

## 임의 중단모형하에서의 평균잔여수명함수의 추정

이 인 석 , 이 우 동<sup>1)</sup>

**요 약** 이 연구에서는 Hjort(1991)에 의해 제안된 누적위험률함수의 비모수적 추정량을 이용하여 무한인 구간까지 정의된 평균잔여수명함수의 추정량을 제안하고 제안된 추정량의 일치성과 점근적 정규성을 밝히고, 모의실험을 통하여 다른 추정량들과 비교하고자 한다.

### 1. 서 론

임의의 기계가  $t$ 시간까지 살았다는 가정하에서 그 기계의 앞으로 남아있는 평균수명시간, 즉, 평균잔여수명(mean residual life)은 의학분야나 공학분야등, 수명시간분석에 관련된 분야에서 관심이 되는 중요한 측도이다. 이러한 평균잔여수명은 생존함수(survival function)의 함수로, 몇몇 연구자들에 의해 생존함수의 비모수적 추정량을 이용한 추정량들이 제안되어 왔다.

임의중단모형(random censoring model)하에서, Yang(1977,1978)은 평균잔여수명에 대한 비모수적 추정량을 제안했으며, 제안된 추정량이 미리 정의된 유한인 구간 상에서 가우시안(Gaussian) 과정으로 수렴함을 밝혔다. 그리고 베이지안 측면에서, Ghorai, Susarla, Susarla 와 Van Ryzin(1980)은 제곱손실오차(squared error loss)를 고려하여 추정량을 제안하였다. Kumazawa(1987)는 Kaplan-Meier추정량을 이용하여 평균잔여수명함수에 대한 추정량을 제안하고 Yang의 연구를 Gill(1983)의 결과를 이용하여 무한인 구간 상으로 확장하여 제안된 추정량의 정규성을 밝혔다.

이 연구에서는 Hjort에 의해 제안된 누적위험률함수의 비모수적 추정량을 이용하여 무한인 구간까지 정의된 평균잔여수명함수의 추정량을 제안하고 제안된 추정량의 점근적 일치성과 점근적 정규성을 밝히고, 모의실험을 통하여 기존의 다른 추정량들과

비교하고자 한다.

## 2. 평균잔여수명함수의 추정

$T_1, T_2, \dots, T_n$ 을  $[0, \infty)$ 상에서 정의된 분포함수(distribution function),  $F$ 로부터 추출된 확률표본(random samples)이라 두자. 그리고 수명함수는

$$\begin{aligned} S(t) &= P\{T_i > t\} \\ &= 1 - F(t) \\ &= \exp(-\Lambda(s)) \end{aligned}$$

라 하자. 여기서  $\Lambda(s)$ 는 누적위험률함수(cumulative hazard function)로서,

$$\Lambda(s) = \int_0^s \lambda(t) dt$$

이고

$$\lambda(t) = \lim_{\Delta t \rightarrow 0} \frac{\Pr\{T \in (t, t + \Delta t) | T > t\}}{\Delta t}$$

이며  $\lambda(t)$ 는 어떤 기계가  $t$ 시간이상 살았다는 가정(조건)하에서 아주 짧은 시간 혹은 순간에 죽을 확률 혹은 위험률(hazard rate)이다. 또,  $C_1, C_2, \dots, C_n$ 을 분포함수로

$1 - G(t)$ 를 갖는 중단확률표본이라 하고  $T_i$ 와  $C_i$ 는 독립이라 가정하자. 그리고

$$X_i \equiv (T_i \wedge C_i) = \min(T_i, C_i), \quad i=1, 2, \dots, n$$

라두고 그 분포함수를  $H(t)$ 라 하면,

$$\begin{aligned} 1 - H(t) &= P\{X_i > t\} \\ &= S(t)G(t) \end{aligned}$$

인 관계가 성립한다.

임의중단 모형에서 관측되는 자료의 형태는  $(X_1, \delta_1), (X_2, \delta_2), \dots, (X_n, \delta_n)$ 이며

여기서

$$\begin{aligned} \delta_i &= I\{T_i < C_i\} \\ &= \begin{cases} 1 & \text{if } T_i < C_i \\ 0 & \text{if } T_i \geq C_i \end{cases} \end{aligned}$$

로서 중단 지시자(indicator function)이다.  $X_{(1)}, X_{(2)}, \dots, X_{(n)}$ 을  $X_1, X_2, \dots, X_n$ 의 순서통계량이라하고  $X_{(i)}$ 의 중단지시자를  $\delta_{(i)}$ 라두자.

위의 임의 중단모형에서 널리 이용되는 수명함수의 비모수적 추정량을 소개하면 아래와 같다. 먼저, Kaplan과 Meier(1958)는 생존함수의 추정량(K-M 추정량)을 아래와 같이 제안하였고

$$\hat{S}_{KM}(t) = \begin{cases} \prod_{i: X_{(i)} \leq t} \left\{ \frac{n-i}{n-i+1} \right\}^{\delta_{(i)}}, & \text{for } t \leq X_{(n)} \\ 0 & \text{if } \delta_{(n)} = 0, \text{ for } t > X_{(n)} \\ \text{undefined} & \text{if } \delta_{(n)} = 1, \text{ for } t > X_{(n)} \end{cases}$$

그후, 여러학자들에 의하여 K-M 추정량에 대한 통계적 성질이 밝혀져 왔는데, 그 중 대표적인 연구자와 연구내용을 간략히 소개해보면, Breslow와 Crowley(1974)는 K-M추정량의 극한분포가 정규분포임을 증명하였고, Peterson(1977)은 K-M 추정량을 두개의 부생존함수로 표시하였으며, Gill(1983)은 K-M 추정량의 점근적성질을 밝혔다. 또, K-M 추정량에 비교될만한 생존함수의 추정량으로는 Nelson(1972)과 Aalen(1978)에 의해 제안된 누적위험률함수에 대한 경험적(empirical) 추정량으로 그들의 추정량은

$$\begin{aligned} \hat{\Lambda}(t) &= \int_0^t \frac{dN(s)}{Y(s)} \\ &= \sum_{i: X_{(i)} \leq t} \frac{\delta_{(i)}}{n-i+1}, \end{aligned}$$

이며, 여기서  $N(t)$ 는  $t$ 보다 작거나 같은 임의중단되지않은 관찰치(uncensored observation)의 개수이며  $Y(t)$ 는  $t$ 보다 작거나 같은 관찰치(observation)의 개수이다. 즉

$$N(t) = \sum_{i=1}^n I_{\{X_{(i)} \leq t, \delta_{(i)} = 1\}}$$

$$Y(t) = \sum_{i=1}^n I_{\{X_{(i)} \geq t\}}$$

이다.

Nelson과 Aalen의 누적위험률함수에 대한 추정량을 이용한 생존함수의 추정량(N-A 추정량)은 분포함수와 누적위험률함수의 관계를 이용하여 다음과 같이 제안할 수 있다.

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





치성과 극한분포를 구하였다.

최근에 Kumazawa(1987)는 임의중단모형에서,  $\tau_F = \sup \{t \mid F(t) < 1\}$ 로 정의하고,

$$e_t^{\tau_F} = \frac{\int_t^{\tau_F} S(s) ds}{S(t)}$$

의 추정문제를 고려하였고 여기서  $\tau_F$ 는  $\infty$ 일 수 있다. 그는 K-M추정량을  $e_t^{\tau_F}$ 에 포함된 생존함수에 대입하여 추정량을 제안하였고 제안된 추정량의 점근분포가 정규분포임을 증명하였다.

본 연구에서는  $e_t^{\tau_F}$ 에 H 추정량을 이용하여 추정량을 제안하고, 점근적 일치성과 정규성을 밝힌다. 먼저

$$T^* = \max_{1 \leq i \leq n} \{X_i\}, \quad \tau_G = \sup \{t \mid G(t) < 1\}$$

그리고

$$\tau_H = \sup \{t \mid H(t) < 1\}$$

이라 정의하자.

만약,  $\tau_G < \tau_F$ 인 관계가 성립한다면  $\tau_G$ 시점을 넘어서는 중단되지 않은 자료의 관측은 불가능할 것이며,  $\tau_G$ 를 넘어서는 시점에서 생존함수 S의 추정은 불가능할 것이다. 그러므로  $\tau_F \leq \tau_G$ 로 가정하면  $\tau_F = \tau_H$ 인 관계가 성립되고, 다음과 같이  $e_t^{\tau_F}$ 의 추정량을 제안하자.

$$\hat{e}_t = \frac{\int_t^{T^*} \hat{S}(s) ds}{\hat{S}(t)}.$$

유사한 방법으로 생존함수의 추정량  $\hat{S}_{KM}(t)$ ,  $\hat{S}_{NA}(t)$ ,  $\hat{S}_{SV}(t)$ 를 이용하여  $e_t^{\tau_F}$ 의 추정량을 각각  $\hat{e}_{t,KM}$ ,  $\hat{e}_{t,NA}$ ,  $\hat{e}_{t,SV}$ 라 두자. 지금부터  $\hat{e}_t$ 의 점근적 성질을 밝혀 보자.

아래의 정리 2.1에서 제안된 추정량의 약 일치성(consistency)을 보였다.

정리 2.1 만약

$$\sqrt{n} \int_{T^*}^{\tau_F} S(s) ds \rightarrow^p 0, \quad \text{as } n \rightarrow \infty.$$

그러면

$$\sup_{0 \leq t \leq T^*} |\hat{e}_t - e_t^{\tau_F}| \rightarrow^p 0, \quad \text{as } n \rightarrow \infty.$$

제안된 추정량의 점근분포를 알아보기 위하여,  $Z_n(t) \equiv \sqrt{n}(S(t) - \hat{S}(t)) / S(t)$ 라 정의하면 확률과정  $Z_n(t)$ 는 보조정리 2.2에 의해서, 확률과정  $S(t)Z_n(t)$ 는 평균이 0이고 공분산함수가  $Cov(Z^*(s), Z^*(t))$ 인 가우시안 확률과정으로 수렴하기 때문에,  $[0, T], H(T) < 1$  상에서 브라운니안 모션 확률과정(Brownian motion process),  $Z \equiv B(V)$ 로 수렴하고 여기서  $V(t) = \int_0^t \frac{dF(s)}{(1-F(s))S(s)}$ 이다. 이 결과와 Gill(1983)의 정리2.2를 이용하여 제안된 추정량의 점근분포를 아래의 정리 2.2에 밝혔다.

정리 2.2  $h(t) = \int_t^{\tau_F} S(s)ds$ 라두고, 만약

$$(1) \sqrt{n} h(T^*) \rightarrow^p 0.$$

$$(2) \int_0^{\tau_H} h^2(s) dV(s) < \infty$$

가 성립한다면,  $0 \leq t \leq T^*$ 에 대하여,  $D[0, \tau_F]$  상에서

$$\sqrt{n}(\hat{e}_t - e_t^{\tau_F}) \rightarrow^D B^*(t) \equiv - \frac{\int_t^{\tau_F} h(s) dZ(s)}{S(t)}, \quad n \rightarrow \infty$$

여기서  $B^*$ 는 평균이 0이고 공분산함수가

$$Cov(B^*(s), B^*(t)) = (S(s)S(t))^{-1} \int_{s \wedge t}^{\tau_F} h^2(u) dV(u).$$

인 가우시안 확률과정이다

증명 함수  $h(t)$ 와 조건 (2)에 의해 Gill(1983)의 정리 2.2의 조건을 만족해주고, Fleming과 Harrington(1991)의 부록에 있는 정리A.1.2를 이용하면,

$$\begin{aligned} \sqrt{n}(\hat{e}_t - e_t^{\tau_F}) &= (S(t)S(t))^{-1} \{ -S(t) \int_t^{T^*} h(s) dZ_n(s) \\ &\quad + S(t)Z_n(T^*)h(T^*) - S(t)\sqrt{nh(T^*)} \} \end{aligned}$$

이다. 그리고 조건 (1)에 의해 위의 정리가 성립함을 보였다.■

위의 정리의 결과를 이용하여  $e_t^{\tau_F}$ 에 대한 점근적 신뢰구간을 구해보자. 먼저 정리2.2



$\hat{e}_{t,NA}$ 의 편의가 가장작고 그 외의 부분에서는  $\hat{e}_{t,KM}$ 의 편의가 더 적었다.

평균제곱오차의 측면에서 볼 때,  $\hat{e}_{t,SV}$ 가 가장 우수하였고  $\hat{e}_t$ 는  $\hat{e}_{t,KM}$ ,  $\hat{e}_{t,NA}$ 보다 더 우수하였다.

(2)수명함수  $F$ 의 분포가 Weibull(1.0, 1.0)인 경우, 평균제곱오차의 측면에서  $\hat{e}_{t,SV}$ 와  $\hat{e}_t$ 는  $\hat{e}_{t,KM}$ ,  $\hat{e}_{t,NA}$ 보다 더욱 우수하였고, 특히, 중단비율이 10%인 경우 표본의 크기에 상관없이  $\hat{e}_t$ 는  $\hat{e}_{t,SV}$ 보다 수명분포의 끝부분에서 더 작은 평균제곱오차를 가졌다.

편의의 측면에서볼때  $\hat{e}_{t,KM}$ ,  $\hat{e}_{t,NA}$ 이  $\hat{e}_{t,SV}$ 와  $\hat{e}_t$ 보다 더 작은 편의를 가졌다.  $\hat{e}_{t,KM}$ 은  $\hat{e}_{t,NA}$ 보다 작은 편의를 가졌다.  $\hat{e}_{t,SV}$ 와  $\hat{e}_t$ 를 비교해볼때 중단 비율이 30%인 경우,  $\hat{e}_t$ 의 편의가  $\hat{e}_{t,SV}$ 보다 작고 중단비율이 10%인 경우는  $\hat{e}_{t,SV}$ 의 편의가  $\hat{e}_t$ 보다 작았다.

(3)수명함수  $F$ 의 분포가 Weibull(2.0, 1.5)일 때, 표본의 크기와 중단비율에 관계없이  $\hat{e}_t$ 와  $\hat{e}_{t,SV}$ 의 평균제곱오차는  $\hat{e}_{t,KM}$ 나  $\hat{e}_{t,NA}$ 보다 작았다. 그러나 편의의 측면에서 볼 때  $\hat{e}_{t,KM}$ ,  $\hat{e}_{t,SV}$ 의 편의가  $\hat{e}_{t,NA}$ 와  $\hat{e}_t$ 보다 작고 특히  $\hat{e}_{t,KM}$ 의 편의가 가장 작은 것을 알 수있었다.

결론적으로 편의나 평균제곱오차의 측면에서는 중단비율이 높다고 생각되는 자료에서는  $\hat{e}_t$ 의 사용이 더 바람직하다고 할 수있다.







## 참 고 문 헌

- (1) Aalen, O. (1978), Nonparametric inference for a family of counting processes. *Annals of Statistics*, 6, 701-726.
- (2) Anderson, P.K. and Borgan,  $\phi$ . (1985), Counting process models for life history data : A review, *Scandinavian Journal of Statistics*, 12, 97-158.
- (3) Br srow, N.E. and Crowley, J. (1974), A large sample study of the life table and product limit estimators under random censorship. *Annals of Statistics*, 2, 435-453.
- (4) Fleming, T.R. and Harrington, D.P. (1991), *Counting Process and Survival Analysis*, John Wiley & Sons Ins., New York.
- (5) Ghorai, J.K., Susarla, A., Susarla, V. and Van Ryzin, J.(1980), Nonparametric estimation of mean residual life time with censored data. *Colloquia Mathematica Societatis Janos Bolyai, Nonparametric Statistical Inference*, 32, 269-291.
- (6) Gill, R.D. (1983), Large sample behavior of the product-limit estimator on the whole line. *Annals of Statistics*, 11, 49-58.
- (7) Hjort, N.L. (1991), Nonparametric Bayes estimators based on beta processes in models for life history data. *Annals of Statistics*, 3, 1259-1294.
- (8) Kaplan, E.L. and Meier, P.(1958), Nonparametric estimation from incomplete observations, *Journal of the American Statistical Association*, 53, 457-481.
- (9) Kumazawa, Y.(1987), A note on an estimator of life expectancy with random censorship. *Biometrika*, 74, 655-658.
- (10) Lee I.S. and Lee, W.D. (1993), Nonparametric Bayes estimation of survival function, Unpublished ph. D. thesis.
- (11) Nelson ,W.B.(1972), Theory and applications of hazard plotting for censored failure data. *Technometrics*, 14, 945-996.
- (12) Peterson, A.V. (1977), Expressing the Kaplan-Meier estimator as a function of empirical survival functions. *Journal of the American Statistical Association*, 72, 854-858.
- (13) Susarla, V. and Van Ryzin, J. (1980), Large sample theory for an estimator of the mean survival time from censored samples. *Annals of Statistics*, 8, 1002-1016.
- (14) Yang, G.L. (1977), Life expectancy under random censorship. *Stochastic process and their applications*, 6, 33-39.
- (15) Yang, G.L. (1978), Estimation of biometric function. *Annals of Statistics*, 6, 112-116.