

# A two-sample test with interval censored competing risk data using multiple imputation

Yuwon Kim<sup>a</sup> · Yang-Jin Kim<sup>a,1</sup>

<sup>a</sup>Department of Statistics, Sookmyung Women's University

(Received January 10, 2017; Revised February 24, 2017; Accepted March 21, 2017)

---

## Abstract

Interval censored data frequently occur in observation studies where the subject is followed periodically. In this paper, our interest is to suggest a test statistic to compare the CIF of two groups with interval censored failure time data in the presence of competing risks. Gray (1988) suggested a test statistic for right censored data that motivated a well-known Fine and Gray's subdistribution hazard model. A multiple imputation technique is adopted to adopt Gray's test statistic to interval censored data. The powers and sizes of the suggested method are investigated through diverse simulation schemes. The main merit of the suggested method is its simplicity to implement with existing software for right censored data. The method is illustrated by analyzing Bangkok's HIV cohort dataset.

Keywords: competing risk, Gray test, interval censored data, multiple imputation

---

## 1. 서론

경쟁 위험 자료는 여러 가지 다른 원인으로 인한 사망 또는 기계 고장이 발생할 경우에 고려되는 생존 자료의 한 유형으로 관련된 예는 의약학, 경제학, 사회학, 보건학 등에서 두 가지 이상의 잠재적인 사건을 경험할 수 있는 경우에 적용될 수 있다. 일반적인 생존 분석 자료가 이항 자료(사건 발생 여부)와 사건 발생 시간과의 결합 자료( $T, \epsilon = I(T < C)$ )라면 경쟁 위험 자료는 다항 자료(원인 1, 원인 2, ..., 원인  $K$ )와 사건 발생 시간과의 결합 자료( $T, \epsilon \times \delta$ )로 구성된다. 여기서  $\delta$ 는 사건 발생 원인 지시함수를 의미한다. 즉, 모든 잠재 사건 변수를  $T_1, \dots, T_K$ 라 할 때, 실제 관측 변수는 그들 중 가장 먼저 일어난 사건의 발생 시간,  $T = \min(T_1, T_2, \dots, T_K)$ 이며  $\delta = (k; T = T_k)$ 로 정의된다.

경쟁 위험 자료 분석에서는 일반적으로 다음 두 가지 통계량, 원인별 위험 함수(cause specific hazard function; CSH),  $\lambda_k(t|z) = \lim_{\Delta \rightarrow 0} (\Pr(t < T < t + \Delta, \delta = k | T \geq t, z)) / \Delta$ 와 누적 분포 함수(cumulative incidence function; CIF),  $F_k(t) = \Pr(T \leq t, \delta = k)$ 가 주로 사용된다. 여기서  $\lim_{t \rightarrow \infty} F_k(t) = \Pr(\delta = k)$ 이며 중도 절단 자료가 없을 경우,  $\sum_{k=1}^K \Pr(\delta = k) = 1$ 의 관계가 성립함으로  $F_k(t)$ 를 종종 subdistribution 함수라 한다. 원인별 위험 함수는 모든 원인이 가능하다는 조건하에서 주어진 시점에서 관심 있는 원인이 발생할 순간 위험률로 해석되며 모수 추정 시 다른 원인에 의한 사건

---

This work was supported by Korea research grant NRF-2014R1A2A2A01003567.

<sup>1</sup>Corresponding author: Department of Statistics, Sookmyung Women's University, 100, Cheongpa-ro 47-gil, Yongsan-gu, Seoul 04310, Korea. E-mail: [yjin@sookmyung.ac.kr](mailto:yjin@sookmyung.ac.kr)

을 중도 절단으로 간주하여 분석하게 된다. 따라서 다른 사건의 발생 여부에 대한 정보를 전혀 고려하지 않음으로써 발생 확률에 대한 추론에 편의를 가져오게 된다. 한편, CIF는 공변량  $Z = z$ 가 주어져 있을 때, 식 (1.1)과 같이 모든 원인별 위험 함수의 정보를 반영하게 된다.

$$F_k(t|z) = \int_0^t \lambda_k(u|z)S(u|z)du = \int_0^t \lambda_k(u|z)e^{-\int_0^u \sum_{l=1}^K \lambda_l(s|z)ds} du. \quad (1.1)$$

따라서  $1 - \exp(-\int_0^t \lambda_k(u|z)du) \neq F_k(t|z)$ 임을 알 수 있다. 이에 CIF와 다음 식 (1.2)의 관계를 가지는 subdistribution 위험 함수  $h_k(t|z)$ 를 정의한다.

$$h_k(t|z) = -\frac{d \log(1 - F_k(t|z))}{dt}. \quad (1.2)$$

Fine과 Gray (1999)는 공변량의 효과를 추정하기 위해  $h_k(t|z)$ 에 비례 위험 모형을 적용한 후, 회귀 계수 추정을 위해 inverse probability censoring weights (IPCW) 기법을 통해 새로운 위험 그룹(risk set)을 정의하였다. 이에 앞서 Gray (1988)는 두 그룹의 CIF를 비교하기 위해 검정 통계량을 제안하는 과정에서 가중화 위험그룹과 subdistribution 위험 함수개념을 소개하였다. 여기서 주 관심은 경쟁 사건을 경험한 대상에 대한 처리이다. 즉, 그들의 위험 그룹 포함 여부 정도는 사건 발생 시간과 우중도 절단 자료의 분포에 의존한다. 이에 근거한 다양한 분석 방법들이 우중도 절단 자료에 대해 제안된 반면에 경쟁 위험을 가진 구간 중도 절단 자료 분석에 대해선 그 연구 범위가 아직까지 넓지 못하다. 특히 대부분의 연구는 제 1유형의 구간 중도 절단 또는 current status 자료에 대한 CIF 추정량들에 관한 것이다 (Jewell과 Kalbfleisch, 2004; Jewell 등, 2003). 본 연구의 목적은 구간 중도 절단 자료에 대한 이표본 검정 통계량을 제안하고자 한다. 이를 위해 다중 대체 방법(multiple imputation method)을 적용할 것이다. 다중 대체 방법은 경쟁 위험 모형과 구간 중도 자료에 적용되어왔다. 경쟁 위험 자료에 대한 다중 대체 방법의 적용은 Goetghebeur와 Ryan (1995) 또는 Ruan과 Gray (2008)에 의해 각각 subdistribution 위험 함수와 원인별 위험 함수에서 적용되었다. 특히, Ruan과 Gray (2008)의 논문에서는 경쟁 위험 사건을 경험함으로써 알 수 없는 중도 절단 시점을 결측 자료(missing data)로 간주한 후 이 문제를 위해 다중 대체 방법을 적용하여 완전 자료로 생성하고자 하였다. Lu와 Tsiatis (2001)는 발생 원인이 결측 자료일 경우, 발생 원인을 추정하기 위해 다중 대체 방법을 적용했다. 구간 중도 절단 자료에 대한 다중 대체 방법의 예로 Pan (2000a)은 근사적 베이지안 방법을 이용해 완전 자료를 생성한 후, 검정 통계량을 적용하였으며 이를 확장하여 회귀 분석을 시행하였다 (Pan, 2000b).

본 논문의 2장에서는 다중 대체 방법이 간략하게 요약된 후 구간 중도 절단 자료에서 두 그룹의 CIF를 비교하기 위해 이 방법의 적용과정이 기술될 것이다. 3장에서는 모의 실험을 통해 제안된 검정통계량의 검정력과 유의 수준을 추정하며 HIV 환자로부터 수집된 구간 중도 절단 자료에서 남녀간의 HIV 누적 발생 함수를 비교하기 위해 제안된 검정 통계량을 적용하였다. 마지막으로 4장에서는 관련된 연구를 요약함으로써 논문을 마무리하고자 한다.

## 2. Two sample test using multiple imputation

먼저 본 논문에서 사용할 기호를 정의하기로 하자.  $\tilde{T}_{jl}$ 은  $j(= 1, 2)$ 번째 그룹의  $l(= 1, \dots, n_j)$ 번째 개체의 생존 시간을 의미하며  $\delta_{jl}$ 은 그 개체가 경험하는 사건 발생 원인 또는 사건 유형을 의미한다. 따라서  $F_{jk}(t) = \Pr(T_{jl} \leq t, \delta_{jl} = k)$ 는  $j$ 번째 그룹의  $k$ 번째 원인에 의한 CIF가 된다. 본 연구의 주요 목적은 다음의 귀무가설을 검정하기 위한 검정 통계량을 제안하고자 한다.

$$H_0 : F_{11} = F_{21} = F_{01}.$$

즉, 두 그룹의 원인 1에 의한 CIF를 비교하고자 한다. 우중도 절단 시간  $C_{jl}$ 과 함께 실제 관측 자료는  $(T_{jl}, \delta_{jl})$ 이며  $T_{jl} = \min(\tilde{T}_{jl}, C_{jl})$ 이 된다.

검정 통계량을 유도하기 위해 두 가지 중요한 통계량이 다음과 같이 정의된다.  $d_{jk}(t) = \sum_{l=1}^{n_j} I(\tilde{T}_{jl} = t, \delta_{jl} = k)$ 를  $t$ 시점에서 그룹  $j$ 에서 일어난 원인  $k$ 에 의한 사건 발생 수를 의미하며  $n_j(t) = I(\tilde{T}_{jl} \geq t) + I(\tilde{T}_{jl} < t, \delta_{jl} \neq 1)$ 는  $t$ 시점에서 그룹  $j$ 에 남아 있는 위험 그룹의 크기를 의미한다. Gray (1988)는 다음과 같이 정의된 가중화된 위험 그룹 함수를 제안하였다.

$$r_j(t) = n_j(t) \frac{1 - \hat{F}_{j1}(t)}{\hat{S}_j(t)}, \quad j = 1, 2,$$

여기서  $\hat{F}_{j1}$ 은 그룹  $j$ 에서 발생한 원인 1의 CIF 추정치이며,  $\hat{S}_j$ 은 그룹  $j$ 에서 발생한 모든 사건을 이용하여 추정된 생존 함수의 Kaplan-Meier 추정량이 된다. 여기서 모든 사건 발생 시간을 순서화한  $t_1 < \dots < t_q$ 를 이용하여 Gray의 검정 통계량을 다음과 같이 정의된다.

$$\hat{G} = \sum_{l=1}^q K(t_l) \left[ \frac{d_{11}(t_l)}{r_1(t_l)} - \frac{d_{\cdot 1}(t_l)}{r_{\cdot}(t_l)} \right] = \sum_{l=1}^q K(t_l) \left[ \hat{h}_{11}(t_l) - \hat{h}_{\cdot 1}(t_l) \right],$$

여기서  $d_{\cdot 1}(t_l) = d_{11}(t_l) + d_{21}(t_l)$ 로  $t_l$ 시점에 두 그룹에서 발생한 원인 1의 사건 합을 의미하며  $r_{\cdot}(t) = r_1(t) + r_2(t)$ 는 두 그룹에 속한 위험 함수의 합이다. 여기서  $\hat{h}_{j1}$ 은 그룹  $j$ 의 원인 1의 subdistribution 위험함수의 추정치를 의미한다. 가중치 함수  $K(t)$ 는 검정력을 높이기 위해 적절하게 선택된 함수로 예를 들어  $K(t) = L(t)r_1(t)$ 에서  $L(t) = [\hat{S}(t)]^\rho$ ,  $\hat{S}(t) = 1 - \hat{F}_{01}(t)$ 가 사용된다. 이 때,  $\rho > 0$ 는 전반부에 두 함수가 차이가 클 때,  $\rho < 0$ 는 후반부에 두 함수가 차이가 큰 경우에 각각 적용됨으로써 더 큰 검정력을 가져오게 된다. 여기서  $\hat{F}_{01}$ 은 쿠무가설하에서 구한 원인 1의 CIF 추정량이 된다.

본 논문에서 분석할 구간 중도 절단자료는  $(V_{jl}, U_{jl}, \delta_{jl}, j = 1, 2; l = 1, \dots, n_j)$ 로 여기서  $0 \leq V_{jl} \leq \tilde{T}_{jl} \leq U_{jl} < \infty$ 의 관계를 가지며, 만약  $U_{jl} = \infty$ 이라면  $\tilde{T}_{jl}$ 는 우중도 절단됨을 의미한다. 중도 절단 자료를 구성하는  $(V_{jl}, U_{jl})$ 은 사건 발생 시간  $\tilde{T}_{jl}$ 와 독립임을 가정한다. 경쟁 위험이 존재하지 않은 구간 중도 절단 자료에 대해 Pan (2000b)은  $H_0: F_1 = F_2$ 를 검정하기 위해 다중 대체 방법을 적용하였다. 여기서  $F_j$ 는 그룹  $j(= 1, 2)$ 의 분포함수로 다중 대체의 첫 번째 단계인 자료 강화(data augmentation)에서는 구간 중도 자료를 대체할 자료를 생성하게 된다. 즉, 구간 중도 자료를 결측 자료로 간주한 후 각 그룹에서 추정된 생존 함수를 이용하여 구간 중도 절단 자료를 대신할 자료를 생성하게 된다. 이렇게 생성된 자료는 더 이상 구간 중도 절단 자료가 아닌 우중도 절단 자료를 포함한 일반적인 생존 자료를 가져오게 된다. 두 번째 단계에서는 이러한 자료를 분석하기 위해 로그 순위 검정통계량과 같은 우중도 생존 자료 분석 방법을 적용한다. 위 두 단계를  $M$ 번 반복함으로써 구한  $M$ 개의 결과치의 평균을 이용하여 검정 통계량을 계산하게 된다. 여기서 표본들 간의 내외 변이를 반영하기 위해 새로운 유형의 분산을 이용하게 된다. 본 연구에서는 Pan (2000b)의 방법을 경쟁 위험 자료로 확장하고자 한다.

먼저 구간 중도 절단자료를 대체할 사건 발생 시간을 생성하기 위해 추정된 CIF를 사용한다. 여기서 구간 중도 절단된 자료에 대해 CIF를 추정하기 위해 Hudgens 등 (2001)의 방법을 적용한다. 그들의 논문에서는 Turnbull (1976)이 제안한 자기 일치 추정량을 확장함으로써 EM 알고리즘의 적용과정을 보였다. 구간 중도 절단 자료에 대한 CIF를 추정하기 위해 먼저 다음의 우도 함수가 적용되었다. 그룹  $j$ 에 대해

$$L^j = \prod_{l=1}^{n_j} [F_{j\delta_{jl}}(u_{jl}) - F_{j\delta_{jl}}(v_{jl})]^{I(\delta_{jl} > 0)} \left[ \sum_{k=1}^K F_{jk}(u_{jl}) - F_{jk}(v_{jl}) \right]^{I(\delta_{jl} = 0)}.$$

위 우도 함수에 근거하여 추정된 nonparametric maximum likelihood estimator (NPMLE)는 원인 별 시간대(time support)에서 정의되며 이는 추정된 생존 함수  $\widehat{S}_j(t) = 1 - \sum_k^K \widehat{F}_{jk}$ 를 적용할 경우 정의되지 못하는 시간대가 존재하게 된다 (Hudgens 등, 2001). 이러한 문제점을 해결하기 위해 모든 원인 별 위험 함수가 공통으로 공유할 수 있는 공통된 시간대(pooled time support)를 이용하여 다음의 pseudolikelihood를 제안하였다.

$$PL^j = \prod_{l=1}^{n_j} [F_{j\delta_{jl}}(u_{jl}) - F_{j\delta_{jl}}(v_{jl})]^{I(\delta_{jl}>0)} [S_j(v_{jl}) - S_j(u_{jl})]^{I(\delta_{jl}=0)}. \quad (2.1)$$

따라서  $F_{jk}$ 를 추정하기 위해 먼저 공통된 시간대를 정의할 필요가 있으며 이를 위해 다음의 equivalence set  $\{Q_r^j = (q_r^j, p_r^j), r = 1, \dots, R_j\}$ 을 정의한다 (Lindsey와 Ryan, 1998). 지시 함수  $\alpha_{ikr}^j = I([q_r^j, p_r^j] \in [v_{jl}, u_{jl}] \cap \delta_{jl} = k)$ 를 정의한 후  $\psi_{kr}^j = F_{jk}(p_r^j+) - F_{jk}(q_r^j-)$ 은  $j$ 번째 그룹의  $r$ 번째 시간대에서  $k$ 번째 원인의 발생 확률을 의미한다. 이 두 함수를 이용하여 식 (2.1)은 다음과 같이 다시 정의된다.

$$\prod_{i=1}^{n_j} \sum_{k=1}^K \sum_{r=1}^{R_j} [\alpha_{ikr}^j \psi_{kr}^j]^{I(\delta_{ji}>0)}.$$

이 때 모수 추정을 위해 EM 알고리즘이 다음과 같이 적용된다.  $I_{ikr}^j = I(T_{jl} \in [q_r^j, p_r^j] \cap \delta_{jl} = k)$ 은 그룹  $j$ 의  $l$ 번째 환자의  $k$ 번째 사건 발생 시간  $T_{jl}$ 이  $[q_r^j, p_r^j]$ 에 포함된 여부를 보여주는 지시함수이다. 하지만  $T_{jl}$ 을 알지 못하기 때문에 미지(unknown) 함수이며 이의 기대값(expected value)을 다음과 같이 구할 수 있다.

$$E(I_{ikr}^j | \psi) = \mu_{ikr}^j(\psi) = \frac{\alpha_{ikr}^j \psi_{kr}^j}{\sum_{r'=1}^{R_j} \sum_{k'=1}^K \alpha_{ik'r'}^j \psi_{k'r'}^j}. \quad (2.2)$$

즉,  $\psi_{kr}^j = F_{jk}(q_r^j) - F_{jk}(p_{r-1}^j)$ 로  $[q_r^j, p_r^j]$ 에서 원인  $k$  사건의 발생 확률을 의미한다. 식 (2.2)를 구하기 위해 먼저 적절한 초기값  $\widehat{\psi}_{kr}^{j,0} = 1/K \times R_j$ 을 대입한 후, 다음의 실패 발생함수를 추정한다.

$$\psi_{kr}^j(\psi^j) = \sum_{i=1}^{n_j} \widetilde{\psi}_r^j \frac{\mu_{ikr}^j}{\sum_{r'} \mu_{ikr'}^j}, \quad \widetilde{\psi}_r^j = \sum_{k=1}^K \psi_{kr}^j. \quad (2.3)$$

적절한 수렴기준을 만족할 때까지, E-step (2.2)와 M-step (2.3)을 반복한다. 마지막 추정된  $\widehat{\psi}_{kr}^j$  ( $r = 1, \dots, R_j; k = 1, \dots, K$ )를 이용하여  $p_l^j < t < q_{l+1}^j$ 에 대해서는  $\widehat{F}_{jk}(t) = \sum_{r=1}^l \widehat{\psi}_{kr}^j$ 이 되며  $t > p_{R_j}^j$ 일 때는  $\widehat{F}_{jk}(t) = \widehat{\psi}_{k1}^j + \dots + \widehat{\psi}_{kR_j}^j$ 가 된다.

이제 다중 대체방법은 다음의 과정을 통해 적용된다. 여기서  $M$ 은 다중 대체의 반복 수를 의미한다.

1)  $m = 1, \dots, M$ 일 때, 다음의 과정을  $M$ 번 반복한다.

(a) 우중도 절단 자료인 경우에  $u_{jl} = \infty$ ,  $T_{jl}^{(m)} = v_{jl}$ 이며  $\delta_{jl}^{(m)} = 0$ 이 된다.

(b) 구간 중도 절단인 경우에  $(v_{jl}, u_{jl})$ , 추정된 CIF,  $\widehat{F}_{jk}$ 로부터  $\{v_{jl} < T_{jl}^{(m)} < u_{jl}\}$ 을 만족하는 조건하에서  $T_{jl}^{(m)}$ 을 생성한다. 여기서  $\delta_{jl}^{(m)} = \delta_{jl}$ 이 된다.

2)  $\{(T_{11}^{(m)}, \delta_{11}^{(m)}), \dots, (T_{1n_1}^{(m)}, \delta_{1n_1}^{(m)})\}, \{(T_{21}^{(m)}, \delta_{21}^{(m)}), \dots, (T_{2n_2}^{(m)}, \delta_{2n_2}^{(m)})\}$ 에 대해  $\widehat{G}^{(m)}$ 을 계산하며 해당 분산  $\widehat{\Sigma}^{(m)}$ 을 구한다.

3) 위 두 통계량을 이용하여

$$\hat{G} = \frac{1}{M} \sum_{m=1}^M \hat{G}^{(m)}, \quad \hat{\Sigma} = \frac{1}{M} \sum_{m=1}^M \hat{\Sigma}^{(m)} + \left(1 + \frac{1}{M}\right) \text{Var} \left( \hat{G}^{(1)}, \dots, \hat{G}^{(M)} \right).$$

4) 귀무 분포  $N(0, \hat{\Sigma})$ 를 이용하여  $p$ -value를 계산한다.

여기서 주의 사항은 위 알고리즘을 시행 시 CIF을 이용할 경우  $\hat{F}_{jk}(v_{jl}) = \hat{F}_{jk}(u_{jl})$ 을 가진 구간 중도 절단 자료  $(v_{jl}, u_{jl})$ 가 발생할 수 있다. 이러한 경우  $T_{jl}^{(m)}$ 를 생성하기 위해 본 논문에서는 선형 스무딩 방법이 적용되었다.

### 3. 모의실험

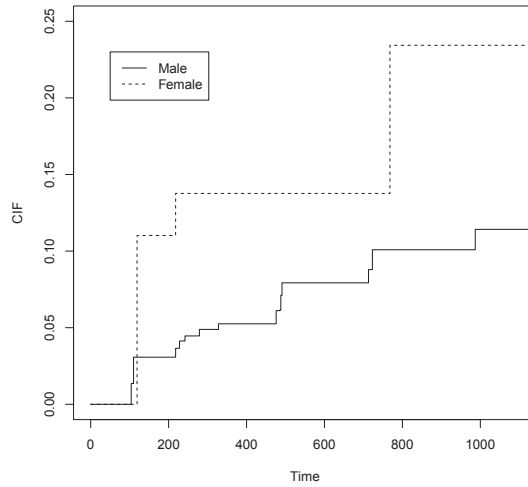
제안한 방법을 평가하기 위해 다양한 모의실험을 시행하였다. 발생 원인이 두 가지( $K = 2$ )인 경쟁 위험모형을 가정하에서 중도절단 비율을 10%와 25%로 설정하고, 각 그룹의 표본의 크기는  $n_j = 100, 200, 300$ 인 자료를 생성하였다. 모의실험에서 알고리즘 반복수는 1,000번에 다중대체 수는  $M = 5$ 를 정하였다. 본 연구에서는 다섯 가지 시나리오를 고려하였다. 여기서 그룹 1에 대한 원인별 위험률을  $(\lambda_{11}, \lambda_{12}) = (0.3, 0.2)$ 로 고정하였고 그룹 2에 대한 원인별 위험률을 다섯 가지 경우로 설정하였다. Figure 3.1은 대립가설 하에서 두 그룹의 원인 1의 CIF를 나타낸 그래프이다. 자료를 생성하는 단계는 다음과 같다 (Beyersmann 등, 2009).

- (1) 각 그룹별 ( $j = 1, 2$ )로  $\lambda_j = \lambda_{j1} + \lambda_{j2}$ 를 이용하여 지수분포로부터 사건 발생 시간  $T_{jl}$ 을 생성한다.
- (2) (1)에서 생성한 발생시간의 원인을 결정하기 위해, 원인 1의 확률( $\lambda_{j1}/(\lambda_{j1} + \lambda_{j2})$ )을 가진 베르누이 확률 변수를 생성한다. 확률 변수가 1의 값을 가질 땐 원인 1에 의해 사건이 발생한다고 간주하고, 그렇지 않으면 원인 2에 의해 사건 ( $\delta_{jl} = 2$ )이 발생했다고 간주한다.
- (3) 중도 절단률을 설정하기 위해 우 중도절단 발생 시간  $C_{jl}$ 를 지수분포로부터 생성한다.  $T_{jl} < C_{jl}$ 인 경우에 다음 단계에서 구간 중도 절단 자료를 생성하고  $T_{jl} > C_{jl}$ 이면 우중도 절단으로 간주한다.
- (4) 이산형 균일분포 로부터 관측 횟수  $R \sim U(10, 15)$ 를 결정한다. 또한 관측시간은  $w_r \sim U(0, 10)$ 로부터 생성하여 크기 순서로  $w_1 < \dots < w_R$ 로 정돈한 후,  $w_{l-1} = v_{jl} < T_{jl} < u_{jl} = w_l$ 을 만족하는 값을  $(v_{jl}, u_{jl})$ 으로 정의한다.

Table 3.1은 다섯 가지 시나리오에 대하여 표본크기, 중도절단 비율, 가중치에 따라 유의수준과 검정력을 정리한 것이다. 귀무가설 하에서 귀무가설을 기각할 확률이 거의 0.05에 가까운 것을 확인할 수 있다. 그룹 2가  $(\lambda_{21}, \lambda_{22}) = (0.3, 0.1)$  또는  $(0.3, 0.5)$ 인 경우에 가중치  $\rho = -1$ 가 가장 큰 검정력을 가졌다. 이는 두 그룹의 차이가 후반부에 더 크기 때문에 나타난 결과이며 Figure 3.1의 위 두 그림을 통해 확인할 수 있다. 반면에  $(\lambda_{21}, \lambda_{22}) = (0.5, 0.2)$  또는  $(0.8, 0.2)$  경우에는 두 그룹의 CIF가 초기에 더 큰 차이를 보이기 때문에  $\rho = 1$ 에서 더 검정력을 가지게 된다. 전체적으로 중도절단비율이 10% 경우보다 25% 경우에 검정력이 떨어졌다.

실제 자료의 적용 예로 1,209명의 HIV 혈청 음성 injecting drug users (IDU)의 자료가 분석된다 (Hudgens 등, 2001). 참여자들은 4개월 마다 병원을 방문하여 혈청변환을 검사하도록 하였으며 그 중 1,124명이 한 번 이상 방문하였고, 그 중 133명이 HIV에 감염된 것으로 확인되었다. 133명의 환자 중 B type의 환자가 27명, E type의 환자가 99명이었으며 나머지 7명의 환자는 type이 밝혀지지 않았다.





**Figure 3.2.** Comparison of CIFs by gender. CIF = cumulative incidence function.

Type이 밝혀지지 않은 7명의 환자는 제외한 1,117명을 대상으로 분석을 시행하였다. 본 연구에서는 E type이 주 관심으로 정하고 성별에 따른 E type의 CIF를 비교하고자 한다. 다중 대체 횟수를  $M = 5$ 와 7로 달리 적용하였으며  $\rho = 0$ 을 적용하였을 때,  $p$ -value는 각각 0.0210, 0.0212로 모두 0.05보다 작았다. 이는 여자와 남자의 E type의 발생 분포가 차이가 있음을 나타낸다. Figure 3.2는 성별에 따른 E type의 추정된 CIF를 나타낸 그래프이다. 이 그림을 통해 여성이 남성보다 E type에 감염될 확률이 더 높은 것을 확인할 수 있다. 또한 다른  $\rho$ 값들을 적용할 경우 구해진  $p$ -value는 위의 결과와 비슷한 값을 보여주었다.

#### 4. 맺음말

본 논문에서는 구간 중도 경쟁 위험 자료를 가지는 두 그룹의 누적 발생 함수를 비교하기 위해 다중 대체 방법을 적용하였다. 제안한 방법은 기존에 개발된 통계 패키지를 이용할 수 있다는 점에서 그 실용성이 높다는 장점을 가진다. 물론 이는 사전에 각 원인별 누적 발생 함수를 추정해야 한다는 조건이 필요하다. 본 연구에서는 제안된 방법의 적절성을 확인하기 위해 여러 가지 상황에서 모의실험을 실시하였다. 귀무가설 하에서는 유의 수준에 근접한 결과를 가져왔으며 여러 가지 대립가설 하에서 여러 가중치를 적용함으로써 CIF 차이의 특성을 잘 반영함으로써 적합한 검정력을 보였다.

본 논문에서 제시된 연구는 여러 가지 방향으로 확장될 수 있다. 경쟁 위험 모형에서는 적절한 가정을 통해 경쟁 사건에 대해 고려할 필요가 있다. Fine와 Gray (1999)는 두 가지 경우를 고려하였다. 예를 들어, administrative censoring과 같이 미리 우중도 절단 시점이 정해진 경우 이 시점을 경쟁 사건의 우중도 절단 시점으로 사용하여 complete censoring data를 가지게 된다. 달리, 랜덤 우중도 절단 시점에 대해서는 우중도 절단 생존 함수를 추정한 후 이를 이용하여 가중치를 구한 후, IPCW방법을 적용할 수 있다. 본 연구의 확장 연구로 IPCW 방법을 구간 중도 절단 경쟁 위험 자료에 적용해보고자 한다. 두 번째 연구는 원인 결측 자료에 다중 대체 방법의 확장하고자 한다 (Bakoyannis 등, 2010). 이 때 우리는 두가지 유형의 결측자료 즉, 구간 중도된 사건 발생 시간과 원인 결측을 동시에 고려할 필요가 있다. 비슷한 연구로 Do와 Kim (2017)은 Klein 과 Andersen (2005)이 제안한 pseudo-value 방법을 구간 중도 절단 자료에 적용하였다.

## References

- Bakoyannis, G., Siannis, F., and Touloumi, G. (2010). Modelling competing risks data with missing cause of failure, *Statistics in Medicine*, **29**, 3172–3185.
- Beyersmann, J., Latouche, A., Buchhols, A. and Schumacher, M. (2009). Simulating competing risk data in survival analysis, *Statistics in Medicine*, **28**, 956–971.
- Do, G. and Kim, Y.-J. (2017). Analysis of interval censored competing risk data with missing causes of failure using pseudo values approach, *Journal of Statistical Computations and Simulations*, In press
- Fine, J. P. and Gray, R. J. (1999). A proportional hazards model for the subdistribution of a competing risk, *Journal of the American Statistical Association*, **94**, 496–509.
- Goetghebeur, E. and Ryan, L. (1995). Analysis of competing risks survival data when some failure types are missing, *Biometrika*, **82**, 821–834.
- Gray, R. J. (1988). A class of k-sample tests for comparing the cumulative incidence of a competing risk, *Annals of Statistics*, **16**, 1141–1154.
- Hudgens, M. G., Satten, G. A., and Longini, I. M. (2001). Nonparametric maximum likelihood estimation for competing risks survival data subject to interval censoring and truncation, *Biometrics*, **57**, 74–80.
- Jewell, N. P. and Kalbfleisch, J. D. (2004). Maximum likelihood estimation of ordered multinomial parameters, *Biostatistics*, **5**, 291–306.
- Jewell, N. P., Van der Laan, M. J., and Henneman, T. (2003). Nonparametric estimation from current status data with competing risks, *Biometrika*, **90**, 183–197.
- Klein, J. P. and Andersen, P. K. (2005). Regression modeling of competing risks data based on pseudo values of the cumulative incidence function, *Biometrics*, **61**, 223–229.
- Lindsey, J. and Ryan, L. (1998). Methods for interval censored data, tutorial in biostatistics, *Statistics in Medicine*, **17**, 219–138.
- Lu, K. and Tsiatis, A. A. (2001). Multiple imputation methods for estimating regression coefficients in the competing risks model with missing cause of failure, *Biometrics*, **57**, 1191–1197.
- Pan, W. (2000a). A multiple imputation approach to cox regression with interval censored data, *Biometrics*, **56**, 199–203.
- Pan, W. (2000b). A two-sample test with interval censored data with multiple imputation, *Statistics in Medicine*, **19**, 1–11.
- Ruan, P. K. and Gray, R. J. (2008). Analyses of cumulative incidence function via non-parametric multiple imputation, *Statistics in Medicine*, **27**, 5709–5724.
- Turnbull, B. W. (1976). The empirical distribution function with arbitrarily grouped censored and truncated data, *Journal of the Royal Statistical Society Series B (Methodological)*, **38**, 290–295.



# 다중대체방법을 이용한 구간 중도 경쟁 위험 모형에서의 이표본 검정

김유원<sup>a</sup> · 김양진<sup>a,1</sup>

<sup>a</sup>숙명여자대학교 통계학과

(2017년 1월 10일 접수, 2017년 2월 24일 수정, 2017년 3월 21일 채택)

---

## 요약

구간 중도 절단 자료는 관측 연구에서 종종 발생하는 생존 자료의 한 유형으로 관심 있는 사건 발생 시간을 정확하게 관측할 수 없는 대신에 이를 포함한 두 관측 시점으로 구성된다. 본 연구의 목적은 경쟁 위험이 구간 중도 절단 자료에서 발생할 경우, 두 그룹의 누적 발생 함수를 비교하기 위한 검정 통계량을 제시하는 것이다. 특히 본 연구에서는 다중 대체 방법을 통해 생성된 자료를 이용하여 검정력과 유의 수준을 구하고자 한다. 모의실험을 통해 제안한 방법이 다양한 경우에서 적절한 결과를 보이는지 검토하였으며 실제 자료 분석의 예로 남녀 그룹의 HIV 발생 함수의 차이를 비교하기 위해 제안한 방법을 적용하였다.

주요용어: 경쟁위험, 구간 중도 절단 자료, 다중 대체 방법, Gray test

---

본 연구는 한국 연구재단의 지원을 받아 수행한 연구임(NRF-2014R1A2A2A01003567).

<sup>1</sup>교신저자: (04310) 서울특별시 용산구 청파로47길 100, 숙명여자대학교 통계학과.

E-mail: yjin@sookmyung.ac.kr