}^n [1 - \lambda_l]^{c_l(\nu)}, \quad (10)$$

where  $c_l(\nu) = (\nu - r) + \sum_{u=l+1}^n f_u + k_l$  and  $k_l = \#\{i : s_i \geq x_l\}$ . Note that, for any given  $\nu$ ,  $c_l(\nu)$  is the number of renewal times known to be greater than  $x_l$ . It can be seen that, for any given  $\nu$ , (10) is maximized by

$$\hat{\lambda}_l(\nu) = \frac{f_l}{f_l + c_l(\nu)}, \quad l = 1, \dots, n. \quad (11)$$

Substituting (11) in (10), we get the profile likelihood function  $L_1(\nu) = L(\nu, \hat{\lambda}(\nu))$  given by

$$L_1(\nu) = \frac{\nu!}{(\nu - r)!} \prod_{l=1}^n \{\hat{\lambda}_l(\nu)\}^{f_l} \prod_{l=1}^n \{(1 - \hat{\lambda}_l(\nu))\}^{c_l(\nu)}. \quad (12)$$

Now, the MLE of  $\nu$ , to be denoted by  $\hat{\nu}_{np}$ , can be obtained by maximizing  $L_1(\nu)$  with respect to  $\nu$ .

The maximizer of (12), however, can be found only numerically and it may not exist. If the renewal distribution is exponential, the profile likelihood can be an increasing function of  $\nu$ , and that happens if and only if  $m = r$ , see Nayak (1988). Our profile likelihood in (12) is a more complex function of  $\nu$  and we could not find exact conditions for the MLE to be infinity, but the following observations indicate that a necessary condition for  $\hat{\nu}_{np} = \infty$  is  $m = r$ . Note that  $[\nu! / (\nu - r)!] \sim \nu^r$ , where  $a(\nu) \sim b(\nu)$  means  $a(\nu)/b(\nu) \rightarrow 1$  as  $\nu \rightarrow \infty$ . Also, in view of (11), the first product in (12) is  $\sim \nu^{-m}$  and the second product is  $\sim e^a$ , where  $a$  is a constant. Thus,  $L_1(\nu) \sim e^a \nu^{(r-m)}$ , which is a decreasing function of  $\nu$  unless  $m = r$ . This means that if  $m > r$ ,  $L_1(\nu)$  cannot grow indefinitely as  $\nu$  increases to infinity. So, a necessary condition for the MLE of  $\nu$  to be  $\infty$  is  $m = r$ . Note that when  $m = r$ ,  $\hat{\nu}_{np} = \infty$  makes some intuitive sense because in this case every observed event discovers a new process, which suggests that the number of component processes is very large. Also, note that  $P(M = R)$  decreases to 0 as  $\tau$  and/or  $\nu$  increases to infinity and hence  $\hat{\nu}_{np} = \infty$  is a rare event in most practical situations.

The nonparametric MLE of  $\nu$  can also be computed using an iterative procedure as described next. For any given  $\tilde{\lambda}$  with  $0 < \lambda_l < 1, l = 1, \dots, n$ , we get

$$\frac{L(\nu + 1, \tilde{\lambda})}{L(\nu, \tilde{\lambda})} = \frac{\nu + 1}{\nu - r + 1} \prod_{l=1}^n (1 - \lambda_l), \quad (13)$$

noting that  $c_l(\nu + 1) = c_l(\nu) + 1$ . Let us denote the product  $\prod_{l=1}^n (1 - \lambda_l)$  by  $\gamma$ . Then, as  $\nu$  increases from  $r$  to  $\infty$ , the ratio in (13) decreases from  $(1 + r)\gamma$  to  $\gamma$ , which is less than 1 (as  $0 < \lambda_l < 1, l = 1, \dots, n$ ). So, the ratio can cross the line  $y = 1$  at most once and from above. If  $(1 + r)\gamma < 1$ , then  $L(\nu, \tilde{\lambda})$  is maximum at  $\nu = r$ . If  $(1 + r)\gamma > 1$ , then as  $\nu$  increases (from  $r$ ),  $L(\nu, \tilde{\lambda})$  first increases and then decreases and the maximum occurs at

$$\nu = \lfloor r[1 - \prod_{l=1}^n (1 - \lambda_l)]^{-1} \rfloor. \quad (14)$$

This is consistent with (2), in view of (9).

The fact that both (11) and (14) are easy to compute suggests that an iterative method for calculating  $\hat{\nu}_{np}$  may be computationally more efficient than maximizing the profile likelihood in (12). A natural approach would be to start with an initial estimate  $\nu^{(0)} > r$  of  $\nu$ , calculate  $\hat{\lambda}^{(1)}$  using (11) and then, obtain a revised estimate  $\nu^{(1)}$  of  $\nu$  using  $\hat{\lambda}^{(1)}$  in (14). Subsequent steps would be similar. At the  $i$ th ( $i \geq 2$ ) stage, calculate  $\hat{\lambda}^{(i)}$  by using  $\nu^{(i-1)}$  in (11) and then calculate  $\nu^{(i)}$  by putting  $\hat{\lambda}^{(i)}$  in (14). Note that since (11) and (14) maximize the likelihood function, for given  $\nu$  and  $\tilde{\lambda}$ , respectively, it follows that

$$L(\nu^{(0)}, \hat{\lambda}^{(1)}) \leq L(\nu^{(1)}, \hat{\lambda}^{(1)}) \leq L(\nu^{(1)}, \hat{\lambda}^{(2)}) \leq L(\nu^{(2)}, \hat{\lambda}^{(2)}) \leq L(\nu^{(2)}, \hat{\lambda}^{(3)}) \leq \dots$$

Naturally, one would continue the iterative process until it converges.

We found that the preceding approach often fails to yield the actual  $\hat{\nu}$  due to its discrete nature. Specifically, in an iterative step, the value of  $r[1 - \prod_{l=1}^n (1 - \lambda_l)]^{-1}$  may

change by a small magnitude so that its integer part remains the same, in which case the value of  $\nu$  in (14) does not change. Thus, the above procedure often gets stuck at some value of  $\nu$  and corresponding  $\hat{\lambda}$ , which are not the true MLE. In our simulation studies, we found that a slight modification of the preceding procedure works well. Specifically, in each step of our iteration, we use the actual value of  $r[1 - \prod_{l=1}^n (1 - \lambda_l)]^{-1}$  (rather than its integer part) to update the value of  $\nu$ . Thus, in intermediate steps  $\nu$  and hence  $c_l(\nu)$  may not be integer valued. When we stop the process, based on convergence, we take the integer part of the latest value  $\nu^{(k)}$  for the MLE of  $\nu$ . Also, in most cases, the starting value of  $\nu^{(0)} = r$  yields  $\nu^{(1)} = r$ . This is because often  $k_n = 0$ , in which case  $\hat{\lambda}_n(r) = 1$ , by (11), and consequently  $\hat{\nu}^{(1)} = r$  by (14). Thus, for the starting value of  $\nu$  one should use a number larger than  $r$ .

Our simulation studies exhibited several desirable properties of the iterative procedure. The method converges very fast, yielding the same value for  $\hat{\nu}_{np}$  as the one obtained by maximizing the profile likelihood. The result does not depend on the choice of initial estimate  $\nu^{(0)}$  (provided that it is larger than  $r$ ). One can use  $\nu^{(0)} = r + 1$ , or an estimate based on some parametric assumption (e.g.,  $\hat{\nu}_{exp}$  based on the exponential distribution). Also, in agreement with the preceding discussion, for all data sets with  $m = r$ , the maximum likelihood method failed, yielding  $\hat{\nu} = \infty$ .

## 4. Asymptotics

We shall appeal to the general results of Dewanji et al. (1995) to discuss some asymptotic properties of the nonparametric estimator  $\hat{\nu}_{np}$ , as  $\nu \rightarrow \infty$ . These results are useful when the number of components ( $\nu$ ) is large. In software reliability context, this seems quite

reasonable, as most programs of practical significance contain thousands of lines of code and contain numerous bugs. More generally, asymptotic properties under  $\nu \rightarrow \infty$  are useful (and customary) for estimators of population size or the number of species, see for example, Samuel (1968), Harris (1968) and Chao and Lee (1993).

We start with a simple case, where the true renewal distribution is a discrete probability distribution with a finite sample space  $\mathcal{S} = \{a_1, a_2, \dots, a_N\}$ , where  $a_1 < a_2 < \dots < a_{N-1} < \tau < a_N$ . Let the probabilities of the sample points be  $\pi_1, \dots, \pi_N$ . Converting the probabilities to discrete hazards, as described earlier, we get a parametric model with the parameters being  $(\nu, \lambda)$ . The general asymptotic results of Dewanji et al. (1995) can be applied to this model. Let  $(\hat{\nu}_*, \hat{\lambda}_*)$  denote the MLE of  $(\nu, \lambda)$  under the assumed discrete model. Note that, under this model, the observed renewal times (i.e.,  $x_1, \dots, x_n$ ) must be a subset of  $\mathcal{S}$  and our nonparametric MLE  $(\hat{\nu}_{np}, \hat{\lambda}_{np})$  coincides with  $(\hat{\nu}_*, \hat{\lambda}_*)$ .

Then, following the result in (4), we have, as  $\nu \rightarrow \infty$ ,

$$[\nu^{1/2}(\hat{\lambda}_{np} - \lambda), \nu^{-1/2}(\hat{\nu}_{np} - \nu)] \sim \mathcal{N}(0, \Sigma^*), \quad (15)$$

where

$$\Sigma^* = \begin{pmatrix} I^*(\lambda) & -\delta^* \\ -(\delta^*)' & F(\tau; \lambda)/\bar{F}(\tau; \lambda) \end{pmatrix}^{-1},$$

with  $\delta^* = \frac{\partial}{\partial \lambda} \log \bar{F}(\tau; \lambda)$  and  $I^*(\lambda) =$  the information matrix based on observations from a single process in time  $\tau$ . A consistent estimator of  $I^*(\lambda)$  is

$$-\frac{1}{\hat{\nu}_{np}} \frac{\partial^2 \log L(\hat{\nu}_{np}, \hat{\lambda}_{np})}{\partial \lambda \partial \lambda'}.$$

This is an  $n \times n$  diagonal matrix whose  $l$ th diagonal entry is given by

$$\frac{1}{\hat{\nu}_{np}} \left( \frac{f_l}{\hat{\lambda}_{l,np}^2} + \frac{c_l(\hat{\nu}_{np})}{(1 - \hat{\lambda}_{l,np})^2} \right),$$

where  $\hat{\lambda}_{l,np}$  is the nonparametric MLE of  $\lambda_l$ . The  $l$ th entry of the  $n \times 1$  vector  $\delta^*$  is  $-1/[1 - \lambda_l]$ . In practice,  $\Sigma^*$  may be consistently estimated by replacing  $\lambda$  by  $\hat{\lambda}_{\sim np}$  and  $I^*(\lambda)$  by its consistent estimate as given above. This leads to the estimate of the asymptotic variance of  $\nu^{-1/2}(\hat{\nu}_{np} - \nu)$  as given by

$$\left[ \frac{F(\tau; \hat{\lambda}_{\sim np})}{\bar{F}(\tau; \hat{\lambda}_{\sim np})} - \hat{\nu}_{np} \sum_{l=1}^n \left( \frac{f_l}{\hat{\lambda}_{l,np}^2} + \frac{c_l(\hat{\nu}_{np})}{(1 - \hat{\lambda}_{l,np})^2} \right)^{-1} (1 - \hat{\lambda}_{l,np})^{-2} \right]^{-1}. \quad (16)$$

In summary, if the renewal distribution is discrete with a finite sample space,  $\hat{\nu}_{np}$  is approximately normally distributed with mean  $\nu$  and an estimate of its variance can be obtained using (16). The asymptotic normality of  $\nu^{-1/2}(\hat{\nu}_{np} - \nu)$  also implies that  $\hat{\nu}_{np}/\nu \xrightarrow{P} 1$  as  $\nu \rightarrow \infty$ .

We note that in general, the distribution of  $\hat{\nu}_{np}$  depends only on a part of the renewal distribution  $F(t)$ . Specifically, since no observed renewal time can be larger than  $\tau$ , the (exact) distribution of  $\hat{\nu}_{np}$  is affected only by  $F(t), 0 \leq t \leq \tau$ ; the right tail of  $F$ , over  $(\tau, \infty)$ , has no effect on any statistical property of  $\hat{\nu}_{np}$ . On the other hand, any cdf  $F(t)$  can be approximated arbitrarily closely over the finite interval  $[0, \tau]$  by a distribution with a finite sample space. Thus, our preceding conclusions about the asymptotic distribution of  $\hat{\nu}_{np}$  should hold more generally, i.e., for any renewal distribution. Specifically, for any renewal distribution, the asymptotic distribution of  $\nu^{-1/2}(\hat{\nu}_{np} - \nu)$ , as  $\nu \rightarrow \infty$ , is normal with mean zero and some variance which can be estimated by (16). Our simulation results, reported in the next section, agree with these conclusions.

Table 1: Empirical Evaluation of the Estimators,  $F(\tau) = 0.90$

Distribution	$\nu = 50$				$\nu = 100$				$\nu = 500$			
	$M$	$\hat{\nu}_{exp}$	$\hat{\nu}_{np}$	$s(\hat{\nu}_{np})$	$M$	$\hat{\nu}_{exp}$	$\hat{\nu}_{np}$	$s(\hat{\nu}_{np})$	$M$	$\hat{\nu}_{exp}$	$\hat{\nu}_{np}$	$s(\hat{\nu}_{np})$
Exponential	115.21	49.51	50.56	3.80	230.23	99.52	100.56	5.50	1151.93	499.64	500.63	12.35
	115.00	50.00	50.00	3.92	230.00	100.00	100.00	5.53	1152.00	500.00	500.00	12.33
	10.72	2.80	3.87	1.36	14.96	3.87	5.43	1.17	34.05	8.61	12.21	1.01
Weibull (IFR)	67.50	77.79	54.16	6.68	134.88	154.96	104.66	9.75	674.69	772.84	504.86	23.46
	67.00	76.00	53.00	7.47	135.00	154.00	103.00	10.64	675.00	772.00	503.00	23.55
	5.36	10.31	6.95	4.58	7.74	13.67	9.56	4.81	17.06	29.05	22.41	3.62
Weibull (DFR)	203.78	45.07	49.89	2.93	407.03	90.53	99.84	4.14	2036.47	454.68	499.78	9.23
	203.00	45.00	50.00	2.92	406.00	91.00	100.00	4.12	2036.00	455.00	500.00	9.23
	22.06	2.17	2.98	0.55	30.72	3.07	4.14	0.52	68.91	6.81	9.20	0.51
Gamma (IFR)	84.65	58.76	51.84	4.91	169.51	117.80	101.92	7.22	847.52	590.38	501.62	16.60
	85.00	59.00	51.00	5.29	169.00	118.00	101.00	7.47	848.00	590.00	501.00	16.58
	7.21	4.67	5.03	2.51	10.26	6.31	6.93	2.37	22.64	14.07	16.44	1.73
Gamma (DFR)	159.50	46.03	50.04	3.21	319.33	92.61	100.15	4.56	1595.48	464.55	499.87	10.16
	159.00	46.00	50.00	3.22	319.00	93.00	100.00	4.55	1595.00	465.00	500.00	10.15
	16.25	2.31	3.27	0.76	22.84	3.23	4.61	0.70	51.68	7.11	10.12	0.67
Log- normal	82.69	59.94	51.97	5.01	165.51	120.14	102.12	7.42	827.67	602.30	502.21	17.14
	83.00	60.00	51.00	5.40	166.00	120.00	101.00	7.67	827.00	602.00	502.00	17.10
	6.63	4.69	5.05	2.55	9.52	6.43	7.09	2.42	21.17	14.33	16.74	1.68

## 5. Simulation Results

Here, we report some simulation results assessing and comparing performances of  $\hat{\nu}_{np}$  and  $\hat{\nu}_{exp}$ . The general simulation setup is as follows. We considered three values of  $\nu$ , namely,  $\nu = 50, 100$ , and  $500$  and renewal distributions ( $F$ ) from four families, which are Exponential, Weibull, Gamma and Log-normal. For Weibull and Gamma, we used distributions with both increasing and decreasing failure rate (IFR and DFR). We kept  $\tau$  fixed at  $50$  and varied the parameters of the distributions so that  $F(\tau; \theta)$  equals  $0.6, 0.8$  and  $0.9$ , respectively. Note that  $F(\tau; \theta)$  is the probability of discovering each source and it is a meaningful measure of incompleteness in data induced by the observation period  $[0, \tau]$ .

In each of the 54 cases considered, we generated 10,000 data sets and for each data set with  $m > r$ , we computed  $\hat{\nu}_{np}$ ,  $\hat{\nu}_{exp}$  and an estimate of the standard deviation of  $\hat{\nu}_{np}$  using (16), which we denote by  $s(\hat{\nu}_{np})$ . As we noted in Section 3, if  $m = r$ , then both  $\hat{\nu}_{np}$  and  $\hat{\nu}_{exp}$  are infinite. In our simulation,  $m = r$  occurred only for  $\nu = 50$  and  $F(\tau) = 0.6$  and under Lognormal (60 cases out of 10,000 simulations), Weibull (IFR) (38 cases) and Gamma (IFR) (2 cases) distributions. We also calculated  $\hat{\nu}_{np}$  using both the profile likelihood and the iterative method, for comparing the two methods. As expected, the estimates yielded by two methods were very close. While both methods took very little time for computation, the iterative method was generally faster. We used  $\hat{\nu}_{np}$ , calculated using the profile likelihood method, to assess the behavior of our nonparametric method.

To summarize the simulation results, we calculated the mean, median and standard deviations of  $\hat{\nu}_{np}$ ,  $\hat{\nu}_{exp}$ ,  $s(\hat{\nu}_{np})$  and  $M$ , based on all cases with  $m > r$ . The results are reported in Tables 1 - 3, where for each distribution, the first row gives the averages, the second row the medians and the third row the standard deviations. We observe the following features in our simulation results. As one would expect, for each  $\nu$  and each renewal distribution, both estimators perform better as  $F(\tau; \theta)$  increases. The nonparametric MLE is nearly unbiased in all cases; for IFR distributions it has a small positive bias, which decreases as  $\nu$  increases. The standard deviation of  $\hat{\nu}_{np}$ , relative to  $\nu$ , decreases as  $\nu$  increases. The estimated standard deviation of  $\hat{\nu}_{np}$  based on (16) is generally close to its simulated values, although its standard deviation is large for the two IFR and the

Table 2: Empirical Evaluation of the Estimators,  $F(\tau) = 0.80$

Distribution	$\nu = 50$				$\nu = 100$				$\nu = 500$			
	$M$	$\hat{\nu}_{exp}$	$\hat{\nu}_{np}$	$s(\hat{\nu}_{np})$	$M$	$\hat{\nu}_{exp}$	$\hat{\nu}_{np}$	$s(\hat{\nu}_{np})$	$M$	$\hat{\nu}_{exp}$	$\hat{\nu}_{np}$	$s(\hat{\nu}_{np})$
Exponential	80.58	49.74	51.10	6.30	161.10	99.65	100.95	8.82	804.77	499.98	501.24	19.59
	80.00	49.00	51.00	6.17	161.00	99.00	101.00	8.69	805.00	500.00	501.00	19.54
	9.03	4.79	6.32	1.90	12.75	6.58	8.94	1.72	28.28	14.52	19.46	1.62
Weibull (IFR)	53.29	91.55	56.19	11.61	106.48	179.64	106.16	16.10	532.47	886.64	505.50	35.85
	53.00	88.00	54.00	11.43	106.00	177.00	104.00	16.01	532.00	884.00	503.00	35.59
	5.12	21.18	12.51	7.65	7.33	25.85	15.76	6.58	16.08	52.09	35.06	5.03
Weibull (DFR)	122.09	41.97	50.18	4.75	243.88	84.38	100.22	6.67	1219.51	423.10	499.70	14.79
	122.00	42.00	50.00	4.67	243.00	84.00	100.00	6.62	1219.00	423.00	499.00	14.76
	15.29	3.21	4.88	0.94	21.84	4.45	6.66	0.90	48.64	9.91	14.81	0.88
Gamma (IFR)	62.51	64.85	52.92	8.48	125.04	129.41	102.94	12.01	625.31	645.32	502.60	26.49
	63.00	64.00	52.00	8.33	125.00	129.00	102.00	11.84	625.00	644.00	502.00	26.39
	6.32	9.18	8.85	3.89	9.05	12.37	11.86	3.32	19.81	26.27	26.50	2.84
Gamma (DFR)	105.34	43.71	50.42	5.15	210.51	87.64	100.27	7.21	1053.13	440.14	500.27	15.99
	105.00	44.00	50.00	5.06	210.00	88.00	100.00	7.16	1053.00	440.00	500.00	15.96
	12.58	3.43	5.12	1.12	18.22	4.86	7.25	1.11	40.06	10.89	15.93	1.06
Log- normal	58.40	72.39	53.82	9.47	116.89	144.32	103.94	13.38	584.80	718.35	503.24	29.66
	58.00	71.00	52.00	9.34	117.00	143.00	102.00	13.25	585.00	717.00	502.00	29.50
	5.67	11.41	9.59	4.75	8.11	15.31	13.28	4.27	18.02	32.23	29.47	3.36

Table 3: Empirical Evaluation of the Estimators,  $F(\tau) = 0.60$

Distribution	$\nu = 50$				$\nu = 100$				$\nu = 500$			
	$M$	$\hat{\nu}_{exp}$	$\hat{\nu}_{np}$	$s(\hat{\nu}_{np})$	$M$	$\hat{\nu}_{exp}$	$\hat{\nu}_{np}$	$s(\hat{\nu}_{np})$	$M$	$\hat{\nu}_{exp}$	$\hat{\nu}_{np}$	$s(\hat{\nu}_{np})$
Exponential	45.82	50.98	53.37	12.71	91.76	100.81	102.96	16.92	458.17	500.56	502.38	36.43
	46.00	50.00	52.00	11.52	92.00	100.00	101.00	16.18	458.00	499.00	501.00	36.14
	6.77	10.73	13.46	5.89	9.55	13.88	17.36	4.83	21.53	29.59	36.43	4.27
Weibull (IFR)	35.28	123.97	58.52	28.39	70.44	230.54	113.74	34.04	352.44	898.40	511.97	65.71
	35.00	106.00	55.00	20.38	70.00	214.00	107.00	28.84	352.00	900.00	506.00	63.97
	4.59	67.50	33.62	33.21	6.52	75.41	39.82	26.78	14.75	15.46	66.24	13.75
Weibull (DFR)	58.61	38.33	51.67	9.41	117.40	76.53	101.48	12.81	586.46	382.66	501.76	27.97
	58.00	38.00	51.00	8.88	117.00	76.00	101.00	12.47	586.00	382.00	501.00	27.82
	9.47	5.80	9.66	2.92	13.43	7.86	12.80	2.61	29.78	17.44	28.54	2.48
Gamma (IFR)	38.46	81.07	58.03	19.51	76.70	155.31	107.38	24.65	383.34	755.54	507.16	50.99
	38.00	75.00	53.00	15.99	77.00	150.00	103.00	22.49	383.00	751.00	504.00	50.22
	5.21	30.99	23.84	18.35	7.51	34.48	26.52	11.78	16.52	63.75	51.16	8.22
Gamma (DFR)	55.20	40.42	51.85	9.96	110.43	80.72	101.92	13.58	551.64	402.27	501.45	29.47
	55.00	40.00	51.00	9.35	110.00	80.00	101.00	13.12	551.00	402.00	501.00	29.30
	8.70	6.50	10.33	3.41	12.42	8.78	13.84	3.04	27.75	19.11	29.42	2.76
Log- normal	35.17	123.91	58.82	28.82	70.43	229.53	113.93	34.10	351.80	898.45	511.86	66.25
	35.00	106.00	55.00	20.45	70.00	215.00	106.00	29.10	352.00	900.00	506.00	64.59
	4.59	66.65	34.64	34.02	6.44	70.57	37.54	26.19	14.60	14.63	67.04	13.65

Lognormal distributions. For given  $\nu$  and  $F(\tau; \theta)$ ,  $\hat{\nu}_{np}$  performs better, in terms of bias and standard deviation, if the underlying renewal distribution yields a larger number of events, i.e., a larger value of  $M$ . Note that for fixed  $\nu$  and  $F(\tau; \theta)$ ,  $E(R) = \nu F(\tau; \theta)$  is the same for all renewal distributions. So, the difference in the value of  $M$  comes from the number of repeat (additional) events in the discovered processes, which is determined by how the probability  $F(\tau; \theta)$  is distributed over  $[0, \tau]$ ; a higher concentration of it near 0 is expected to yield a larger number of repeat events, i.e., a larger value of  $(M - R)$ . For example, Weibull (DFR) yields larger  $M$  compared to Weibull (IFR) and better estimates of  $\nu$ . If the true distribution is Exponential, as expected,  $\hat{\nu}_{exp}$  is better than  $\hat{\nu}_{np}$  and the standard deviation of  $\hat{\nu}_{np}$  is about 30% larger than that of  $\hat{\nu}_{exp}$ .

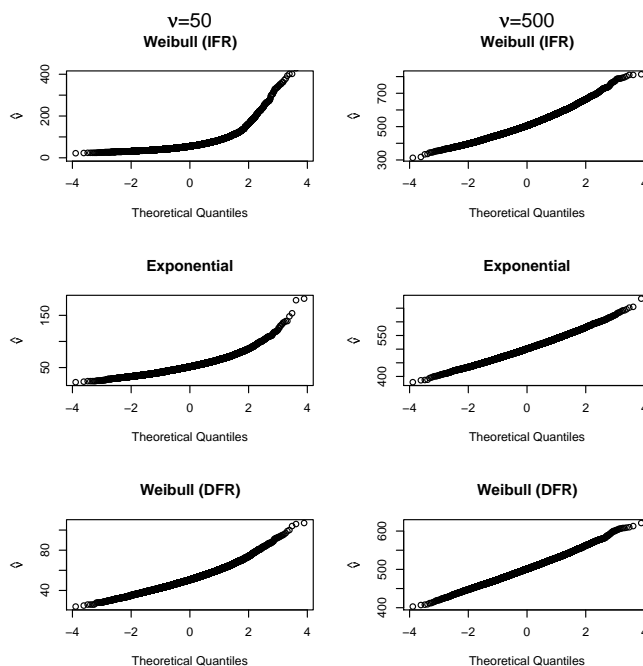


Figure 1: Normal Quantile Plots for different distributions:  $F(\tau = 0.6)$

We also examined the distributions of  $\hat{\nu}_{np}$ . For  $\nu = 50$ , the distributions are visibly positively skewed. For  $\nu = 500$ ,  $\hat{\nu}_{np}$  is approximately normally distributed. For  $\nu = 100$ ,

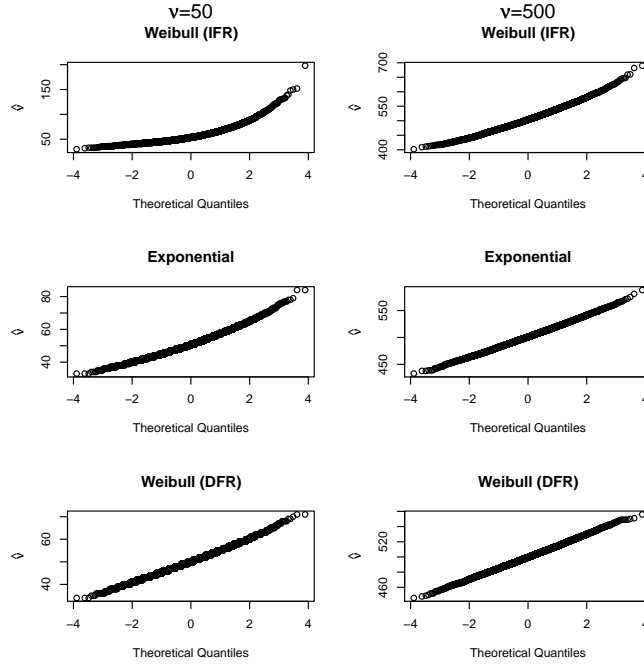


Figure 2: Normal Quantile Plots:  $F(\tau = 0.8)$

approximate normality holds reasonably well for  $F(\tau; \theta) = 0.8, 0.9$ . We present some Q-Q plots in Figures 1 and 2. In Figure 1, Normal quantile plots for  $\hat{\nu}_{np}$  for Weibull (DFR), Exponential and Weibull (DFR), respectively from top to bottom, are presented for  $F(\tau) = 0.6$ . Left column represents the cases for  $\nu = 50$  and the right column for  $\nu = 500$ . Similarly, in Figure 2, we present the similar normal quantile plots but for  $F(\tau) = 0.8$ . As we discussed in Section 4,  $\hat{\nu}_{np}$  converges to a normal distribution as  $\nu \rightarrow \infty$ , but our simulation results show that the convergence is somewhat slow. Also, the convergence is faster for DFR distributions.

As one might expect,  $\hat{\nu}_{exp}$  is biased when the renewal distribution is not Exponential. What may be a new finding from our simulation is that the bias is positive (and substantially so for  $\nu = 50$ ) for IFR distributions and negative for DFR distributions. A rudimentary explanation of this phenomenon is as follows. Under the Exponential model,

$E(R) = \nu[1 - \exp(-\lambda\tau)]$  and  $E(M) = \nu\lambda\tau$ . Equating  $R$  and  $M$  to their expected values and applying the method of moments, an estimator of  $\nu$  is the solution of the equation

$$R = \nu(1 - e^{-M/\nu}). \quad (17)$$

It can be seen that the right hand side of (17) is an increasing function of both  $M$  and  $\nu$ . So, for fixed  $R$ , the solution of (17) for  $\nu$ , say  $\hat{\nu}$ , is a decreasing function of  $M$ , i.e.,  $\hat{\nu}$  would be large if  $M$  is close to  $R$ . Since IFR distributions generate fewer repeat events (and hence smaller  $(M - R)$ ), compared to an Exponential distribution with the same detection probability  $F(\tau; \lambda)$ , we would expect  $\hat{\nu}$  to overestimate  $\nu$  when the underlying distribution is IFR. Similarly,  $\hat{\nu}$  is expected to be negatively biased if the renewal distribution is DFR. Finally, we should note that while  $\hat{\nu}$  and  $\hat{\nu}_{exp}$  are not the same, they should be fairly close, especially because one of the estimating equation, viz.,  $M = \nu\lambda\tau$ , is common to the two methods.

## 6. Concluding Remarks

In this work, we suggest a nonparametric method to estimate the number of components in a superposition of independent and identically distributed renewal processes, where the common renewal distribution is assumed arbitrary. As it turns out, the method involves much less computation than a parametric method, except under exponential renewal distribution, and yields a consistent estimator with asymptotically normal distribution. But, the most attractive feature of our estimator is that it is fairly unbiased. When an assumed parametric model is correct, one naturally expects the parametric MLE to have smaller variance than a nonparametric estimator and our numerical results also exhibit this phenomenon for exponential renewal distribution. In that case, the ratio of the two

standard deviations, viz.,  $sd(\hat{\nu}_{np})/sd(\hat{\nu}_{exp})$ , is about 1.25 for  $F(\tau) = .6$ , 1.33 for  $F(\tau) = .8$  and 1.40 for  $F(\tau) = .9$ , indicating that the efficiency loss decreases as  $F(\tau)$  decreases. We believe, in many practical situations  $F(\tau)$  is small, in which case the robustness of  $\hat{\nu}_{np}$  well justifies the price of moderately higher variance.

It should be noted that the nonparametric method is not of much help for predicting future events, especially discovering new processes, because our nonparametric method yields an estimate of  $\bar{F}(\tau)$ , see equation (9), but no estimate of the distribution  $F$  beyond time  $\tau$  is available. However, this would be a common feature (and deficiency) of all nonparametric methods. In software reliability applications, it is of natural interest to estimate the reliability of the software after the testing period. Specifically, one may wish to estimate the probability of failure-free operation of the software for a length of time  $t$ , after removing all errors detected during testing. Theoretically, this probability is given by  $[\bar{F}(\tau + t)/\bar{F}(\tau)]^{\nu-r}$ , but it cannot be estimated in a nonparametric framework, lacking an estimate of  $\bar{F}(\tau + t)$ . In contrast, the reliability is estimable in a parametric set-up, where the parametric form of the distribution enables one to estimate  $\bar{F}(\tau + t)$ . Nevertheless, one may consider  $\nu - r$ , the remaining number of errors in the software, as a measure of reliability, the estimate of which is readily available along with its variance estimate.

## 7. References

1. Becker, N.G. (1984). Estimating population size from capture recapture experiments in continuous time. *Australian Journal of Statistics* **26**, 1-7.
2. Chao, A. and Lee, S.-M. (1993). Estimating population size for continuous-time

- capture-recapture models via sample coverage. *Biometrical Journal* **35**, 29–45.
3. Dewanji, A., Nayak, T.K. and Sen, P.K. (1995). Estimating the number of components of a system of superimposed renewal processes, *Sankhya A* **57**, 486–499.
  4. Kaplan, E.L. and Meier, P. (1958). Nonparametric estimation from incomplete observations, *Journal of American Statistical Association* **53**, 457–481.
  5. Lloyd, C. J. (1994). Efficiency of martingale methods in recapture studies. *Biometrika* **81**, 305–315.
  6. Nayak, T.K.(1988). Estimating population size by recapture debugging, *Biometrika* **75**, 113–120.
  7. Nayak, T.K. (1991). Estimating the number of component processes of a superimposed process, *Biometrika* **78**, 75–81.
  8. Wilson, K.R. and Anderson, D.R. (1995). Continuous-time capture-recapture population estimation when capture probabilities vary over time. *Environmental and Ecological Statistics* **2**, 55–69.
  9. Xi, L., Yip, P.S.F., and Watson, R. (2007). A unified likelihood-based approach for estimating population size in continuous-time capture-recapture experiments with frailty. *Biometrics* **63**, 228–236.
  10. Yip, P., Huggins, R.M. and Lin, D.Y. (1996). Inference for capture-recapture experiments in continuous time with variable capture rates. *Biometrika* **83**, 477–483.
  11. Yoshida, O.S., Leite J.G. and Bolfarine, H. (1996). Bayes' estimation of the number

of component processes of a superimposed process. *Prob. Eng. Inform. Sci.* **10**, 443–461.