

What to do today (2022/03/22)?

Part IV. Advanced Topics

- ▶ **Part IV.1 Counting Process Formulation** (Revisits to KM estm, Logrank test, and Cox PH model)
 - ▶ *Part IV.1.1 Theoretical Preparation*
- ▶ **Part IV.1.2 Counting Process Formulation in LIDA and Applications: Revisits to KM, Logrank, Cox PH**
 - ▶ *Part IV.1.2A Formulation*
 - ▶ *Part IV.1.2B Revisit to Logrank Test*
 - ▶ *Part IV.1.2C Revisit to Cox's Partial Likelihood Approach*
 - ▶ **Part IV.1.2D Revisit to Nelson-Aalen Estimator and Breslow Estimator**
 - ▶ **Part IV.1.2E Revisit to Kaplan-Meier Estimator**
 - ▶ **Part IV.1.2F Others**
- ▶ *Part IV.2 Selected Recent Topics in LIDA*
- ▶ *Part IV.3 Beyond Lifetime Data Analysis*

Part IV.1.2D Revisit to Nelson-Aalen Estimator and Breslow Estimator

Recall the notation:

$$N(t), Y(t), A(t), M(t) = N(t) - A(t)$$

Thus $dM(t) = dN(t) - Y(t)dH(t)$.

Assign $dM(t) = 0 \implies d\hat{H}(t) = \frac{dN(t)}{Y(t)}$:

Nelson-Aalen Estimator for the cumulative hazard function

$$\hat{H}(t) = \int_0^t \frac{1}{Y(u)} dN(u)$$

Part IV.1.2D Revisit to Nelson-Aalen Estimator and Breslow Estimator

Remarks:

- ▶ The observed failure times $0 \leq V_1 < \dots < V_J$,

$$\hat{H}(t) = \sum_{V_j \leq t} \frac{\# \text{ failures at } V_j}{\# \text{at risk at } V_j} = \sum_{V_j \leq t} \frac{d_j}{N_j}.$$

- ▶ Define $0/0 = 0$ or replace $1/Y_i(u) = I(Y_i(u) > 0)/Y_i(u)$.
- ▶ $E[dN_i(t)|Y_i(t)] = Y_i(t)h(t)dt$

Nelson-Aalen Estimator's Asymptotics

$[\hat{H}(t) - H(t)]$ is

$$\begin{aligned} & \int_0^t \frac{I(Y_{\cdot}(u) > 0)}{Y_{\cdot}(u)} dN_{\cdot}(u) - \int_0^t \frac{I(Y_{\cdot}(u) > 0) Y_{\cdot}(u)}{Y_{\cdot}(u)} dH(u) - \int_0^t I(Y_{\cdot}(u) = 0) dH(u) \\ &= \int_0^t \frac{I(Y_{\cdot}(u) > 0)}{Y_{\cdot}(u)} dM_{\cdot}(u) - \int_0^t I(Y_{\cdot}(u) = 0) dH(u) \end{aligned}$$

By Martingale CLT, in distn

$$\sqrt{n} \int_0^t \frac{I(Y_{\cdot}(u) > 0)}{Y_{\cdot}(u)} dM_{\cdot}(u) \rightarrow Gaussian(0, \sigma^2(t))$$

and $\sqrt{n} \int_0^t I(Y_{\cdot}(u) = 0) dH(u) \rightarrow 0$ in Pr 1.

Thus $\sqrt{n} [\hat{H}(t) - H(t)] \rightarrow Gaussian(0, \sigma^2(t))$.

Plus $\sup_{t>0} |\hat{H}(t) - H(t)| \rightarrow 0$ in Pr 1.

Remarks:

- ▶ $\sigma^2(t) = \int_0^t \frac{1}{P(U \geq u)} dH(u)$ and $\hat{\sigma}^2(t) = \int_0^t \frac{nI(Y_{\cdot}(u) > 0)}{Y_{\cdot}(u)} d\hat{H}_{NA}(u)$
with $n \gg 1$.
- ▶ **Fleming-Harrington Estimator** for survivor function
 $\hat{S}(t) = \exp\{-\hat{H}(t)\}$:

$$\sqrt{n}[\hat{S}(t) - S(t)] \approx S(t)\sqrt{n}[\hat{H}(t) - H(t)]$$

- ▶ **Breslow Estimator** for the baseline hazard function in the Cox PH model:

$$\hat{H}_0(t; \beta) - H_0(t) = \int_0^t \frac{1}{\sum_{l=1}^n Y_l(t) e^{\beta z_l}} dN_{\cdot}(u)$$

with β replaced by the PMLE $\hat{\beta}$

- ▶ Note that
 $\hat{H}_0(t; \hat{\beta}) = [\hat{H}_0(t; \hat{\beta}) - \hat{H}_0(t; \beta)] + [\hat{H}_0(t; \beta) - H_0(t)]$

Part IV.1.2D Revisit to Nelson-Aalen Estimator and Breslow Estimator: Confidence Band of $H(\cdot)$

Recall that $\hat{H}_{NA}(t) = \int_0^t \frac{I(Y_{\cdot}(u))}{Y_{\cdot}(u)} dN_{\cdot}(u) \sim Gaussian(H(t), \frac{\sigma^2(t)}{n})$ approximately.

Pointwise 95% CI: $\forall t > 0, \hat{H}_{NA}(t) \pm 1.96 \sqrt{\frac{\hat{\sigma}^2(t)}{n}}$

95% Confidence Band: the lower, upper boundaries $L(t), U(t)$ satisfy $P(L(t) \leq H(t) \leq U(t) : t \in (0, \infty)) = 95\%$.

e.g. $\hat{H}_{NA}(t) \pm c \sqrt{\frac{\hat{\sigma}^2(t)}{n}}$ with c determined by

$$P\left(\sup_{t>0} \left| \frac{\hat{H}_{NA}(t) - H(t)}{\sqrt{\hat{\sigma}^2(t)/n}} \right| \leq c\right) = 95\%.$$

How to compute the critical value c ?

► **Approach I** (based on the Brownian Motion)

Recall $\sqrt{n} \frac{\hat{H}_{NA}(\cdot) - H(\cdot)}{\hat{\sigma}(t)} \rightarrow W\left(\frac{\sigma^2(\cdot)}{\sigma^2(t)}\right)$ in distn: $W(\cdot)$ the standard Brownian motion.

Thus, by the continuous mapping theorem,

$$\sup_{0 \leq s \leq t} \sqrt{n} \left| \frac{\hat{H}_{NA}(\cdot) - H(\cdot)}{\hat{\sigma}(t)} \right| \rightarrow \sup_{0 \leq s \leq t} \left| W\left(\frac{\sigma^2(\cdot)}{\sigma^2(t)}\right) \right| = \sup_{0 \leq u \leq 1} |W(u)|$$

in distribution.

From the table of Brownian Motion, find c such that
 $P(\sup_{0 \leq u \leq 1} |W(u)| \leq c) = 95\%.$

► **Approach II** (resampling) (Lin et al, 1993, Biometrika)

Recall to choose c such that $P(\sup_{0 \leq t < \infty} |Q_n(t)| \leq c) = 95\%$:

$Q_n(t) = \sqrt{n} \frac{\hat{H}_{NA}(\cdot) - H(\cdot)}{\hat{\sigma}(t)}$ is about

$$\begin{aligned} & \frac{\sqrt{n}}{\hat{\sigma}(t)} \left[\int_0^t \frac{I(Y_\cdot(u) > 0)}{Y_\cdot(u)} (dN_\cdot(u) - Y_\cdot(u)dH(u)) \right] \\ &= \frac{\sqrt{n}}{\hat{\sigma}(t)} \sum_{i=1}^n \int_0^t \frac{I(Y_\cdot(u) > 0)}{Y_\cdot(u)} dM_i(u) \end{aligned}$$

Define $\tilde{Q}_n(t)$ as follows

$$\frac{\sqrt{n}}{\hat{\sigma}(t)} \sum_{i=1}^n \int_0^t \frac{I(Y_\cdot(u) > 0)}{Y_\cdot(u)} dN_i(u) Z_i$$

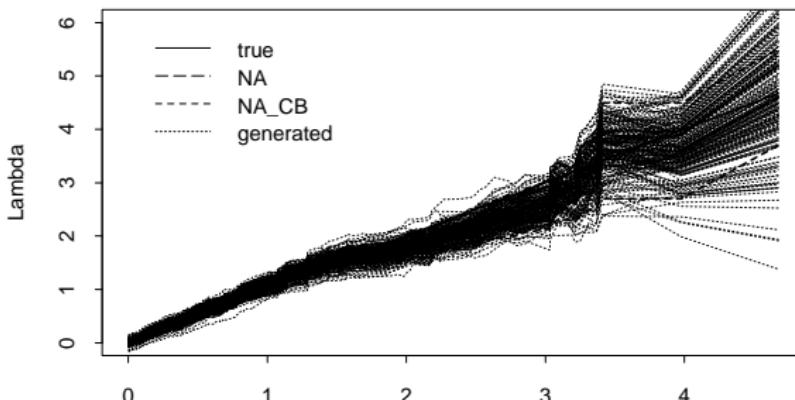
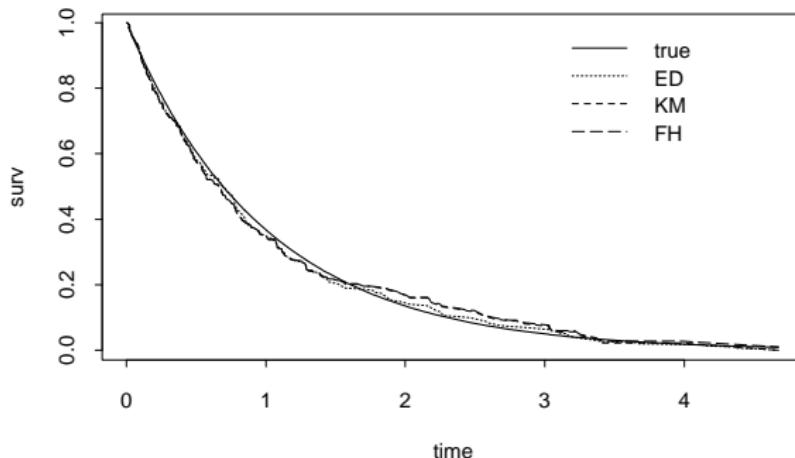
with Z_1, \dots, Z_n iid from $N(0, 1)$ and $\perp\!\!\!\perp$ the data:

$\tilde{Q}_n(\cdot)$ converges weakly to $Gaussian(0, \sigma^2(t))$ (i.e. asymptotically equivalent to $\sigma(t)Q_n(\cdot)$).

Algorithm. At l th time, $l = 1, \dots, M$,

- ▶ Step A. Generate $(Z_1, \dots, Z_n)^{(l)}$ as an iid sample from $N(0, 1)$ and $\perp\!\!\!\perp$ the data.
- ▶ Step B. Evaluate $\tilde{Q}_n^{(l)}(\cdot)$ over $[0, T^*]$.
- ▶ Step C. Calculate $W^{(l)} = \sup_{0 \leq s \leq T^*} |\tilde{Q}_n^{(l)}(s)|$.

Choose the critical c to be the 95% quantile of $W^{(1)}, \dots, W^{(M)}$.



Part IV.1.2E Revisit to Kaplan-Meier Estimator

Recall, with right-censored event times due to indpt censoring,

$$\hat{S}_{KM}(t) = \prod_{V_j < t} \left(1 - \frac{d_j}{N_j}\right)$$

Distinct observed event times $0 \leq V_1 < \dots < V_J \leq V_{J+1}$: $d_j = \#$ of observed event time at V_j , $N_j = \#$ of individuals at risk at V_j .
Plus $\hat{S}_{KM}(t) - S(t)$ is

$$\prod_{s < t} \left(1 - \frac{\Delta N(s)}{Y(s)}\right) - S(t) = -S(t) \int_0^t \frac{\hat{S}_{KM}(s-) I(Y(s) > 0)}{S(s) Y(s)} dM(s) + B(t)$$

- ▶ $\sqrt{n}B(t) \rightarrow 0$ in Pr 1 as $n \rightarrow \infty$.
- ▶ $\sqrt{n}[\hat{S}_{KM}(t) - S(t)] \rightarrow Gaussian(0, \sigma^2(t))$ in distn.
Thus $\sup_{t > 0} |\hat{S}_{KM}(t) - S(t)| \rightarrow 0$ (*uniform consistency*).

Part IV.1.2E Revisit to Kaplan-Meier Estimator

Note that $\sqrt{n}[\hat{S}_{KM}(t) - S(t)]$ can be approximated by

$$Q_n(t) = -S(t) \int_0^t \frac{\sqrt{n}I(Y_{\cdot}(s) > 0)}{Y_{\cdot}(s)} dM_{\cdot}(s)$$

Define

$$\tilde{Q}_n(t) = -\hat{S}_{KM}(t) \sum_{i=1}^n \int_0^t \frac{\sqrt{n}I(Y_{\cdot}(s) > 0)}{Y_{\cdot}(s)} dN_i(s) Z_i$$

$Z_1, \dots, Z_n \sim N(0, 1)$ iid and $\perp\!\!\!\perp$ the right-censored data.

\implies the resampling method to determine the critical value c for the CB (alternative to the nonparametric bootstrap approach)

Part IV.1.2E Revisit to Kaplan-Meier Estimator

Remarks:

► Product-Integral

$$\Pi_{s \in (0, t]} (1 + dX) = \lim_{\max |t_i - t_{i-1}| \rightarrow 0} \prod (1 + \Delta X(t_i))$$

$$S(t) = \exp\{-\Lambda(t)\} = \Pi_{s \in (0, t]} (1 - d\Lambda(s)) = \Pi_{s \in (0, t]} (1 - \lambda(s)ds)$$

► Applications of CB

- Model checking
- Two-Sample Problem: $H_0 : S_1(\cdot) = S_2(\cdot)$
 - nonparametric comparison?
CB of $\Delta(t) = S_1(t) - S_2(t)$ or $\Delta(t) = \Lambda_1(t) - \Lambda_2(t)$
e.g. Hu and Lagakos (Biometrika, 1999); Zhao et al
(Biostatistics, 2009)
 - semiparametric comparison?
Assume $\lambda_2(t) = \theta\lambda_1(t)$ and test on $H_0 : \theta = 1$
- K-Sample Problem?

What to study next?

Part IV. Advanced Topics

- ▶ *Part IV.1 Counting Process Formulation* (Revisits to KM estm, Logrank test, and Cox PH model)
- ▶ **Part IV.2 Selected Recent Topics in LIDA**
 - ▶ **Part IV.2.1 Alternatives to Cox PH model**
 - ▶ **Part IV.2.2 Multivariate event times**
 - ▶ *Part IV.2.3 More incomplete data structures*
 - ▶ *Part IV.2.4 Missing covariates in regression*
- ▶ *Part IV.3 Beyond Lifetime Data Analysis*