

NONPARAMETRIC BAYES ESTIMATOR OF SURVIVAL FUNCTIONS FOR DOUBLY/INTERVAL CENSORED DATA

Mai Zhou

University of Kentucky

Abstract: The non-parametric Bayes estimator with Dirichlet process prior of a survival function based on right censored data was considered by Susarla and Van Ryzin (1976) and many others. We obtain the non-parametric Bayes estimator of a survival function when data are right, left or interval censored. The resulting Bayes estimator with Dirichlet process prior has an explicit formula. In contrast, there is no explicit formula known for the non-parametric maximum likelihood estimator (NPMLE) with such data. In fact, we show that the NPMLE with doubly/interval censored data cannot, in general, be the limit of Bayes estimators for *any* sequence of priors. Several examples are given, showing that the NPMLE and the non-parametric Bayes estimator may or may not be the same, even when the prior is 'non-informative'.

Key words and phrases: Dirichlet process prior, non-informative prior, NPMLE, square error loss.

1. Introduction, Notation and Preliminary

Suppose lifetimes X_1, \dots, X_n are non-negative and i.i.d. with a distribution $F(\cdot)$. However, these lifetimes are subject to censoring. In the case of right censoring, we only observe

$$Z_i = \begin{cases} X_i, & \text{if } X_i \leq C_i, \\ C_i, & \text{if } X_i > C_i, \end{cases} \quad \text{and} \quad \Delta_i = \begin{cases} 1, & \text{if } X_i \leq C_i, \\ 0, & \text{if } X_i > C_i, \end{cases} \quad (1.1)$$

where C_i are the (right) censoring times.

A generalization of right censoring is double censoring (Chang and Yang (1987), Gu and Zhang (1993)). In the case of double censoring we only observe

$$Z_i = \begin{cases} X_i, & \text{if } Y_i \leq X_i \leq C_i, \\ C_i, & \text{if } X_i > C_i, \\ Y_i, & \text{if } X_i < Y_i, \end{cases} \quad \text{and} \quad \Delta_i = \begin{cases} 1, & \text{if } Y_i \leq X_i \leq C_i, \\ 0, & \text{if } X_i > C_i, \\ 2, & \text{if } X_i < Y_i. \end{cases} \quad (1.2)$$

Here (C_i, Y_i) , $i = 1, \dots, n$, are the left and right censoring times, with $C_i > Y_i$. Let the observations be arranged such that Z_1, \dots, Z_k are the uncensored

observations, i.e., $\Delta_1 = 1, \dots, \Delta_k = 1$. Notice that (Z_1, \dots, Z_k) is (X_1, \dots, X_k) , while Z_{k+1}, \dots, Z_n are the (either right or left) censored observations.

In the Bayesian estimation of $F(\cdot)$, we need not make assumptions about the distributions of the left and right censoring times C_i and Y_i . The calculations are conditioned on the observed censoring times. Thus the observations can be described in three parts Z_1, \dots, Z_k where $X_i = Z_i$; Z_{k+1}, \dots, Z_m where $X_i > Z_i$; and Z_{m+1}, \dots, Z_n , where $X_i < Z_i$.

Next we discuss interval censored data. The current status data, or case 1 interval censored data, consist of an observed ‘‘inspection’’ time T_i and the information whether X_i is larger than or less than T_i (the status of X_i , see Huang and Wellner (1996)):

$$T_i, \quad \Delta_i = \begin{cases} 0, & \text{if } X_i > T_i, \\ 2, & \text{if } X_i < T_i. \end{cases}$$

Usually the ‘‘inspection’’ times T_i are assumed i.i.d. Similar to the discussion above, this i.i.d. assumption does not make a difference in the Bayesian analysis and therefore the current status data is a special case of (1.2), where all the observations are either left or right censored, i.e., $k = 0$.

In case 2 of interval censoring, we assume X_1, \dots, X_k are observed exactly (k non-random, and possibly zero), and only the observations X_{k+1}, \dots, X_n are interval censored. Then we observe $n - k$ intervals. With some abuse of notation, they are denoted by $[L_j, Z_j)$ for $j = k + 1, \dots, n$. We know that $L_j \leq X_j < Z_j$.

Again notice that we do not need to make assumptions about the distribution of L_j or Z_j . Therefore this fits both case 2 and case k of interval censoring in Huang and Wellner (1996). To make the notation consistent with the doubly censored case, we let $Z_1 = X_1, \dots, Z_k = X_k$ for directly observable outcomes, interval censored outcomes are $[L_j, Z_j)$ for $j = k + 1, \dots, n$. Notice that when $[L_j, Z_j) = [a, \infty)$, one has right censored data, and when $[L_j, Z_j) = [0, a)$, left censored data.

In Bayesian analysis, the probability $F(\cdot)$ is random. We assume in this paper that $F(\cdot)$ is distributed as a Dirichlet process with parameter α , a measure on the real line. Under the Dirichlet process prior assumption, the probability measure $P(A) = \int_A dF$ has the following property: given any partition of real line A_1, \dots, A_u , the joint distribution of the random vector $(P(A_1), \dots, P(A_u))$ has a Dirichlet distribution with parameter given by $\alpha(A_1), \dots, \alpha(A_u)$. For more discussion and properties of Dirichlet process prior, see Ferguson (1973), Susarla and Van Ryzin (1976) and Ferguson, Phadia and Tiwari (1993). Another possibility is to work with the cumulative hazard functions $H(t)$. A beta process

prior on the space of the cumulative hazard function was introduced by Hjort (1990). While using a beta process prior for right censored data works well, it has no advantage over the Dirichlet process prior for doubly/interval censored data: the likelihood of the data does not simplify by using the hazard function with doubly censored data.

Using squared error loss, Susarla and Van Ryzin (1976) obtained the Bayes estimator for $F(\cdot)$ under a Dirichlet process prior when data are only subject to right censoring. They also showed that when the weight parameter, α , of the Dirichlet process prior approaches zero, the non-parametric Bayes estimator reduces to the Kaplan-Meier estimator, the NPMLE. Some later papers studied the consistency of the Bayes estimator (Susarla and Van Ryzin (1978)) and the posterior distribution (Ghosh and Ramamoorthi (1995)). Huffer and Doss (1999) used Monte Carlo methods to compute the nonparametric Bayes estimator.

We obtain the Bayes estimator of $1 - F(\cdot)$ when data are subject to both right and left censoring, or are subject to interval censoring. The large sample properties of this Bayes estimator are not discussed here, though it is not unreasonable to expect that it is consistent. However, we show that for *any* sequence of priors the nonparametric Bayes estimators under squared error loss *cannot* always converge to the corresponding NPMLE with doubly censored data. This is a bit surprising since, in most cases, MLEs are limits of Bayes estimators.

The Bayes estimator we obtain is more complicated than those with only right censored data, especially when there are many left censored or interval censored observations. Nevertheless, it has an explicit formula that can be easily programmed. In contrast, the nonparametric maximum likelihood estimator (NPMLE) in the case of doubly censored data or interval censored data does not have an explicit formula and its computation requires an iterative method. See Turnbull (1974), Chen and Zhou (2003) and Fay (1999). Besides, the Bayes estimator is always uniquely defined while the NPMLE is often only defined up to an equivalent class. This non-uniqueness of the NPMLE makes many important statistics like the mean estimator difficult to define. The Bayes estimator is also smoother than the NPMLE. On the other hand, there are consistency results for the NPMLE (Gu and Zhang, (1993), Groeneboom and Wellner (1992) and Huang and Wellner (1996)) but we know very little of the consistency of the Bayes estimators beyond the right censored, R^1 data case. In fact, R. Pruitt gave an example of an inconsistent Bayes estimator with Dirichlet process prior for right censored data in R^2 .

To minimize the amount of new notation, we follow Susarla and Van Ryzin's (1976), hereafter SV, and we use their convention that all observations are positive. Obviously we can extend this to the case where observations have support in $(-\infty, \infty)$ without much difficulty.

2. Bayes Estimator with Right, Left/Interval Censored Observations

The Bayes estimator of $1 - F(\cdot)$ under squared error loss of SV is the conditional expectation of $1 - F$ given all the observations. Similar to SV the conditional expectation is computed in two steps: first, given all the uncensored observations we find the conditional *distribution* of $1 - F$; second, given all the censored observations we compute the conditional *expectation*, where the distribution of those lifetimes before censoring is given in the first step.

The following theorem specifies the conditional distribution of $F(\cdot)$ given all the uncensored observations, which accomplishes the first step.

Theorem 1. *The posterior distribution of the random probability measure P given $(\Delta_1 = 1, Z_1), \dots, (\Delta_k = 1, Z_k)$ is the Dirichlet process with parameter $\beta = \alpha + \sum_{i=1}^k \delta_{Z_i}$, where δ_a is a unit measure on the point a .*

Proof. The proof of this theorem is similar to SV (1976) and Ferguson (1973). We only sketch the proof for the doubly censored case. Furthermore, we only give those calculations that differ from the proof of Theorem 4 of SV (1976), the rest of the proof is the same as theirs and is not repeated here.

From (1) of Chang (1990), the probability of $(\Delta = 1, X = u)$ is $(S_C(u) - S_Y(u))dP(X \leq u) = dG(u)$, say. Recall the marginal distribution of X is $\alpha(u)/\alpha(R^+)$. We compute

$$\begin{aligned} & \int_{[\Delta=1, Z \in A]} D(\cdot | \alpha(B_1) + \delta_u(B_1), \dots, \alpha(B_l) + \delta_u(B_l)) dG(u) \\ &= \int_{[u \in A]} D(\cdot | \alpha(B_1) + \delta_u(B_1), \dots, \alpha(B_l) + \delta_u(B_l)) (S_C(u) - S_Y(u)) d \frac{\alpha(u)}{\alpha(R^+)} \\ &= \sum_{j=1}^l D(y_1, \dots, y_l | \alpha_1^{(j)}, \dots, \alpha_l^{(j)}) \int_{[u \in A \cap B_j]} (S_C(u) - S_Y(u)) d \frac{\alpha(u)}{\alpha(R^+)}. \end{aligned}$$

On the other hand,

$$\begin{aligned} & \mathcal{P}\{P(B_i) \leq y_i, i = 1, \dots, l; \Delta = 1, Z \in A\} \\ &= \int_{u=0}^{\infty} \mathcal{P}\{P(B_i) \leq y_i, i = 1, \dots, l; X \in [u, u + du) \cap A\} (S_C(u) - S_Y(u)) \\ &= \sum_{j=1}^l \int \frac{\alpha(B_j \cap A \cap [u, u + du))}{\alpha(R^+)} (S_C(u) - S_Y(u)) D(y_1, \dots, y_l | \alpha_1^{(j)}, \dots, \alpha_l^{(j)}) \\ &= \sum_{j=1}^l D(y_1, \dots, y_l | \alpha_1^{(j)}, \dots, \alpha_l^{(j)}) \int_{B_j \cap A} \frac{S_C(u) - S_Y(u)}{\alpha(R^+)} d\alpha(u), \end{aligned}$$

which is same as above.

Now, the conditional expectation of $1 - F(u) = P[u, \infty)$ is computed given the remaining $n - k - 1$ censored observations: $Z_{k+1}, \Delta_{k+1}, \dots, Z_n, \Delta_n$ in the doubly censored case; $[L_{k+1}, Z_{k+1}), \dots, [L_n, Z_n)$ in the interval censored case. Notice the original X_{k+1}, \dots, X_n is now a random sample from a Dirichlet process with parameter β . Let E_β denote the expectation with respect to this Dirichlet process.

2.1. Bayes estimator with one interval/left censored observation

To fix ideas and enhance readability, we first present in detail the Bayes estimator with many right censored observations but only one interval censored observation, denoted by $[L_w, Z_w)$. (If $L_w = 0$ then this is left censored.) The general case with many interval/left censored observations will be given later.

As in SV Corollary 1, the conditional expectation, E_β , of $1 - F(u) = P[u, \infty) = P(X \geq u)$ given all the right censored data and one interval censored observation is

$$\hat{S}_D(u) = \frac{E_\beta\{P[u, \infty)P[L_w, Z_w) \prod_{\text{right-censored}} P[Z_i, \infty)\}}{E_\beta\{P[L_w, Z_w) \prod_{\text{right-censored}} P[Z_i, \infty)\}}.$$

This is also the desired Bayes estimator of $1 - F(u)$. We abbreviate the subscript of *right-censored* to $r - c$ and *left-censored* to $l - c$ and *interval-censored* to $i - c$. Straightforward calculation yields

$$\begin{aligned} \hat{S}_D(u) &= \frac{E_\beta P[u, \infty)\{P[L_w, \infty) - P[Z_w, \infty)\} \prod_{r-c} P[Z_i, \infty)}{E_\beta\{P[L_w, \infty) - P[Z_w, \infty)\} \prod_{r-c} P[Z_i, \infty)} \\ &= \frac{E_\beta\{P[u, \infty)P[L_w, \infty) \prod_{r-c} P[Z_i, \infty)\} - E_\beta\{P[u, \infty)P[Z_w, \infty) \prod_{r-c} P[Z_i, \infty)\}}{E_\beta\{P[L_w, \infty) \prod_{r-c} P[Z_i, \infty)\} - E_\beta\{P[Z_w, \infty) \prod_{r-c} P[Z_i, \infty)\}} \\ &= \frac{E_\beta \textcircled{1} - E_\beta \textcircled{2}}{E_\beta \textcircled{3} - E_\beta \textcircled{4}}. \end{aligned}$$

The last four expectations are all of the same type and can be computed explicitly by the Lemma below.

Given a set of positive numbers $0 < a_{k+1} < a_{k+2} < \dots < a_m < \infty$, consider the partition of R^+ into intervals $[0, a_{k+1}), [a_{k+1}, a_{k+2}), \dots, [a_m, \infty)$. By Theorem 1 the random vector $P[0, a_{k+1}), P[a_{k+1}, a_{k+2}), \dots, P[a_m, \infty)$ has a Dirichlet distribution with parameter vector $(\beta_{k+1}, \dots, \beta_{m+1})$ where $\beta_{k+1} = \beta[0, a_{k+1}), \dots, \beta_{m+1} = \beta[a_m, \infty)$. The measure β is given as before by $\beta = \alpha + \sum_{\text{uncensored}} \delta_{Z_i}$.

Lemma 1. (Susarla and Van Ryzin) *With the notation above, we have*

$$E_\beta \prod_{i=k+1}^m P[a_i, \infty) = \prod_{i=0}^{m-k-1} \left(\frac{i + \sum_{j=0}^i \beta_{m+1-j}}{i + \beta(R^+)} \right) = \prod_{i=0}^{m-k-1} \left(\frac{i + \beta[a_{m-i}, \infty)}{i + \beta(R^+)} \right).$$

Proof. This is essentially Lemma 2 (a) of SV (1976) with some extra simplifications.

When $\alpha(R^+) = 0$ the expression on the right hand side of Lemma 1 is still well defined unless there are no uncensored observations in the sample. In Example 2 of Section 3, there are no uncensored observation in the sample and we do not discuss the limit of the Bayes estimator as $\alpha(R^+) \rightarrow 0$ there.

Remark. It is clear from the definition of β that when $\alpha(R^+) \rightarrow 0$, β is integer valued. This implies that the expectation in Lemma 1 has a rational number value (finite product of rational numbers) as $\alpha(R^+) \rightarrow 0$.

2.2. Many interval/left censored observations

When the data contain many interval and many right censored observations, the Bayes estimator of $1 - F(u) = P(X \geq u)$ given all the data (censored or uncensored) is

$$\hat{S}_D(u) = \frac{E_\beta\{P[u, \infty) \prod_{i-c} P[L_w, Z_w) \prod_{r-c} P[Z_i, \infty)]\}}{E_\beta\{\prod_{i-c} P[L_w, Z_w) \prod_{r-c} P[Z_i, \infty)\}}. \quad (2.1)$$

When data contains many left and many right censored observations, the Bayes estimator of $1 - F(u)$ is

$$\hat{S}_D(u) = \frac{E_\beta\{P[u, \infty) \prod_{l-c}[1 - P[Z_w, \infty)] \prod_{r-c} P[Z_i, \infty)\}}{E_\beta\{\prod_{l-c}[1 - P[Z_w, \infty)] \prod_{r-c} P[Z_i, \infty)\}}. \quad (2.2)$$

Because left censored observation is a special case of interval censored observation as pointed out in the previous section, we only present in detail below the Bayes estimator with many interval/right censored observations.

Let us recall the identity

$$\prod_{i=1}^m (b_i - a_i) = \sum y_1 \cdots y_m, \quad (2.3)$$

where y_i is either b_i or $-a_i$ and the summation is over all possible 2^m choices. The integer m is defined as $m = \#\{i - c\}$ = number of interval censored observations.

By using (2.3), we can write $\prod_{i-c} P[L_w, Z_w) = \prod_{i-c} \{P[L_w, \infty) - P[Z_w, \infty)\} = \sum P_1 \cdots P_m$, where each P_w is either $P[L_w, \infty)$ or $-P[Z_w, \infty)$, and the summation is over all 2^m different choices.

To make the expression more specific we introduce some notation. Define vectors $\xi = (\xi_1, \dots, \xi_m)$ where each $\xi_i =$ either 0 or 1. Given m interval censored observations, $[L_i, Z_i)$, we define 2^m sets of numbers $\{c_i(\xi), i = 1, \dots, m\}$ where $c_i(\xi) = L_i$ if $\xi_i = 0$ otherwise $c_i(\xi) = Z_i$. With each set $\{c_i(\xi), i = 1, \dots, m\}$, associate a sign: if the set contains an even number of Z_i 's then the sign is positive, if the set contains odd number of Z_i 's the sign is negative.

With these definition we can write $\prod_{i-c} P[L_w, Z_w) = \sum P_1 \cdots P_m = \sum_{\xi} \pm \prod_{i=1}^m P[c_i(\xi), \infty)$, where the summation is over all 2^m different ξ 's, and \pm is the associated sign.

Finally, we define new sets of numbers by adding r ($r = \#\{r - c\}$) right censored observations Z_1, \dots, Z_r to $\{c_i(\xi), i = 1, \dots, m\}$: $\{b_j(\xi), j = 1, \dots, m + r\} = \{c_i(\xi), i = 1, \dots, m\} \cup \{Z_1, \dots, Z_r\}$. For any sets of real numbers b_1, \dots, b_k , we denote by $b_{(-1)}, \dots, b_{(-k)}$ the reversely ordered numbers (descending). So, $b_{(-i)}^+(\xi), i = 1, \dots, m + r$ is a set of $m + r$ numbers ordered from largest to smallest.

With these sets of numbers defined, the denominator of (2.1) can be written as

$$\sum E_{\beta} \left(P_1 \cdots P_m \times \prod_{r-c} P[Z_i, \infty) \right) = \sum_{\xi} \left(\pm \prod_{i=1}^{m+r} \frac{i - 1 + \beta[b_{(-i)}^+(\xi), \infty)}{i - 1 + \beta(R^+)} \right),$$

where the summation is over 2^m different ξ 's. We can similarly compute the numerator of (2.1) except there is one more term, $P[u, \infty)$, included with the right censored observations. Define $\{b_j^+(\xi)\} = \{c_i(\xi), i = 1, \dots, m\} \cup \{Z_1, \dots, Z_r, u\}$.

Theorem 2. *The nonparametric Bayes estimator of the survival function $S(u) = 1 - F(u)$ with right censored and interval censored data under a Dirichlet process prior is*

$$\begin{aligned} \hat{S}_D(u) &= \frac{\sum E_{\beta} \{P_1 \cdots P_m \times \prod_{r-c} P[Z_i, \infty) \times P[u, \infty)\}}{\sum E_{\beta} \{P_1 \cdots P_m \times \prod_{r-c} P[Z_i, \infty)\}} , \\ &= \frac{\sum_{\xi} (-1)^{\sum \xi_s} \prod_{i=1}^{m+r+1} \frac{i - 1 + \beta[b_{(-i)}^+(\xi), \infty)}{i - 1 + \beta(R^+)}}{\sum_{\xi} (-1)^{\sum \xi_s} \prod_{i=1}^{m+r} \frac{i - 1 + \beta[b_{(-i)}(\xi), \infty)}{i - 1 + \beta(R^+)}} . \end{aligned} \quad (2.4)$$

The sums in (2.4) are over all 2^m possible ξ 's.

Admittedly the two summations above involves 2^m terms when there are m interval censored observations. Also, in the summations, there are both positive and negative terms that will cancel to a large extend. Rounding errors will be magnified if we use (2.4) directly. Our purpose here is to show that an explicit formula exists for the Bayes estimator. Simplifications/alternative formulae are desirable and will be pursued in the future.

Remark. From Lemma 1 and Theorem 2, we can infer that the limit of the Bayes estimator (2.4) when the α measure approaches zero is a step function, at least for $u < \text{maximum observed value}$. This is because all the E_{β} involved

will be step functions according to Lemma 1. We can also infer that when the α measure approaches zero, the Bayes estimator (2.4) is a rational, since the E_β involved are all rational.

3. Examples

The examples presented here are hand-calculated or are obtained by using software we developed (Example 2 and the NPMLE in Example 1). We pay close attention to the limit of the Bayes estimator when $\alpha \rightarrow 0$ in the Dirichlet prior (non-informative prior), and compare the estimator with the NPMLE. The software used here are packaged as **R** (<http://www.r-project.org/>) packages and can be found at <http://www.ms.uky.edu/~mai/research/>. The software for computing NPMLE is also available at this site.

To minimize additional notation, we recycle the notation used by SV as much as possible. Assume $Z_{(k+1)}, \dots, Z_{(m)}$ are the ordered, distinct censored (both right and left/interval) times among the sample (1.1). Assume there are no ties among the left/interval and right censored observations (but ties within right censored observations are allowed). At each censored observation $Z_{(i)}, k+1 \leq i \leq m$, let λ_i be the number of right censored observations that equal $Z_{(i)}$. Thus if there are two right censored observations equal to $Z_{(i)}$, then $\lambda_i = 2$. If $Z_{(j)}$ is a left censored observation then $\lambda_j = 0$. To make the notation consistent, we define $Z_{(k)} = 0$ and $Z_{(m+1)} = \infty$.

Let $N(u)$ be the number of uncensored and right censored observations that are larger than or equal to u , i.e., $N(u) = \sum_{j=1}^m I_{[Z_j \geq u]} + \sum_{i=k+1}^m \lambda_i I_{[Z_{(i)} \geq u]}$, and let $N^+(u) = N(u+)$.

We reproduce SV's Bayes estimator (based only on the uncensored and right censored observations of the sample (1.1)) in a slightly modified form: For $Z_{(l)} \leq u < Z_{(l+1)}$ with $k \leq l \leq m+1$,

$$\hat{S}(u) = \frac{\alpha(u, \infty) + N^+(u)}{\alpha(R^+) + N^+(0)} \times \prod_{j=k+1}^l \left\{ \frac{\alpha[Z_{(j)}, \infty) + N(Z_{(j)})}{\alpha[Z_{(j)}, \infty) + N(Z_{(j)}) - \lambda_j} \right\}. \quad (3.1)$$

We have changed two things: we added the nodes $Z_{(j)}$ for left/interval censored observations, though with zero λ_j 's; n is replaced by $N^+(0)$.

Example 1. Here is an example with one left censored observation and four right censored observations. These are the data used by SV (1976) but with an added left censored observation at $Z = 4$.

The ordered observations with their censoring indicators are listed below in Table 1.

Table 1. Data with one left and four right censored observations.

Z'_i 's :	0.8	1.0	2.7	3.1	4	5.4	7.0	9.2	12.1
Δ :	1	0	0	1	2	1	0	1	0

Let the Bayes estimator of SV based only on uncensored and right censored observations be $\hat{S}(u)$, i.e., as defined in (3.1). Our estimator that takes into account one left censored observation can be written as follows. For $u > Z_{left} = 4$, after tedious but straightforward simplification we get

$$\hat{S}_D(u) = \hat{S}(u) \times \frac{\frac{\alpha[0,1]+1}{\alpha(R^+)+9} + \frac{\alpha[1,2.7]}{\alpha[1,\infty)+7} \times \frac{\alpha[1,\infty)+8}{\alpha(R^+)+9} + \frac{\alpha[2.7,4)+1}{\alpha[2.7,\infty)+6} \times \frac{\alpha[1,\infty)+8}{\alpha(R^+)+9} \times \frac{\alpha[2.7,\infty)+7}{\alpha[1,\infty)+7}}{\frac{\alpha[0,1)+1}{\alpha(R^+)+8} + \frac{\alpha[1,2.7)}{\alpha[1,\infty)+6} \times \frac{\alpha[1,\infty)+7}{\alpha(R^+)+8} + \frac{\alpha[2.7,4)+1}{\alpha[2.7,\infty)+5} \times \frac{\alpha[1,\infty)+7}{\alpha(R^+)+8} \times \frac{\alpha[2.7,\infty)+6}{\alpha[1,\infty)+6}}.$$

For u in other time intervals, the estimator can be similarly expressed as the product of $\hat{S}(u)$ and some other term, the details are omitted. The plot of the Bayes estimator is given in Figure 1.

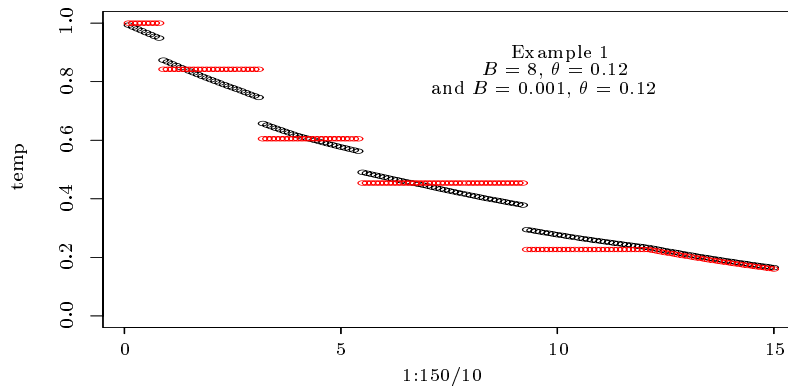


Figure 1. Plot for Example 1.

Next we compute the limit of the Bayes estimator. When $\alpha \rightarrow 0$, the SV estimator, $\hat{S}(u)$, has as a limit the Kaplan-Meier estimator S_{KM} . For $9.2 \leq u < 12.1$, the limit of our estimator is $S_{KM} \times (70/81) = (7/8) \times (4/5) \times (3/4) \times (1/2) \times (70/81) = 0.2268519$. For u in other intervals the limit can be computed similarly.

We plot the estimator with $\alpha(u, \infty) = B \exp(-\theta u)$. The plot shows estimators for $B = 8, \theta = 0.12$ and $B = 0.001, \theta = 0.12$. The latter is indistinguishable in appearance with the limit just calculated.

Computation of the NPMLE for doubly censored data can be done by EM type iteration (see Turnbull (1974) and Chen and Zhou (2003)). For the data in Table 1 we obtain the values in Table 2:

Table 2. NPMLE and limit of Bayes estimator for data in Table 1.

t	0-0.8	0.8-3.1	3.1-5.4	5.4-9.2	9.2-12.1
NPMLE	1	0.8457284	0.6028477	0.4521358	0.2260679
Limit Bayes	1	0.8425926	0.6049383	0.4537037	0.2268519

The differences between the NPMLE and the limit of the nonparametric Bayes estimator are small but real. The likelihood of the distribution in Table 2 is larger than those of the limit of the Bayes estimator: 3.70674×10^{-5} vs. 3.704924×10^{-5} .

Example 2. We took the first ten observations from the breast cosmesis data with radiation of Finkelstein and Wolfe, as reported by Fay (1999). Out of the ten, there are four right censored observations, four interval censored observations and two left censored observations (i.e., interval censored with left-ends as 0). Data: $[45, \infty)$, $[6, 10)$, $[0, 7)$, $[46, \infty)$, $[46, \infty)$, $[7, 16)$, $[17, \infty)$, $[7, 14)$, $[37, 44)$, $[0, 8)$.

We computed the nonparametric Bayes estimator with $\alpha(u, \infty) = B \exp(-\theta u)$. The resulting estimator with $B = 8$ and $\theta = 0.3$ is computed using the software we developed and is plotted in Figure 2.

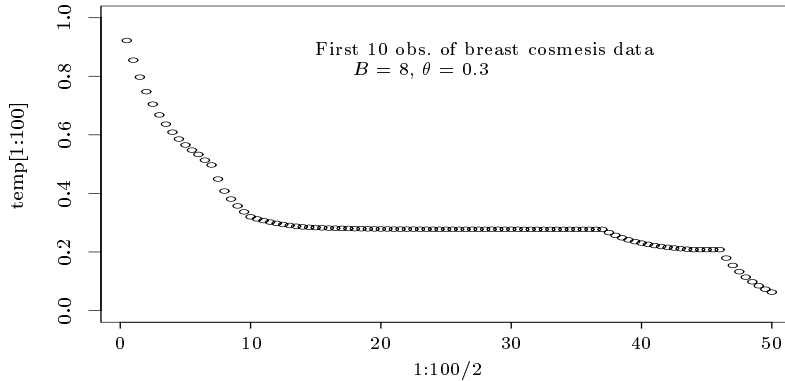


Figure 2. Plot for Example 2.

In the following two examples, the Bayes estimators are obtained with formula (2.4) and then we let $\alpha(R^+) \rightarrow 0$ to obtain the limit. The NPMLE's are also calculated, not by software but analytically.

Example 3. Here we took a small example with one left and one right censored observation. The NPMLE and the limit of the non-parametric Bayes estimator turn out to be exactly the same.

Table 4. Data with one left and one right censored observation.

Z'_i 's :	$Z_{(1)}$	$Z_{(2)}$	$Z_{(3)}$	$Z_{(4)}$	$Z_{(5)}$
Δ :	1	0	2	1	1
Jump of $1 - \hat{F}(u)$	0.4	0	0	0.3	0.3

Example 4. The order of two censoring indicators in the above table are switched and the NPMLE is different from the limit of Bayes estimator. The limit of Bayes estimator is not self-consistent either. To calculate the NPMLE, we first note for this data the NPMLE $F(\cdot)$ has only three jumps at $Z_{(1)}, Z_{(3)}$ and $Z_{(5)}$. Denote the jumps size by p_1, p_2, p_3 . By symmetry we must have $p_1 = p_3$. Using the constraint $\sum p_i = 1$, we can reduce the likelihood, $L = p_1(p_2 + p_3)p_2(p_1 + p_2)p_3$, to a function of p_2 only. Straightforward calculation shows that $p_2 = \sqrt{5}/5$ maximizes the likelihood. Therefore $p_1 = p_3 = (5 - \sqrt{5})/10$, which is the entry 0.2763932 in the table.

$Z_i^s :$	$Z_{(1)}$	$Z_{(2)}$	$Z_{(3)}$	$Z_{(4)}$	$Z_{(5)}$
$\Delta :$	1	0	1	2	1
Jump of limit Bayes	0.28000	0	0.44000	0	0.28000
Jump of NPMLE	0.2763932	0	0.4472136	0	0.2763932

Remark. Example 4 shows that with positive probability, the NPMLE $1 - \hat{F}(\cdot)$ for doubly/interval censored data can take irrational values. A closer look at the example provides some insight as why the estimators are different, as described in the next section.

Remark. Example 3 and 4 reveal two different situations. The difference is that left and right censored data overlap (right censored observation is smaller then the left censored observation) in Example 4. The overlap in Example 3 is not real, since there is no probability mass inside the overlap.

4. Limit of Bayes and NPMLE

In this section we formally summarize some results concerning the limit of the Bayes estimator and the NPMLE in the doubly/interval censored data case. The argument below is valid for *any* prior, not just the Dirichlet process prior.

Theorem 3. *Suppose a sequence of priors $\pi_v; v = 1, 2, \dots$, is such that the (non-parametric) Bayes estimators $1 - \hat{F}_v(\cdot)$ under squared error loss, converge to the Kaplan-Meier estimator whenever the data has only right censoring. Then this same sequence of Bayes estimators cannot converge, in general, to the NPMLE for interval/doubly censored data.*

Proof. The Bayes estimator under squared error loss can be written as

$$1 - \hat{F}_v(u) = \frac{E_\pi P[u, \infty) L_F(\text{data})}{E_\pi L_F(\text{data})},$$

where $L_F(\text{data})$ is the likelihood of the data when its distribution is F . The assumption of the Theorem for right censored data says that, as $v \rightarrow \infty$, we

always have

$$\frac{E_\pi\{P[u, \infty) \prod_{r-c} P[x_i, \infty) \prod_{uncensor} P(\{x_j\})\}}{E_\pi \prod_{r-c} P[x_i, \infty) \prod_{uncensor} P(\{x_j\})} \rightarrow 1 - F_{K-M}(u). \quad (4.1)$$

Notice the Kaplan-Meier estimator, $F_{K-M}(u)$, is always rational valued.

Now we look at a particular sample configuration with just one left censored observation, for example the data in Example 4. The Bayes estimator for these data can be written as

$$\frac{E_\pi P[u, \infty) P(\{Z_1\}) P[Z_2, \infty) P(\{Z_3\}) P[0, Z_4) P(\{Z_5\})}{E_\pi P(\{Z_1\}) P[Z_2, \infty) P(\{Z_3\}) P[0, Z_4) P(\{Z_5\})}.$$

Let us use the notation $P(Z) = P(\{Z\})$, and $P(Z^+) = P[Z, \infty)$. Write $P[0, Z_4) = 1 - P[Z_4, \infty) = 1 - P(Z_4^+)$ and expand to get

$$= \frac{E_\pi P(Z_1) P(Z_3) P(Z_5) P(Z_2^+) P(u^+) - E_\pi P(Z_1) P(Z_3) P(Z_5) P(Z_2^+) P(Z_4^+) P(u^+)}{E_\pi P(Z_1) P(Z_3) P(Z_5) P(Z_2^+) - E_\pi P(Z_1) P(Z_3) P(Z_5) P(Z_2^+) P(Z_4^+)}. \quad (4.2)$$

If we divide the numerator of (4.2) by $E_\pi P(Z_1) P(Z_3) P(Z_5) P(Z_2^+) P(u^+)$ then, as $v \rightarrow \infty$, the numerator will converge to the limit $1 - [1 - F_{K-M}^*(Z_4)]$ according to (4.1). Here the Kaplan-Meier estimator is based on three uncensored observations: Z_1, Z_3, Z_5 and two right censored observations, Z_2 and u .

Similarly, if we divide the denominator of (4.2) by $E_\pi P(Z_1) P(Z_3) P(Z_5) P(Z_2^+)$ then it has the limit $1 - [1 - F_{K-M}^{**}(Z_4)]$, where the Kaplan-Meier estimator is based on three uncensored observations, Z_1, Z_3, Z_5 and one right censored observation Z_2 .

In other words, multiply (4.2) by

$$\frac{E_\pi P(Z_1) P(Z_3) P(Z_5) P(Z_2^+)}{E_\pi P(Z_1) P(Z_3) P(Z_5) P(Z_2^+) P(u^+)} \quad (4.3)$$

to produce a rational limit. The factor (4.3) itself has a rational limit as $v \rightarrow \infty$ ($= [1 - F_{K-M}^{**}(u)]^{-1}$). This implies that the limit of (4.2), as $v \rightarrow \infty$, is

$$\frac{F_{K-M}^*(Z_4)}{F_{K-M}^{**}(Z_4)} \times [1 - F_{K-M}^{**}(u)].$$

But that cannot be the NPMLE as Example 4 shows the NPMLE is irrational.

This also serves as the proof for the interval censored case, since the left censored observation is just $[0, Z_{(4)})$ interval censored.

Corollary 1. *There is no sequence of priors, π_v such that the resulting sequence of Bayes estimators under squared error loss, $1 - F_v(\cdot)$, always converges to the NPMLE in the interval/doubly censored data case.*

Proof. Suppose, to the contrary, there is such a sequence of priors. Since right censoring is a special case of double/interval censoring (zero left censoring or all intervals are of the form $[a_i, \infty)$), this sequence of estimators must converge to the Kaplan-Meier estimator with such data. But by Theorem 3, such sequence cannot converge to the NPMLE for doubly/interval censored data case in general.

5. Discussion

The formula (2.4) has 2^m terms when there are m interval censored observations. While we do not have a formal proof that the computation of the Bayes estimator cannot be reduced to polynomial order, it is not hard to see that the computation is equivalent to

$$\int \cdots \int \prod_j \left(\sum_{r=1}^j x_r \right) \prod_j (1 - \sum_{r=1}^j x_r) (1 - \sum_j x_j)^{\beta_m} \prod x_j^{\beta_j} \prod dx_j$$

on the region $x_j > 0$ and $\sum x_j \leq 1$.

Remark. The irrational value of NPMLE with doubly censored data also implies that the EM algorithm, if started from the Kaplan-Meier estimator, cannot converge in a finite number of steps in general.

Acknowledgement

I thank C. Srinivasan for many helpful discussions.

References

- Chang, M. (1990). Weak convergence of a self-consistent estimator of the survival function with doubly censored data. *Ann. Statist.* **18**, 391-404.
- Chang, M. and Yang, G. L. (1987). Strong consistency of a nonparametric estimator of the survival function with doubly censored data. *Ann. Statist.* **15**, 1536-1547.
- Chen, K. and Zhou, M. (2003). Non-parametric hypothesis testing and confidence intervals with doubly censored data. *Lifetime Data Anal.* **9**, 71-91.
- Fay, M. (1999). Splus functions for nonparametric estimate and test for interval censored data. <http://lib.stat.cmu.edu/S/interval.tar.gz>
- Ferguson, T. (1973). A Bayesian analysis of some nonparametric problems. *Ann. Statist.* **1**, 209-230.
- Ferguson, T., Phadia, E. G. and Tiwari, R. C. (1993). Bayesian nonparametric inference. In *Current Issues in Statistical Inference: Essays in honor of D. Basu* (Edited by M. Ghosh and P. K. Pathak), 127-150. IMS Lecture Notes Monograph series **34**.

- Ghosh, J. K. and Ramamoorthi, R. V. (1995). Consistency of Bayesian inference for survival analysis with or without censoring. In *Analysis of Censored Data* (Edited by H. Koul and J. V. Deshpande), 95-103. IMS Lecture Notes Monograph Series **27**.
- Groeneboom, P. and Wellner, J. (1992). *Information Bounds and Nonparametric Maximum Likelihood Estimation*. Birkhauser, Basel.
- Gu, M. G. and Zhang, C. H. (1993). Asymptotic properties of self-consistent estimators based on doubly censored data. *Ann. Statist.* **21**, 611-624.
- Hjort, N. (1990). Nonparametric Bayes estimators based on beta processes in models for life history data. *Ann. Statist.* **18**, 1259-1294.
- Huang, J. and Wellner, J. (1996). Interval censored survival data: a review of recent progress. In *Proceedings of the First Seattle Symposium in Biostatistics* (Edited by D. Y. Lin and T. Fleming), 123-170.
- Huffer, F. and Doss, H. (1999). Software for Bayesian analysis of censored data using mixtures of Dirichlet priors. Preprint.
- Kaplan, E. and Meier, P. (1958), Non-parametric estimator from incomplete observations. *J. Amer. Statist. Assoc.* **53**, 457-481.
- Susarla, V. and Van Ryzin, J. (1976), Nonparametric Bayesian estimation of survival curves from incomplete observations. *J. Amer. Statist. Assoc.* **71**, 897-902.
- Susarla, V. and Van Ryzin, J. (1978), Large sample theory for a Bayesian nonparametric survival curves estimator based on censored samples. *Ann. Statist.* **6**, 755-768.
- Turnbull, B. (1974). Nonparametric estimation of a survivorship function with doubly censored data. *J. Amer. Statist. Assoc.* **69**, 169-173.

Department of Statistics, College of Art & Sciences, University of Kentucky, 849 Patterson Office Tower, Lexington, KY 40506-0027, U.S.A.

E-mail: mai@ms.uky.edu

(Received July 2001; accepted September 2003)