# Kaplan–Meier
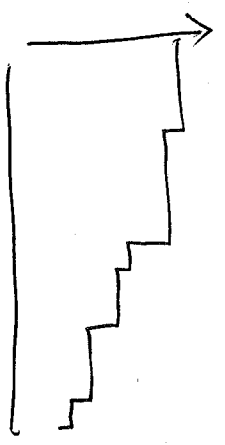
Product limit estimator = limit of LTE when interval length is small, such that at most one distinct event occurs (no ties).



$$\hat{KM}(t) = \hat{S}(t) = \overline{\prod_{x \le t} \{1 - \hat{m}(x)\}} \quad \text{where } \hat{m}(x) = \frac{d(x)}{n(x)}$$
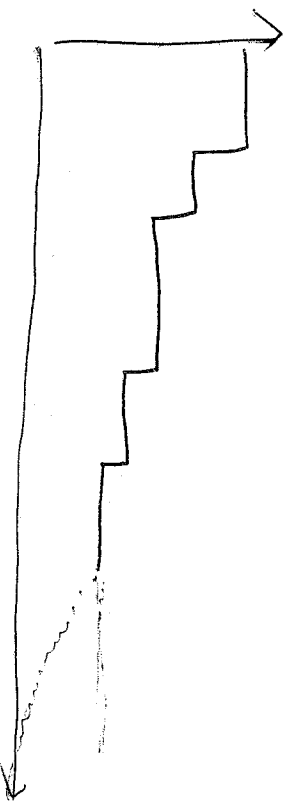
Step size $n(x)$ → Note! here in a 6.7.

← persons @ risk @ the pt $x$.

Definition based rate $P(x \le T < x + \Delta x \mid T \ge x) \approx h(x)\Delta x$
$$\Delta x \to 0$$

## NOTE

- If $t_{max}$ such that $\delta = 1$ (real event time)
  $\Rightarrow \hat{KM}(t) = 0 \quad t > t_{max}$

- If $t_{max}$ s.t. $\delta = 0$ (censored), then $\hat{KM}(t)$ not defined beyond $t_{max}$

  (Alts)
  → extend $\hat{S}(t) = \hat{S}(t_{max})$ (infinite survival for censored (unr))
  → $\hat{S}(t) = \exp\left[\frac{t}{t_{max}} \ln \hat{S}(t_{max})\right]$

KM and redistribution to the right; we obtain estimates of $S(t)$

Remember the LTE; we obtain estimates of $S(t)$

by adjusting/accounts for looking @ end/begining or

mid-interval.

Another way of looking at this.

$\Rightarrow$ Back to basics; of no censoring

$$\hat{S}(t) = \sum_{i=1}^{n} \frac{1}{n} 1\{t_i > t\}$$

$$= 1 - \sum_{i=1}^{n} \frac{1}{n} 1\{t_i \le t\}$$

stepsize $\frac{1}{n}$.

(but) what :) we now find out that a

$t_i$ is a censored observation?

$\Rightarrow$ adjust denominator to reflect this

This is called the "redistribute to the right" algorithm.

That is, @ $t_i = (T_i, \delta_i = 0)$, take the mass $\frac{1}{n}$

and redistribute equally over all $t_j > t_i$ (remaining time

points). From now on, the step size $\frac{1}{n(x)}$

If the next t you consider is $\delta = 1$ (real event), step size $\frac{1}{n}$ into

is used, o/w you again push the mass equally to

the right ...

The "reliability to the right" = KM.

Since $KM(t)$ is the limit of the LTE, variance formulae look very similar.

Greenwood's formula

$$V(\hat{S}_{km}(t)) = \left(\hat{S}_{km}(t)\right)^2 \sum_{x \leq t} \frac{d(x)}{n(x)(n(x)-d(x))}$$

## Nelson-Aalen estimator.

① $\dfrac{d(x)}{n(x)}$ is small $\Rightarrow \exp\left\{-\dfrac{d(x)}{n(x)}\right\} \approx \left\{1 - \dfrac{d(x)}{n(x)}\right\}$

and so $\widehat{km}(t) = \hat{S}_{km}(t) = \prod_{x \leq t}\left\{1 - \dfrac{d(x)}{n(x)}\right\} \approx \prod_{x \leq t} \exp\left(-\dfrac{d(x)}{n(x)}\right)$

$\qquad = \exp\left\{-\sum_{x \leq t}\dfrac{d(x)}{n(x)}\right\}$

Definition of Basic Concepts

Cumulative Hazard $H(t) = -\ln S(t)$

Alt 1: $\hat{H}_{km}(t) = -\ln \hat{S}_{km}(t)$

Alt 2: $\hat{H}_{NA}(t) = \sum_{x \leq t}\dfrac{d(x)}{n(x)}$

$\qquad \hat{S}_{NA}(t) = \exp\left(-\hat{H}_{NA}(t)\right)$

Turns out, the second alternative

$\hat{H}_{NA}(t) = \sum_{x \leq t} \frac{d(x)}{n(x)}$ has better small-sample properties.

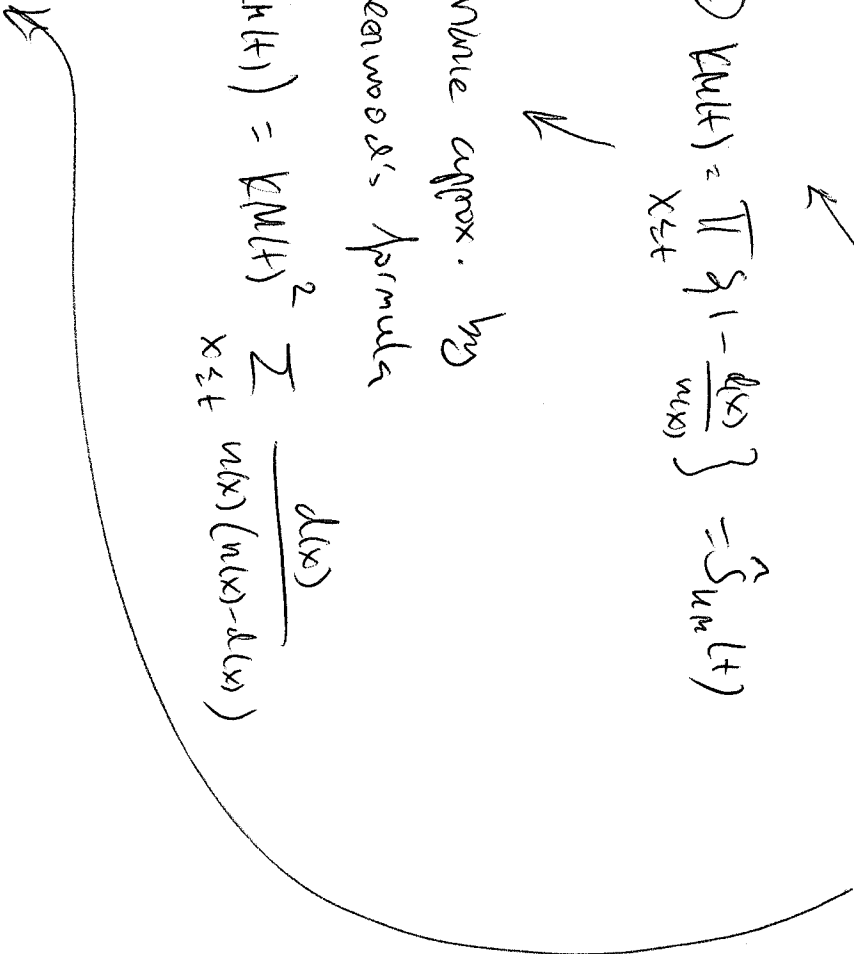So, we have two alternative estimates of survival

① $\hat{K}_M(t) = \prod_{x \leq t} \left\{ 1 - \frac{d(x)}{n(x)} \right\} = \hat{S}_{KM}(t)$

Varance approx. by Greenwood's formula

$\hat{V}(\hat{K}_M(t)) = \hat{K}_M(t)^2 \sum_{x \leq t} \frac{d(x)}{n(x)(n(x)-d(x))}$

② Nelson - Aalen

$\hat{H}_{NA}(t) = \sum_{x \leq t} \frac{d(x)}{n(x)} \Rightarrow \hat{S}_{NA}(t) = \exp(-\hat{H}_{NA}(t))$

To get Variance approx for $\hat{S}_{NA}(t)$, we again rely on results from counting processes (once lahd increments), and the delta-method.

after some
$\xrightarrow{\text{derivation}}$

$$\hat{V}(\hat{H}_{Na}(t)) = \sum_{x \leq t} \frac{d(x)}{n(x)} \frac{\left(1 - \frac{d(x)}{n(x)}\right)}{n(x) - 1}$$

if ⅋ assume
$=$ that $d(x) = 0$ or $1$ only (no ~~ties~~)

$$\sum_{x \leq t} \frac{d(x)}{(n(x))^2}$$

$\left(\begin{array}{c}\text{unbiased} \\ \text{estimator of} \\ \hat{V}(\hat{H}_{Na}(t))\end{array}\right)$  ⑭

Summary:

① $\hat{S}_{en}(t) = KM(t) = \prod_{x \leq t} \left\{1 - \frac{d(x)}{n(x)}\right\}$

$\hat{V}(\hat{S}_{en}(t)) = \hat{S}_{en}(t)^2 \sum_{x \leq t} \frac{d(x)}{n(x)(n(x) - d(x))}$

② $\hat{H}_{Na}(t) = \sum_{x \leq t} \frac{d(x)}{n(x)}$

$\hat{V}(\hat{H}_{Na}(t)) = \sum_{x \leq t} \frac{d(x)}{n(x)^2}$ , $\hat{S}_{Na}(t) = \hat{S}_{Na}(t)^2 \exp\left(-\hat{H}_{Na}(t)\right)$

$\longrightarrow \hat{V}(\hat{S}_{Na}(t)) = \hat{S}_{Na}(t)^2 \sum_{x \leq t} \frac{d(x)}{n(x)^2}$

Note, estimates of variance only really OK for large samples. Also, be cautious when $n(x)$ is small (tail of data).

# Confidence intervals for $S(t)$

Counting process theory $\Rightarrow \hat{K}_u(t)$ and $\hat{S}_{Na}(t) \xrightarrow{n \to \infty}$ Gaussian processes

④ remember this is a large-sample result.

$\rightarrow$ alternative methods for setting up $CI$ for $S(t)$

① [ Most commonly used ]

KM + Greenwood

$$\left[ \hat{S}_{u}(t_0) \pm Z_{1-\alpha/2} \sqrt{\hat{V}_{Greenwood}(\hat{S}_{u}(t_0))} \right]$$

② Use KM to estimate $\hat{S}_{u}(t_0)$ and $\hat{H}_{u}(t_0) = -\ln \hat{S}(t)$

$\Rightarrow \hat{V}(\hat{H}_{u}(t_0)) = \sum\limits_{x: t_x \leq t_0} \dfrac{d(x)}{n(x)(n(x) - d(x))}$

$\Rightarrow CI$ for $S(t_0)$ by exponentiating $CI$ bounds

of $\hat{H}_{u}(t_0)$

$$\exp \left[ \hat{H}_{u}(t_0) \pm Z_{1-\alpha/2}^{*} \sqrt{\hat{V}(\hat{H}_{u}(t_0))} \right] \quad [ \qquad ]$$

③ Use the Nelson-Aalen estimator of $H(t)$, Taylor expand to get $\hat{S}_{NA}(t)$ and use the delta-method

$$\left[ \hat{S}_{NA}(t_0) \pm z_{1-\alpha/2} \sqrt{\hat{S}_{NA}(t_0)^2 \sum_{x \leq t_0} \frac{d(x)}{n(x)^2}} \right]$$

④ [ Usually works well in practice ]

Use the Nelson-Aalen estimator of $H(t)$, establish CI for $H(t)$, exponentiate the CI bounds

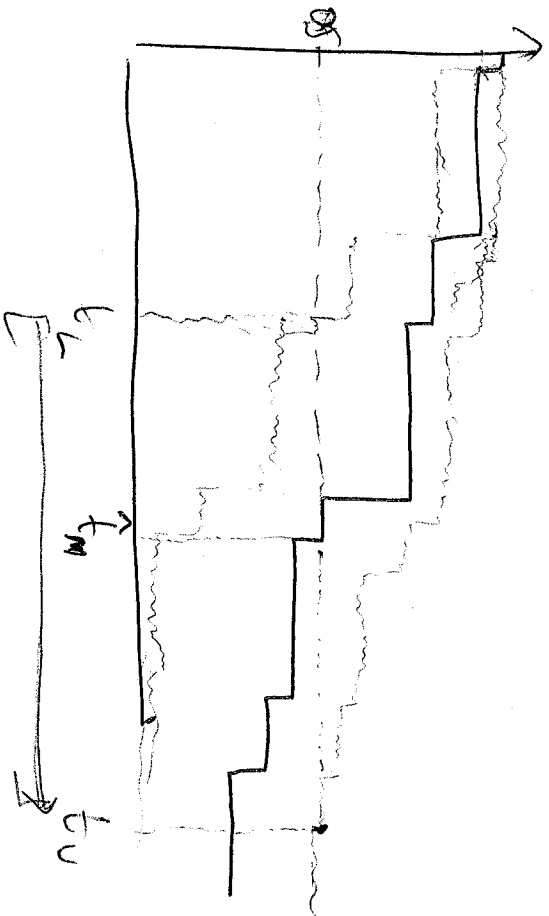$$\exp\left[ \hat{H}_{NA}(t_0) \pm z_{1-\alpha/2} \sqrt{\sum_{x \leq t_0} \frac{d(x)}{n(x)^2}} \right]$$

NOTE : These CI are <u>pointwise</u>

→ another lecture (part of) we'll talk about
confidence <u>bands</u> (simultaneous CI).

③

⑥

Another quantity of interest — median survival time

$$t_M = \{t : S(t) = \tfrac{1}{2}\}$$



Find $t_{m,U}$ s.t. $\hat{S}(t_{m,U}) + z_{\alpha/2}\, SE(\hat{S}(t_{m,U})) = \tfrac{1}{2}$

and $t_{m,L}$ s.t. $\hat{S}(t_{m,L}) - z_{\alpha/2}\, SE(\hat{S}(t_{m,L})) = \tfrac{1}{2}$

$\Rightarrow$ An approx $1-\alpha$ CI for $t_m$ is thus
$$[\hat{t}_{m,L},\ \hat{t}_{m,U}]$$

(Note; $P(\hat{t}_{m,L} < t_m < \hat{t}_{m,U}) = P(S(\hat{t}_{m,U}) < \tfrac{1}{2} < S(\hat{t}_{m,L}))$

$= 1 - \left( P(S(\hat{t}_{m,U}) > \tfrac{1}{2}) + P(S(\hat{t}_{m,L}) < \tfrac{1}{2}) \right)$

$\qquad\qquad\left( \text{S is a decreasing function so by def } S(t_m) = \tfrac{1}{2} \right)$

$\left( \text{Since; } S(\hat{t}_{m,L}) > S(\hat{t}_{m,U}) \text{ for } \hat{t}_L < \hat{t}_U \right.$

$\quad$ and $P(S(\hat{t}_L) > \tfrac{1}{2}) = P(S(\hat{t}_U) > \tfrac{1}{2},\ S(\hat{t}_L) > \tfrac{1}{2})$

$\quad$ $P(S(\hat{t}_L) < \tfrac{1}{2}) = P(S(\hat{t}_U) < \tfrac{1}{2},\ S(\hat{t}_L) < \tfrac{1}{2})$

$\Rightarrow$

$$\Rightarrow P(S(\hat{t}_U) > \tfrac{1}{2}) + P(S(\hat{t}_L) < \tfrac{1}{2})$$

$$= P(S(\hat{t}_U) > \tfrac{1}{2}, S(\hat{t}_L) > \tfrac{1}{2}) + P(S(\hat{t}_L) < \tfrac{1}{2}, S(\hat{t}_U) < \tfrac{1}{2})$$

$$= P(\text{both above}) + P(\text{both below})$$

$$\geq 1 - P(\text{one above, one below})$$

$$\geq 1 - P(S(\hat{t}_U) < \tfrac{1}{2} < S(\hat{t}_L))$$

$\hat{t}_L$, so $P(\hat{t}_L < t_n < \hat{t}_U) = 1 - \left(P(S(\hat{t}_U) > \tfrac{1}{2}) + P(S(\hat{t}_L) < \tfrac{1}{2})\right)$

Now, if $\hat{t}_U$ is the solution to $S(t_U) + Z_{\alpha/2} SE(\hat{S}(t_U)) = \tfrac{1}{2}$

then $\hat{t}_U = t_U \Rightarrow P(S(\hat{t}_U) > \tfrac{1}{2}) = P(S(t_U) > \hat{S}(\hat{t}_U))$

$$= P\left(S(t_U) + Z_{\alpha/2} SE(\hat{S}(\hat{t}_U)) = \tfrac{1}{2}\right)$$

$$= P\left(\frac{\hat{S}(\hat{t}_U) - S(t_U)}{SE(\hat{S}(\hat{t}_U))} < -Z_{\alpha/2}\right)$$

$$= P\left(Z < -Z_{\alpha/2}\right) = \alpha/2 \qquad \left(\hat{S} \to \text{bum proies}\right)$$

and similarly

$$P(S(\hat{t}_L) < \tfrac{1}{2}) = \alpha/2$$

$$\Rightarrow P(\hat{t}_L < t_n < \hat{t}_U) = 1 - \alpha$$