

점진적 임의증단법에서 생존함수의 비모수적 추정에 관한 연구*

박 병구** 이 광호***

요 약

신뢰수명의 검정이나 임상실험에서 대상물에 대한 관측치를 충분히 얻기 위해서는 많은 시간과 경비가 필요하거나 현실적인 제한으로 인하여 관측이 불가능한 경우가 흔히 있다. 이러한 이유로 점진적 임의증단법에 의하여 얻어진 관측치를 이용한 생존함수의 추정은 현실적으로 매우 중요하다고 하겠다.

연구에서는 점진적 임의증단된 자료를 기초로 스플라인함수를 이용하여 생존함수의 비모수적 최우추정량을 제안하고 그것의 성질과 효율성을 비교, 연구한다.

1. 서 론

여러가지 신뢰성이나 생존분석연구에서 생존함수를 추정하는 문제가 대두되는데, 그것은 필요한 대상물에 대한 임상실험이나 수명검정을 통하여 자료를 수집함으로써 시작된다. 그러나, 현실적인 여러가지 제약조건(시간, 경비, 대상물의 특징 등)으로 인하여 대상물에 대한 실험 혹은 관측이 중단되는 경우가 많이 있는데, 이런 경우의 자료를 일반적으로 불완전 자료라 부르며 중단되는 특징에 따라 여러가지 다른 이름으로 불리고 있다.

불완전 자료를 이용한 생존함수의 추론은 문제의 다양성과 복잡성으로 인하여 Kaplan과 Meier(1958)의 연구이래 비모수적 방법으로 폭넓게 연구되고 있다. 특히 임의 증단된 자료를 이용한 생존함수의 추론은 Breslow와 Crowley(1974)를 포함하여 Foldes, Rejto와 Winter(1980), Suzuki(1985), Tsai(1986), Wong(1987) 등에 의해 연구되었다.

실제로 임의 증단에 의한 표본추출은 공학에서의 신뢰성 분석문제나 의학분야에서의

* 이 연구는 88년도 한국과학재단 연구비 지원에 의한 결과임 (881-0105-006-1)

** 경북대학교 자연과학대학 통계학과

*** 영남대학교 이과대학 통계학과

생존분석연구 등에서 널리 적용되고 있다. 그러나 이러한 분야에서의 관측치는 대개가 시간이 경과함에 따라 순차적으로 얻어지는 특징을 지니기 때문에 임의 중단법으로 표본추출하는 데에도 몇 가지 애로점이 나타날 수 있다. 즉, 현실적인 여건으로 인하여 너무 빨리 중단하게 되면 표본크기가 너무 적게 되어 추정을 어렵게 하거나 신뢰도를 떨어지게 할 수 있으며, 극단적으로는 전혀 관측치를 얻지 못하게 될 수도 있다. 이에 반하여 충분한 관측치를 얻기위하여 중단시간을 충분히 길게 정하면 그 시간에 수반되는 경비가 과다하게 소요될 수 있고 또 필요이상의 실험 단위의 파괴내지는 학생이 수반될 수 있으며, 또한 실험 단위에 대한 사전정보가 전혀 없는 경우에는 이 충분한 시간동안에 실험단위의 특징에 따라 극히적은 관측치만을 얻게 될지도 모르는 것이다.

이러한 이유로 인하여, 임상실험과 수명검정문제 등에서 점진적 임의 중단법 (progressively random censoring scheme)의 필요성이 주장되어 왔다. 이 방법에는 관측치가 실험의 시작부터 수집되고 그 결과를 순차적으로 누적하여 추정 혹은 검정에 이용한다. 만약 실험의 어느 단계에서 현 상태의 누적 관측치가 어떤 분명한 통계적 결정을 보장할 수 있다면 실험은 현 단계에서 마치게 된다. 따라서 점진적 임의 중단법은 실험비용과 실험단위의 학생 모두의 감소와 더불어 비교적 짧은 실험기간이 요구된다.

점진적 임의 중단법으로 Cohen(1963, 1965)은 모수모형에 대한 생존함수의 최우추정량을 제안했고 Chatterjee와 Sen(1973)은 일반적인 순위검정법을 제안하였다. Davis(1978)는 생존함수의 비모수적 베이즈추정량을 제안하였다. 점진적 임의 중단법의 다른 응용에 대해 Halperin과 Ware(1974), Sen(1976), Gardiner와 Sen(1978), Gardiner와 Susarla(1982)등이 연구했었다.

본 연구에서는 점진적 임의 중단법으로 생존함수의 새로운 비모수적추정량을 제안하고 그 성질을 연구하고자 한다. 즉, 스플라인함수를 사용하여 위험률함수를 근사적으로 구하고 그것과 생존함수의 관계를 이용하여 생존함수의 비모수적 최우추정량을 구하고자 한다. 만약 실험종결단계(m)가 완전한 표본(n)과 일치하면 제안된 통계량은 Whittmore와 Keller(1986)의 통계량과 일치한다. 끝으로, 점진적 임의 중단법하에서 제안된 통계량의 성질을 모의 실험을 통해 비교. 연구하며 실제의 수치적 예를 들어 적용시켜 보고자 한다.

2. 점진적 임의 중단법

이 절에서는 점진적 임의 중단법에서의 표본추출과정과 생존함수에 대한 정의를 한다. 생존시간 T 와 중단시간 C 는 각각 연속분포함수 $F(\cdot)$ 와 $G(\cdot)$, 그리고 확률밀도함수 $f(\cdot)$ 와 $g(\cdot)$ 를 각각 가지며 T 와 C 는 서로 독립이라 한다. 실제 관측시간 Y 와 중단 지시함수 δ

는 각각의 실험단위에 대하여 다음의 식으로 표시된다.

$$Y = \min(T, C),$$

$$\delta = I[T \leq C] = \begin{cases} 1, & \text{if } T \leq C \text{ (사망)}, \\ 0, & \text{if } T > C \text{ (중단).} \end{cases}$$

따라서, $(Y_1, \delta_1), (Y_2, \delta_2), \dots, (Y_n, \delta_n)$ 는 n 개의 항목을 실험할 때의 관측 가능한 자료이다. $Y_{(1)}, Y_{(2)}, \dots, Y_{(n)}$ 은 Y_1, Y_2, \dots, Y_n 의 순서통계량이고 δ_i^* 은 $Y_{(i)}$, $i=1, 2, \dots, n$, 과 연관된 중단지시함수라 하자. 지금 m 단계에서 실험이 종결되었다 하자. 여기에서 $m \in \{1, 2, \dots, n\}$ 이다. 그러면 실제로 기록된 자료는 $(Y_{(1)}, \delta_1^*), (Y_{(2)}, \delta_2^*), \dots, (Y_{(m)}, \delta_m^*)$ 와 나머지 $n-m$ 개의 항목의 중단시간과 생존시간이 모두 $Y_{(m)}$ 을 초과한다는 것뿐이다. 즉 $Y_{(j)} > Y_{(m)}$, $j=m+1, \dots, n$, 이다. 식을 간단하게 하기 위해 $Y_{(i)}$ 와 δ_i^* 를 Y_i 와 δ_i 로 쓰기로 한다.

이와 같이 관측된 자료를 이용하여 다음과 같이 정의된 생존함수와 위험률함수, 누적위험률함수를 추정하고자 한다.

$$\text{생존함수} : S(t) = P(T > t) = 1 - F(t). \quad (2.1)$$

$$\text{위험률함수} : \lambda(t) = \frac{f(t)}{1 - F(t)} \quad (2.2)$$

$$\begin{aligned} \text{누적위험률함수} : \Lambda(t) &= \int_0^t \lambda(\mu) d\mu \\ &= - \int_0^t d \log(1 - F(\mu)). \end{aligned} \quad (2.3)$$

여기서 생존함수와 누적위험률함수의 관계는 다음과 같다.

$$S(t) = \exp\{-\Lambda(t)\}. \quad (2.4)$$

Kaplan과 Meier(1958), Nelson(1969)는 생존함수 $S(t)$ 의 비모수적 추정량을 임의 중단법(random censoring)된 정보를 기초로 다음과 같은 \hat{S}_{KM} 과 \hat{S}_{NA} 로 각각 제안하고 그 성질을 연구하였다.

$$\hat{S}_{KM}(t) = \prod_{y \leq t} \left(\frac{n-i}{n-i+1} \right) \delta$$

$$\hat{S}_{NA}(t) = \exp\{-\hat{\Lambda}_{NA}(t)\}.$$

$$\text{단 } A_{NA}(t) = \sum_{y \leq t} \frac{\delta}{n-i+1} \text{이다.}$$

[비고 2.1] 시각 t 에서 $n-i$ 가 클 때는 $S_{KM}(t)$ 와 $S_{NA}(t)$ 는 거의 일치하며 모두 다음과 같은 단점을 내포하고 있다.

- (1) 위험률함수의 추정량을 직접 구할 수 없다.
- (2) 대단히 유용한 정보가 될지도 모르는 정확한 중단시간을 사용하고 있지 않고 있다.

3. 생존함수의 비모수적 최우추정량

이 절에서는 스플라인함수를 사용하여 생존함수의 비모수적 최우추정량을 구하고자 한다. 먼저 $[0, T]$ 상에서 생존함수를 추정하기 위해 점진적 임의 중단법으로 얻어진 자료 (y_i, δ_i) , $i=1, 2, \dots, n$ 의 우도함수를 이용하고자 한다. 각 표본의 확률은 다음과 같다.

$$\begin{aligned} \Pr(y_i=t, \delta_i=0) &= \Pr(C_i=t, T_i > C_i) \\ &= g(t) (1-F(t)), \\ \Pr(y_i=t, \delta_i=1) &= \Pr(T_i=t, T_i \leq C_i) \\ &= f(t) (1-G(t)). \end{aligned}$$

위의 식들은 다음과 같이 표현할 수 있다.

$$\Pr(y_i=t, \delta_i) = [f(t)\{1-G(t)\}]^{\delta_i} [g(t)\{1-F(t)\}]^{1-\delta_i}$$

$m < i \leq n$ 에 대한 자료로부터 이용 가능한 유일한 정보는 그들의 중단시간과 생존시간이 모두 y_m 을 초과한다는 것뿐이다. 즉 $y_i > y_m$, $i = m+1, \dots, n$, 이다. 따라서 (y_i, δ_i) , $i = 1, 2, \dots, n$, 의 우도함수는 다음과 같다.

$$\begin{aligned} &\prod_{i=1}^m \{1-G(y_i)\}^{\delta_i} g(y_i)^{1-\delta_i} \{1-G(y_m)\}^{n-m} \\ &\cdot \prod_{i=1}^m f(y_i)^{\delta_i} \{1-F(y_i)\}^{1-\delta_i} \{1-F(y_m)\}^{n-m} \end{aligned} \quad (3.1)$$

F 와 G 는 서로 독립이므로 $G(\cdot)$ 와 $g(\cdot)$ 에는 관심의 대상인 어떤 모수도 포함되어 있지 않는다. 따라서 식 (3.1)의 첫 항은 무시할 수 있으며 다음의 결과를 얻는다.

[정리 3.1] 절진적 임의증단법에서 자료 $(Y_1, \delta_1), (Y_2, \delta_2), \dots, (Y_m, \delta_m)$ 과 $Y_i > Y_m, i = m+1, \dots, n$, 의 우도함수는

$$L = \prod_{i=1}^m f(y_i)^{\delta_i} \{1-F(y_i)\}^{1-\delta_i} \{1-F(y_m)\}^{n-m}$$

이다.

식 (2.1)과 (2.2)에서의 생존함수와 위험률함수의 관계식으로부터 정리 2.1의 우도함수는

$$L(\lambda) = \prod_{i=1}^m \lambda(y_i)^{\delta_i} S(y_i) S(y_m)^{n-m}$$

으로 나타낼 수 있으며 식 (2.3)와 (2.4)으로부터 $L(\lambda)$ 의 로그-우도함수는

$$\log L(\lambda) = \sum_{i=1}^m \left\{ \delta_i \log \lambda(y_i) - \int_0^{y_i} \lambda(u) du \right\} - (n-m) \int_0^{y_m} \lambda(u) du \quad (3.2)$$

가 된다.

이제부터 식(3.2)에서의 위험률함수 $\lambda(\cdot)$ 을 스플라인함수를 이용하여 근사적으로 구해 보자.

[정의 3.1] $0 = t_0 < t_1 < \dots < t_k < t_{k+1} = T$ 라 하고 $\theta_j(t), j = 0, 1, \dots, k+1$, 는 알려진 t 의 함수라 하면 t 의 스플라인함수 λ_Q 를

$$\lambda_Q(t) = \sum_{j=0}^{k+1} \theta_j(t) \lambda_j \quad (3.3)$$

로 정의한다. 단 $\lambda_j \equiv \lambda(t_j)$ 이다. 실제로 스플라인함수 λ_Q 는 위험률함수를 보간하고 있으며 알려진 함수 $\theta_j(t)$ 에 의해 결정되어진다.

[정리 3.2] 정리 3.1의 조건에서 위험률함수와 누적위험률함수, 생존함수의 최우추정량은 각각 다음과 같이 주어진다.

$$\hat{\lambda}_Q(t) = \sum_{j=1}^k \theta_j(t) d_j \gamma_j^{-1}, \quad (3.4)$$

$$\hat{\Lambda}_Q(t) = \sum_{j=1}^k \alpha_j(t) d_j \gamma_j^{-1}, \quad (3.5)$$

$$S_0(t) = \exp\{-\Lambda_0(t)\}, \quad (3.6)$$

단 $j = 1, 2, \dots, k$ 대해서 d_j 는 시각 t_j 에서의 사망한 항목의 갯수이며 $\theta_j(t)$ 는 시각 t 의 알려진 함수이고, $a_j(t)$ 와 γ_j 는 식 (3.8)과 (3.10)에서 각각 주어진다.

[증명] 식 (3.3)에서 주어진 λ_0 의 정의로 부터 함수 λ_0 의 적분값은 근사적인 누적위험률함수 $\Lambda(t)$ 가 된다. 즉, λ_0 함수의 적분은

$$\begin{aligned} \int_0^t \lambda_0(u) du &= \sum_{j=0}^{k+1} \lambda_j \int_0^t \theta_j(u) du \\ &= \sum_{j=0}^{k+1} a_j(t) \lambda_j \end{aligned} \quad (3.7)$$

이다. 여기에서 계수

$$a_j(t) = \int_0^t \theta_j(u) du \quad (3.8)$$

이며 $\lambda(\cdot)$ 에 의존하지 않는 알려진 함수이다.

여기에서 t_j , $j = 1, 2, \dots, k$, 는 k 개의 서로 다른 사망시간(dead time)으로 잡자. 만약, 로그-우도함수 식 (3.2)에서 $\lambda(\cdot)$ 의 적분을 λ_0 의 적분으로 치환하여 식 (3.8)을 이용하면

$$L(\lambda_0) = \sum_{j=0}^{k+1} (d_j \log \lambda_j - \gamma_j \lambda_j) \quad (3.9)$$

가 된다. 단

$$\gamma_j = \sum_{i=1}^m a_j(y_i) + (n-m)a_j(y_m) \quad (3.10)$$

이다. 식 (3.9)은 평균이 $\gamma_j \lambda_j$, $j = 0, 1, \dots, k+1$, 인 $k+2$ 개의 독립인 포아송 변수 d_j 의 로그-우도함수이며 그 해는

$$\hat{\lambda}_j = d_j \gamma_j^{-1} \quad (3.11)$$

이 된다. 따라서, 이 값 $\hat{\lambda}_j$ 을 식 (3.3)과 (3.7)에서 사용하면 위험률함수와 누적 위험률함수, 생존함수에 대한 최우추정량을 각각 구할 수 있다.

[비고 3.1] 정리 3.2에서 얻어진 추정량들은 식 (3.3)에 주어진 $\theta_j(t)$ 에 전적으로 의존한다. 즉 특별한 함수 $\theta_j(t)$ 를 선택함에 따라 위험률함수와 생존함수의 추정량들이 결정된다.

[정의 3.2] 임의의 측정가능한(measurable) 함수 h 에 대해 함수 $\delta(\cdot)$ 가

$$\int h(x)\delta(\tau-x)dx = h(\tau) \quad (3.12)$$

을 만족할 때 함수 $\delta(\cdot)$ 를 Dirac-delta함수라 한다.

[파를정리 3.1] 정리 3.2의 조건에서 $\theta_j(t) = \delta(t-t_j)$, $j = 0, 1, \dots, k+1$, 이면 누적위험률함수와 생존함수의 최우추정량은 다음과 같이 주어진다.

$$\begin{aligned} \Lambda_{Q0}(t) &= \sum_{t_j \leq t} d_j \gamma_j^{-1} \\ S_{Q0}(t) &= \exp\{-\Lambda_{Q0}(t)\}. \end{aligned} \quad (3.13)$$

단

$$\gamma_j = \sum_{i=1}^m I(y_i \geq t_j) + (n-m)I(y_m \geq t_j)$$

이다.

[증명] 식 (3.5)에서의 스플라인함수는 주어진 조건에 의해

$$\lambda_{Q0}(t) = \sum_{j=0}^{k+1} \delta(t-t_j) \lambda_j$$

이고, 식 (3.8)과 (3.12)로 부터

$$\begin{aligned} a_j(t) &= \int_0^t \delta(u-t_j) du \\ &= \int I(t \geq u) \cdot \delta(u-t_j) du \\ &= I(t \geq t_j) \end{aligned}$$

가 된다. 더우기 식 (3.7)으로부터 γ_j 을 얻을 수 있다. 그러므로 식 (3.4) 및 (3.5), (3.6), (3.12)로 부터 본 정의의 결과를 얻을 수 있다.

[비고 3.2] 식 (3.13) 을 자세히 보면 \hat{S}_{Q0} 는 실험종결단계(m)가 n 과 같고 대등한 생존시간(tied survival time)이 없을 때는 2절에서 소개된 Nelson의 추정량 \hat{S}_{AN} 이 됨을 알 수 있다.

만약, $\theta_j(t)$ 를

$$\theta_j(t) = \begin{cases} I(t_{j-1} < t \leq t_j), & j = 1, 2, \dots, k+1 \\ 0, & j = 0 \end{cases} \quad (3.14)$$

과 같이 정의하면 [파름정리 3.1]의 증명에서와 유사한 방법으로 다음의 결과를 얻을 수 있다.

[파름정리 3.2] 정리 3.2의 조건에서 $\theta_j(t)$ 를 식 (3.5)와 같이 하면 누적위험률 함수와 생존함수의 최우추정량은 각각

$$\hat{\Lambda}_{Q1}(t) = \sum_{j=1}^k a_j(t) d_j \gamma_j^{-1}$$

와

$$\hat{S}_{Q1}(t) = \exp\{-\hat{\Lambda}_{Q1}(t)\}$$

이다. 단 $a_0(t) = 0$ 이고 $a_j(t)$ 와 γ_j 는 $j=1, 2, \dots, k+1$ 에 대해 각각

$$a_j(t) = \begin{cases} 0, & 0 \leq t \leq t_{j-1}, \\ t - t_{j-1}, & t_{j-1} < t \leq t_j, \\ t_j - t_{j-1} \equiv \Delta_j, & t_j < t \leq T, \end{cases} \quad (3.15)$$

와

$$\gamma_j = \sum_{i=1}^m a_j(y_i) + (n-m)a_j(y_m) \quad (3.16)$$

이다.

다음으로, $\theta_j(t)$ 를

$$\theta_j(t) = \begin{cases} I(t_j < t \leq t_{j+1}), & j = 0, 1, \dots, k, \\ 0, & j = k+1, \end{cases} \quad (3.17)$$

과 같이 정의하면 다음의 결과를 얻는다.

[파름정리 3.3] 정리 3.2의 조건에서 $\theta_j(t)$ 를 (3.17)식과 같이 하면 누적위험률 함수와 생존함수의 최우추정량은 각각

$$\hat{\Lambda}_{Q2}(t) = \sum_{j=1}^k a_{j+1}(t) d_j \gamma_j^{-1} \quad \text{와} \quad \hat{S}_{Q2}(t) = \exp\{-\hat{\Lambda}_{Q2}(t)\}$$

이다. 단 $j=1, 2, \dots, k$ 에 대해서 $a_j(t)$ 는 (3.15)에 주어진 것이고, $a_{k+1}(t)=0$ 이며 γ_j 는 다음과 같다.

$$\gamma_j = \sum_{i=1}^m a_{j+1}(y_i) + (n-m)a_{j+1}(y_m).$$

여기에서 Λ_{Q2} 는 우측 연속 계단함수이고 높이는 $d_j \lambda_j^{-1}$ 을 가진다.

마지막으로, $\theta_j(t)$ 를

$$\theta_j(t) = \begin{cases} I[t_0 < t < t_i] \cdot [\Delta_i - a_i(t)]/\Delta_i, & j = 0, \\ I[t_{j-1} \leq t \leq t_j] \cdot a_j(t)/\Delta_j \\ \quad + I[t_j \leq t \leq t_{j+1}] \cdot [\Delta_{j+1} - a_{j+1}(t)]/\Delta_{j+1}, & j = 1, 2, \dots, k, \\ I[t_k \leq t \leq t_{k+1}] \cdot a_{k+1}(t)/\Delta_{k+1}, & j = k+1 \end{cases} \quad (3.18)$$

와 같이 정의하면 다음의 결과를 얻을 수 있다.

[파률정리 3.4] 정리 3.2의 조건에서 식 $\theta_j(t)$ 를 (3.18)와 같이 정의하면 누적위험률함수와 생존함수의 최우추정량은 각각 다음과 같이 주어진다.

$$\Lambda_{Q3}(t) = \sum_{j=1}^k a_j(t) d_j \gamma_j^{-1},$$

$$S_{Q3}(t) = \exp\{-\Lambda_{Q3}(t)\}.$$

단 $a_j(t)$ 와 γ_j 는 각각

$$a_j(t) = a_{j+1}(t) + \frac{1}{2} \left\{ \frac{a_j(t)^2}{\Delta_j} - \frac{a_{j+1}(t)^2}{\Delta_{j+1}} \right\}$$

와

$$\begin{aligned} \gamma_j &= \sum_{i=1}^m \left[a_{j+1}(y_i) + \frac{1}{2} \left\{ \frac{a_j(y_i)^2}{\Delta_j} - \frac{a_{j+1}(y_i)^2}{\Delta_{j+1}} \right\} \right] \\ &\quad + (n-m)[a_{j+1}(y_m) + \frac{1}{2} \left\{ \frac{a_j(y_m)^2}{\Delta_j} - \frac{a_{j+1}(y_m)^2}{\Delta_{j+1}} \right\}] \end{aligned}$$

이다.

4. 적용사례

이 절에서는 부분적인 자료 $\{(y_i, \delta_i); 1 \leq i \leq m, y_i \geq y_m; j = m+1, \dots, n\}$ 을 가지고 점진적 임의 중단하에서 제안한 생존함수의 추정값을 계산한다. 다음의 실제 자료는 Gardiner와 Susarla(1982)에서 인용했다. 위스콘신대학교 중앙종양학그룹에서 흑색종

위의 도표로부터 각 경우에서의 제안된 추정량의 값은 서로 비슷하다. 그리고 KM추정량과 제안된 추정량의 값의 차이도 거의 없음을 볼 수 있다. 즉, 점진적 임의중단법에서 종결단계(■)을 줄여 표본위 크기가 작아도 제안된 추정량들은 이용가치가 있음을 보여주고 있다.

5. 몬테칼로 연구

이제까지 점진적 임의 중단된 자료를 기초로 스플라인함수를 이용하여 생존함수의 몇 가지 추정량을 제안하였으며, 이 절에서는 몬테칼로 방법을 이용하여 제안된 추정량들 \hat{S}_{Q0} , \hat{S}_{Q1} , \hat{S}_{Q2} , \hat{S}_{Q3} 과 KM추정량의 효율성을 비교하고자 한다.

단순하게 나타내기 위해 추정량 S_{Q0} , S_{Q1} , S_{Q2} , S_{Q3} 을 각각 Q_0 , Q_1 , Q_2 , Q_3 이라 하자. 추정량 Q_2 는 우측연속계단함수이고 Q_1 은 좌측연속계단함수인 점을 제외하고는 두 추정량은 거의 같은 형태이다. 즉, Q_1 이 생존함수를 과대 추정한다면 Q_2 는 과소 추정할 것이고 반대로 Q_1 이 과소 추정 한다면 Q_2 는 과대 추정하게 될 것이다. 따라서 여기서는 Q_2 추정량은 고려하지 않겠다.

(표 1)에서는 이 절에서 사용되는 세 가지 형태의 생존분포(F)와 중단분포(G)를 나타냈 다. 생존분포는 각각 위험률함수가 증가, 감소 그리고 상수함수인 경우가 되도록 선택된것이고, 중단분포는 중단되는 관측치가 많은 경우(70%)와 보통인 경우(50%)가 되도록 고려된 것이다.

각 경우에서 표본크기(n)가 20이고 실험종결단계(■)가 16, 18, 20일 때 500번 시행하여 표 2, 3, 4에 나타내었으며 표본크기가 30과 40일 때도 유사한 결과를 얻을 수 있었다. 표의 수치들은 5개 시간구간 $(0.0, 0.5]$, $(0.5, 1.0]$, $(1.0, 1.5]$, $(1.5, 2.0]$, $[0.0, 2.0]$ 에서 시간이 0.005단위로 변할 때마다 편의와 평균제곱오차를 계산하여 평균한 것이다.

표2를 통하여 다음의 결과를 얻을수 있다.

- (1) 생존함수의 침값이 큰 경우에는 점진적 중단법이 대단히 유익한 표본추출법이 되는 것 같다.
- (2) 중단된 관측치가 많을 때는 $(1.0, 2.0]$ 구간의 모든 경우에 대하여

Q3추정량이 Q0와 Q1 보다 더 작은 편의와 평균제곱오차를 가지는 것 같다. 그러나 (0, 0, 0.5]구간에서는 Q1은 더 적은 편의와 평균제곱오차를 가진다.

- (3) 중단된 관측치가 보통일 때는 Q3추정량이 (1.0, 2.0]구간에서 다른 모든 추정량보다 더 작은 편의를 가지고 (0, 0, 1.0]구간에서는 모든 경우에 대해 Q3가 KM추정량보다 더 작은 편의와 평균제곱오차를 가지는 것 같다.
- (4) (1.0, 2.0]구간에서 중단된 관측치가 많을 때는 보다 더 작은 평균제곱오차를 가지고, 중단된 관측치가 보통일 때는 Q3가 KM추정량 보다 더 작은 편의를 가진다.
- (5) 시간구간 (0, 0, 0.5]에서는 실험종결 단계(■)가 작더라도, 즉 자료를 부분적으로 이용하더라도 추정량들 모두는 유사한 편의와 평균제곱오차를 가지고 있다.

표3과 표4를 통하여서도 표2에서 얻은 결과와 유사한 결과들을 찾을 수 있다.

표 1. 수명분포와 중단분포

수명 분포	중단 분포	중단율
Weibull (1, 0.5)	Exponential (0.5)	70
	Exponential (1.0)	50
Weibull (1, 1.0)	Exponential (0.5)	70
	Exponential (1.0)	50
Weibull (1, 1.5)	Exponential (0.5)	70
	Exponential (1.0)	50

6. 결론

본 연구에서는 점진적 임의증단법으로 신뢰함수(생존함수)의 비모수적 추정문제를 다룬다. 스플라인함수를 이용하여 위험함수를 근사적으로 구하고, 생존함수와 위험률함수와의 관계식으로부터 생존함수의 비모수적 최우추정량을 제안하였다. 즉, 몇몇의 특정한 스플라인함수를 가정하여, Nelson형 추정량 및 두 종류의 계단함수형 추정량, 한 종류의 연속형 추정량을 얻었다.

4 절에서 실제의 자료를 가지고 제안된 추정량과 Kaplan-Meier(KM)추정량의 값을 구해 본 결과, 연속형 추정량(Q3)의 값은 실험종결단계(■)가 완전한 자료의 크기(n)보다 작더라도 KM추정량과 차이가 거의 없음을 알 수 있었다.

5 절에서 제안한 추정량 (Q0, Q1, Q3)과 KM추정량의 효율성을 몬테칼로 방법으로 비교하여 다음의 결과를 얻을 수 있었다.

- (1) 연속형 추정량(Q3)는 다른 추정량보다 편의와 최소제곱오차 면에서 대체로 더 유효하였다.
- (2) 실험종결단계(■)가 전체 자료의 크기(n)와 같을 때라도 연속형 추정량 (Q3)은 KM추정량보다 대체로 더 유효하였다.
- (3) 생존함수의 참 값이 비교적 클 때 점진적 임의증단법에서의 추정량들은 실험증단 단계(■)가 자료의 크기(n)보다 작을 때라도 편의와 평균제곱오차들은 거의 변화가 없었다.

참 고 문 헌

1. Breslow, N. and Crowley, J. (1974). A Large Sample Study of Life Table and Product Limit Estimates under Random Censorship. *The Annals of Statistics*, 2, 437-453.
2. Chatterjee, S. K. and Sen, P. K. (1973). Nonparametric Testing under progressive Censoring. *Calcutta Statistical Association, Bulletin.*, 22, 13-50.
3. Cohen, A. C. (1963). Progressively Censored Samples in Life Testing. *Technometrics*, 5, 327-339.
4. Cohen, A. C. (1965). Maximum Likelihood Estimation in the Weibull Distribution based on Complete and on Censored Samples. *Technometrics*, 7, 579-588.

5. Davis, C. E. (1978). A Two Sample Wilcoxon test for Progressively Censored Data. *Communications in Statistics - Theory and Method*, A7, 389-398.
6. Földes, A., Rejtö, L. and Winnnter, B. B. (1980). Strong Consistency Properties of Nonparametric Estimators for Randomly Censored data, I: The Product Limit Estimator. *Periodica Mathematica Hungarica*, 11, 233-250.
7. Gardiner, J. C. and Sen, P. K. (1978). Asymptotic Normality of a Class of Time-sequential Statistics and Applications. *Communications in Statistics-Theory and Method*, A7, 373-388.
8. Gardiner, J. C. and Susarla V. (1982). A Nonparametric Estimator of the Survival Function under Progressive Censoring. *IMS, Lecture notes: Special Topics Meeting on Survival Analysis*, 26-40.
9. Halperin, M. and Ware, J. (1974). Early Decision in a Censored Wilcoxon two Sample Test for Accumulation Survival Data. *Journal of the American Statistical Association*, 69, 414-422.
10. Kaplan, E. L. and Meier P. (1958). Nonparametric Estimation from Incomplete Observations. *Journal of the American Statistical Association*, 53, 457-481.
11. Nelson, W. B. (1969). Harzard Plotting for Incomplete Failure Data. *Journal of Quality Technology*, 1, 27-52.
12. Sen, P. K. (1976). Weak Convergence of Progressively Censored Likelihood Ratio Statstic and Its Role in Asymptotic Theory of Life Testing. *Annals of Statistics*, 4, 1247-1257.
13. Suzuki, K. (1985). Nonparametric Estimation of Lifetime Distributions from a Record of Failures and Follow-ups. *Journal of the American Statistical Association*, 80, 68-72.
14. Tsai, W. Y. (1986). Estimation of Survival Curves from Dependent Censorship Models via a Generalized self-consistent Property whit Nonparametric Bayesian Estimation Application. *The Annals of Statistics*, 14, 238-249.
15. Wang, M. (1987). Product Limit Estimates : A Generalized Maximum Likelihood Study. *Communications in Statistics- Theory and Method*, 16, 3117-3132.
16. Whittemore, A. S. and Keller, J. B. (1986). Survival Estimation using Splines. *Biometrics.*, 42, 496-506.
17. Wong, W. H. (1986). Theory of partial likelihood. *The Annals of Statistics*, 14, 88-123.