

Estimating a population cumulative incidence under calendar time trends

Stefan N. Hansen^{*1}, Morten Overgaard¹, Per K. Andersen², and Erik T. Parner¹

¹*Section for Biostatistics, Aarhus University*

²*Section of Biostatistics, University of Copenhagen*

BMC Medical Research Methodology

Additional file 1

The variance of the Kaplan–Meier estimator of the event risk within the study period

We derive the weighted average of the stratum-wise Kaplan–Meier estimators as a non-parametric maximum likelihood estimator and its asymptotic variance when the π_i 's are considered unknown.

Suppose we have a sample consisting of n_i observations from the i th stratum with $n = \sum_i n_i$ being the total sample size. Due to right-censoring we observe $t_{ij} = \tilde{t}_{ij} \wedge c_{ij}$, where \tilde{t}_{ij} is the actual time-to-event and c_{ij} is the censoring time, and $\delta_{ij} = \mathbf{1}_{\tilde{t}_{ij} \leq c_{ij}}$ for $i = 1, \dots, k$ and $j = 1, \dots, n_i$ with observations being independent. Under independent right-censoring the likelihood is

$$\begin{aligned} L(\pi_1, \dots, \pi_k, S_1, \dots, S_k) &= \prod_{i=1}^k \prod_{j=1}^{n_i} \pi_i f_i(t_{ij})^{\delta_{ij}} S_i(t_{ij})^{1-\delta_{ij}} \\ &= \prod_{i=1}^k \pi_i^{n_i} \prod_{j=1}^{n_i} f_i(t_{ij})^{\delta_{ij}} S_i(t_{ij})^{1-\delta_{ij}} \end{aligned}$$

with (S_1, \dots, S_k) and (f_1, \dots, f_k) being the conditional survival functions and conditional densities in the k strata respectively and (π_1, \dots, π_k) being the probabilities of belonging to the strata. The corresponding non-parametric likelihood function [2] is thus

$$\tilde{L}(\pi_1, \dots, \pi_k, S_1, \dots, S_k) = \prod_{i=1}^k \pi_i^{n_i} \prod_{j=1}^{n_i} \Delta S_i(t_{ij})^{\delta_{ij}} S_i(t_{ij})^{1-\delta_{ij}}$$

with $\Delta S_i(t_{ij}) = S_i(t_{ij-}) - S_i(t_{ij})$ being the jump at t_{ij} . Maximizing this non-parametric likelihood function yields the estimates $\hat{\pi}_i = n_i/n$ which is the observed proportion and $\hat{S}_i = (\hat{S}_i(t))_{t \geq 0}$, the Kaplan–Meier estimator in the i th stratum, for $i = 1, \dots, k$.

Let $S_w(t_1) = P(T \leq U)$ with $U = \sum_{i=1}^k t_i \mathbf{1}_{B=i}$ and let

$$\hat{S}_w(t_1) = \sum_{i=1}^k \frac{n_i}{n} \hat{S}_i(t_i)$$

denote the Kaplan–Meier estimator of the within-study event risk. If $\hat{\boldsymbol{\pi}} = (\hat{\pi}_1, \dots, \hat{\pi}_k)^\top$ is the vector of observed proportions and $\hat{\mathbf{S}} = (\hat{S}_1(t_1), \dots, \hat{S}_k(t_k))^\top$ is the vector of Kaplan–Meier estimates evaluated at the end of follow-up times t_1, \dots, t_k , then by properties of the non-parametric maximum likelihood estimates we have

$$\sqrt{n} \left(\begin{pmatrix} \hat{\boldsymbol{\pi}} \\ \hat{\mathbf{S}} \end{pmatrix} - \begin{pmatrix} \boldsymbol{\pi} \\ \mathbf{S} \end{pmatrix} \right) \xrightarrow{\mathcal{D}} \mathcal{N}_{2k}(\mathbf{0}, \begin{pmatrix} \boldsymbol{\Sigma}_1 & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\Sigma}_2 \end{pmatrix}),$$

^{*}Corresponding author; Address: Bartholins Allé, DK-8000 Aarhus C – E-mail: stefanh@ph.au.dk

where $\boldsymbol{\pi} = (\pi_1, \dots, \pi_k)^\top$, $\mathbf{S} = (S_1(t_1), \dots, S_k(t_k))^\top$. Here

$$\boldsymbol{\Sigma}_1 = \text{diag}(\boldsymbol{\pi}) - \boldsymbol{\pi}\boldsymbol{\pi}^\top, \quad \boldsymbol{\Sigma}_2 = \text{diag}(\sigma_1^2/\pi_1, \dots, \sigma_k^2/\pi_k),$$

for some $\sigma_1^2, \dots, \sigma_k^2 > 0$.

The delta method [1] on $g(x_1, \dots, x_{2k}) = \sum_{i=1}^k x_i x_{k+i}$ now yields

$$\sqrt{n}(\widehat{S}_w(t_1) - S_w(t_1)) \xrightarrow{\mathcal{D}} \mathcal{N}_1(0, \sigma_w^2)$$

with

$$\sigma_w^2 = (\mathbf{S}^\top \quad \boldsymbol{\pi}^\top) \begin{pmatrix} \boldsymbol{\Sigma}_1 & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\Sigma}_2 \end{pmatrix} \begin{pmatrix} \mathbf{S} \\ \boldsymbol{\pi} \end{pmatrix} = \mathbf{S}^\top \boldsymbol{\Sigma}_1 \mathbf{S} + \boldsymbol{\pi}^\top \boldsymbol{\Sigma}_2 \boldsymbol{\pi}.$$

The variance of this Kaplan–Meier estimator is thus approximately given by

$$\text{Var}(\widehat{S}_w(t_1)) \approx \frac{\sigma_w^2}{n} = \frac{1}{n} \left[\sum_{i=1}^k \pi_i S_i(t_i)^2 - \left(\sum_{i=1}^k \pi_i S_i(t_i) \right)^2 + \sum_{i=1}^k \pi_i \sigma_i^2 \right].$$

If $\widehat{\text{Var}}(\widehat{S}_i(t_i))$ is Greenwood's estimate of the variance σ_i^2/n_i of the Kaplan–Meier estimator in the i th group, then σ_i^2 can be estimated by $n_i \widehat{\text{Var}}(\widehat{S}_i(t_i))$. Thus, we may estimate the variance by

$$\widehat{\text{Var}}(\widehat{S}_w(t_1)) = \frac{1}{n} \left[\sum_{i=1}^k \frac{n_i}{n} \widehat{S}_i(t_i)^2 - \left(\sum_{i=1}^k \frac{n_i}{n} \widehat{S}_i(t_i) \right)^2 \right] + \sum_{i=1}^k \frac{n_i^2}{n^2} \widehat{\text{Var}}(\widehat{S}_i(t_i)).$$

References

- [1] Lehmann, E.L.: Elements of Large-Sample Theory. Springer-Verlag, New York (1999)
- [2] van der Vaart, A.W.: Asymptotic Statistics. Cambridge University Press, Cambridge (2000)