

# What to do today (2022/02/03)?

*Part I. Preliminaries*

*Part II. Parametric Inference*

**Part III. Nonparametric/Semi-parametric Inference**

*Part III.1. Introduction and Overview: Motivation*

**Part III.2. Kaplan-Meier Estimator**

*Part III.3. Nonparametric Tests*

*Part III.4. Cox Proportional Hazards Function*

*Part IV. Advanced Topics*

## Part III.2. Kaplan-Meier Estimator: Discrete Case

Consider a discrete event time  $T \sim F(\cdot)$  with mass points at  $0 \leq V_1 < \dots < V_K$

- ▶  $h_1 = P(T = V_1); h_j = P(T = V_j | T > V_{j-1})$
- ▶ For  $t \in [V_j, V_{j+1}), S(t) = P(T > t) = P(T > V_j) = P(T > V_j | T > V_{j-1})P(T > V_{j-1}) = \prod_{l=1}^j (1 - h_l)$
- ▶  $P(T = V_j) = h_j S(V_{j-1})$

$$L(F) = \prod_{i=1}^n f(u_i)^{\delta_i} S(u_i)^{1-\delta_i} = \prod_{j=1}^K h_j^{n_j} (1 - h_j)^{N_j - n_j}$$

with  $n_j = \# \text{ who fail at } V_j = \sum \delta_i I(u_i = V_j),$

$N_j = \# \text{ who at risk at } V_j = \sum I(u_i \geq V_j)$

## Part III.2. Kaplan-Meier Estimator: Discrete Case

⇒ the MLE  $\hat{h}_j = \frac{n_j}{N_j}$ . Thus

$$\hat{S}(t) = \begin{cases} 1 & t \leq V_1 \\ \prod_{l=1}^j (1 - \hat{h}_l) & V_j < t \leq V_{j+1} \\ 0 & t > V_K \end{cases}$$

*left-continuous*

**Example.**  $V_j : 2, 3, 5, 7, 9, 11, 16, 18, 20$ ; Data (n=10): 2, 2, 3+, 5, 5+, 7, 9, 16, 16, 18+

## Part III.2. Kaplan-Meier Estimator: A General Case

In general,  $F \in \mathcal{F} = \{\text{all cdfs}\}$

With the right-censored data, the likelihood function

$$L(F) = \prod_{i=1}^n dF(u_i)^{\delta_i} [1 - F(u_i)]^{1-\delta_i}$$

Maximize  $L(F)$  as  $F(\cdot)$  having only masses at the distinct observed event times:  $0 = V_0 \leq V_1 < \dots < V_J \leq V_{J+1}$   
⇒ the Kaplan-Meier estimator (left-continuous)

$$\hat{S}(t) = \prod_{j: V_j < t} \left(1 - \frac{n_j}{N_j}\right) = \begin{cases} 1 & t \leq V_1 \\ \prod_{l=1}^j (1 - \hat{h}_j) & V_j < t \leq V_{j+1} \\ ? & t > V_{J+1} \end{cases}$$

## Part III.2. Kaplan-Meier Estimator: More ...

### Remarks

- ▶ the general case vs the discrete case
- ▶ the observed sample mean (reduced sample estm) vs KM estm
- ▶ asymptotic properties of  $\hat{S}_{KM,n}(\cdot)$ 
  - ▶  $\hat{S}_{KM,n}(\cdot) \rightarrow S(\cdot)$  a.s. uniformly  
i.e.  $\sup_{0 < t < \infty} |\hat{S}_{KM,n}(t) - S(t)| \rightarrow 0$  a.s.
  - ▶ weak convergence  $\sqrt{n}(\hat{S}_{KM,n}(t) - S(t)) \rightarrow$  Gaussian Process with mean zero in distribution
- ▶ interval estm
  - ▶ pointwise CI: for  $t > 0$ ,  $\hat{S}_{KM}(t) \pm 1.96 \sqrt{\hat{Var}(\hat{S}_{KM}(t))}$
  - ▶ confidence band:  
 $P(S(t) \in (\hat{S}_L(t), \hat{S}_U(t)) : t \in (0, \infty)) \geq 95\%$   
how to compute CB?

Recall a pointwise CI: for  $t > 0$ ,  $\hat{S}_{KM}(t) \pm 1.96 \sqrt{\hat{Var}(\hat{S}_{KM}(t))}$

**an alternative way to construct a CI for  $S(t)$ :**

- ▶ to obtain a CI for  $\log S(t)$  first

$$\log \hat{S}_{KM}(t) \pm 1.96 \sqrt{\hat{V}(\log \hat{S}_{KM}(t))}$$

- ▶  $\hat{V}(\log \hat{S}_{KM}(t)) \approx \sum_{l=1}^j Var[\log(1 - \hat{h}_l)]$  for  $t \in [V_j, V_{j+1})$

- ▶  $Var[\log(1 - \hat{h}_l)] \approx Var(\hat{h}_l) \frac{1}{(1 - \hat{h}_l)^2}$  by the  $\Delta$ -method.

- ▶  $Var(\hat{h}_l) \approx \frac{1}{N_l} \frac{n_l}{N_l} \left(1 - \frac{n_l}{N_l}\right)$

- ▶ to obtain a CI for  $S(t)$  as

$$\exp \left\{ \log \hat{S}_{KM}(t) \pm 1.96 \sqrt{\hat{V}(\log \hat{S}_{KM}(t))} \right\} = \hat{S}_{KM}(t) e^{\pm 1.96 \sqrt{\hat{V}(\log \hat{S}_{KM}(t))}}$$

Thus, the alternative pointwise CI: for  $t > 0$ ,

$$(\hat{S}_{KM}(t)e^{-1.96\sqrt{\hat{Var}(\hat{S}_{KM}(t))}}, \hat{S}_{KM}(t)e^{1.96\sqrt{\hat{Var}(\hat{S}_{KM}(t))}})$$

**any other alternative constructions for a CI of  $S(t)$ ?**

- ▶ the logit transformation?  
*(proportional odds failure time models)*
  
- ▶ the probit transformation?

## Part III.2. Kaplan-Meier Estimator: Applications

- ▶ for comparing two populations' distn with censored data  
e.g.  $\sup_{t>0} |\hat{S}_{1,KM}(t) - \hat{S}_{2,KM}(t)|$ ? an extension of the Kolmogorov-Smirnov test statistic  $\sup_{t>0} |F_{1,n}(t) - F_{2,m}(t)|$   
*no need to specify the population distributions into parametric models*
- ▶ for justifying actuarial life table

## Part III.2. Kaplan-Meier Estimator: Applications

- ▶ for assessing parametric goodness-of-fit with censored data
  - ▶ e.g. is  $T \sim NE(\lambda)$  ( $H(t) = \lambda t$ )?  
     $\Rightarrow$  to check if  $\log S(t) = -\lambda t$ ?  
    using the scatter plot of  $\log \hat{S}(t)$  vs  $t$ : is  $\log \hat{S}(t)$  linear function of  $t$ ?
  - ▶ e.g. is  $T \sim Weibull(\lambda, \rho)$  ( $H(t) = \lambda t^\rho$ )?  
     $\Rightarrow$  to check if  $\log(-\log S(t)) = \log \lambda + \rho \log t$ ?  
    using the scatter plot of  $\log(-\log \hat{S}(t))$  vs  $\log t$ : look for linearity?

# What to study next?

*Part I. Preliminaries*

*Part II. Parametric Inference*

## Part III. Nonparametric/Semi-parametric Inference

- ▶ *Part III.1. Introduction and Overview*
- ▶ *Part III.2. Kaplan-Meier Estimator*
- ▶ **Part III.3. Nonparametric Tests**
- ▶ *Part III.4. Cox Proportional Hazards Function*

*Part IV. Advanced Topics*