

chapter ③: Estimation of incomplete data.

① Complete data:

① Empirical distribution function:

Data: x_1, \dots, x_n .

The empirical distribution function (edf) is given by:

$$F_n(x) = \frac{\# \text{ of values } y \leq x}{n}$$

Ex: $5, 2, 5, 10, 1$. Compute the edf?

$$F_5(x) = \begin{cases} 0 & x < 1 \\ 1/5 & 1 \leq x < 2 \\ 2/5 & 2 \leq x < 5 \\ 4/5 & 5 \leq x < 10 \\ 1 & x \geq 10 \end{cases}$$

* $y_1 < y_2 < \dots < y_k$ (we arrange x_i 's without repetition)
 $s_i =$ number of repetitions of y_i in the sample.

$$r_i = s_i + s_{i+1} + \dots + s_k$$

* we can write:

$$F_n(x) = \begin{cases} 0 & x < y_1 \\ 1 - \frac{r_i}{n} & y_{i-1} \leq x < y_i, i=2, \dots, k \\ 1 & x \geq y_k \end{cases}$$

Ex:

~~5, 2, 5, 10, 1~~

i	y_i	s_i	r_i
1	1	1	5
2	2	1	4
3	5	2	3
4	10	1	1

$$F_S(x) = \begin{cases} 0 & x < 1 \\ 1 - \frac{4}{5} = \frac{1}{5} & 1 \leq x < 2 \\ 1 - \frac{2}{5} = \frac{3}{5} & 2 \leq x < 5 \\ 1 - \frac{1}{5} = \frac{4}{5} & 5 \leq x < 10 \\ 1 & x \geq 10 \end{cases}$$

(2) Nelson-Aalen estimate:

* $H(x)$ = Cumulative hazard rate

$$H(x) = \int_0^x h(t) dt = \int_0^x -\frac{S'(t)}{S(t)} dt = -\ln(1 - F(x))$$

$$F(x) = 1 - e^{-H(x)}$$

* The Nelson-Aalen estimate for the cumulative hazard rate function is:

$$H(x) = \begin{cases} 0 & x < y_1 \\ \frac{s_1}{r_1} + \dots + \frac{s_{i-1}}{r_{i-1}} & y_{i-1} \leq x < y_i, i=2, \dots, k \\ \frac{s_1}{r_1} + \dots + \frac{s_k}{r_k} & x \geq y_k \end{cases}$$

Ex: Compute the Nelson-Aalen estimate for

previous example:

$$\frac{29}{60} + 1 = \frac{89}{60}$$

(2)

$$H(x) = \begin{cases} 0 & x < 1 \\ \frac{1}{5} & 1 \leq x < 2 \\ \frac{1}{5} + \frac{1}{4} = \frac{9}{20} & 2 \leq x < 5 \\ \frac{1}{5} + \frac{1}{4} + \frac{2}{3} = \frac{29}{60} & 5 \leq x < 10 \\ 1 & x \geq 10 \end{cases}$$

(3) Empirical distribution for grouped data:

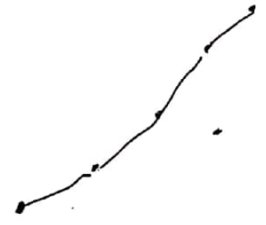
	# of observations
$(c_0, c_1]$	n_1
$(c_1, c_2]$	n_2
\vdots	\vdots
$(c_{k-1}, c_k]$	n_k
(c_k, ∞)	0

$$n = n_1 + \dots + n_k$$

$$F_n(c_j) = \frac{1}{n} \sum_{i=1}^j n_i ; F_n(c_0) = 0$$

$$F_n(x) = \frac{c_j - x}{c_j - c_{j-1}} F_n(c_{j-1}) + \frac{x - c_{j-1}}{c_j - c_{j-1}} F_n(c_j)$$

$$c_{j-1} \leq x \leq c_j$$



Ex:

Given Interval	# of observations	$F_{50}(\cdot)$
$(0, 2]$	25	$25/50 = 0.5$
$(2, 10]$	10	$35/50 = 0.7$
$(10, 100]$	10	$45/50 = 0.9$
$(100, 1000]$	5	$50/50 = 1$

$n = 50$

Find the edf?

$$F_{50}(x) = \frac{2-x}{2} F_{50}(0) + \frac{x-0}{2} F_{50}(2) = 0.25x$$

$$x \leq 2:$$

$$F_{50}(x) = \frac{10-x}{8} F_{50}(2) + \frac{x-2}{8} F_{50}(10)$$

$$2 \leq x \leq 10:$$

$$= \frac{0.5(10-x) + 0.7(x-2)}{8}$$

$$= \frac{0.2x + 3.6}{8} = 0.025x + 0.45$$

$$10 \leq x \leq 100:$$

$$F_{50}(x) = \frac{100-x}{90} F_{50}(10) + \frac{x-10}{90} F_{50}(100)$$

$$= \frac{(100-x)0.7 + (x-10)0.9}{90}$$

$$= \frac{0.2x + 61}{90} = 0.0022x + 0.67$$

(3)

$$100 \leq x \leq 1000 :$$

$$F_{50}(x) = \frac{1000-x}{900} F_{50}(100) + \frac{x-100}{900} F_{50}(1000)$$

$$= \frac{(1000-x)0.9 + (x-100)1}{900}$$

$$= \frac{0.1x + 800}{900} = 0.0001x + 0.888$$

$$F_{50}(x) = \begin{cases} 0.25x & 0 \leq x \leq 2 \\ 0.025x + 0.45 & 2 \leq x \leq 10 \\ 0.0022x + 0.67 & 10 \leq x \leq 100 \\ 0.0001x + 0.88 & 100 \leq x \leq 1000 \end{cases}$$

II

Incomplete data :

data: $x_1, x_2, x_3, \dots, x_n$

→ y_1, y_2, \dots, y_k

$$s_j = \sum_{i=1}^n 1_{(x_i = y_j)}$$

$$r_j = \sum_{i=1}^n 1_{(x_i \geq y_j)} + \sum_{i=1}^n 1_{(u_i \geq y_j)} - \sum_{i=1}^n 1_{(d_i \geq y_j)}$$

where : $u_i =$ censored point for the i^{th} observation.

$d_i =$ truncated point for the i^{th} observation.

Ex :

Ex:

i	d_i	x_i	u_i
1	0	-	4
2	0	0.5	-
3	0	-	1
4	0	-	4
5	1	-	4
6	1.2	2	-
7	1.5	2	-
8	2	-	3
9	2.5	-	4
10	3.1	3.2	-

Complete the table:

j	r_j	f_j	g_j
1	0.5	1	$4+6-6=4$
2	2	2	$3+5-3=5$
3	3.2	1	$1+4-0=5$

① Kaplan-Meier Method (product limit):
 $0 \leq t \leq \tau_1$

$$S_n(t) = \begin{cases} 1 & 0 \leq t \leq \tau_1 \\ \prod_{i=1}^{j-1} \left(1 - \frac{d_i}{r_i}\right) & \tau_{j-1} \leq t \leq \tau_j, j=2, \dots, k \\ \prod_{i=1}^k \left(1 - \frac{d_i}{r_i}\right) & \text{or } 0; t \geq \tau_k \end{cases}$$

Ex: Use the previous table to compute the Kaplan-Meier survival function:

→

⑤

$$S_{10} = \begin{cases} 1 - \frac{c_1}{r_1} = 1 - \frac{1}{4} = \frac{3}{4} & y_1 \leq t < y_2 \\ \frac{3}{4} \left(1 - \frac{c_2}{r_2}\right) = \frac{3}{4} \left(1 - \frac{2}{5}\right) = \frac{9}{20} & y_2 \leq t < y_3 \\ \frac{3}{20} \left(1 - \frac{c_3}{r_3}\right) = \frac{3}{20} \left(1 - \frac{1}{5}\right) = \frac{9}{25} & t \geq y_3 \end{cases}$$

Ex: 2, 3, 3, 5, 5⁺, 6, 7, 7⁺, 9, 10⁺.

where "+" indicates that the loss exceeded the policy limit.

Calculate $P(\text{loss} > 8)$ using Kaplan-Meier estimated edf?

i	d _i	r _i	u _i
1	0	2	-
2	0	3	-
3	1	3	-
4	1	5	5
5	1	6	-
6	1	7	-
7	1	7	7
8	1	9	-
9	1	9	-
10	1	10	10

i	y _i	s _i	r _i
1	2	1	7+3=10
2	3	2	6+3=9
3	5	1	4+3=7
4	6	1	3+2=5
5	7	1	2+2=4
6	9	1	1+1=2

$$S_{10}(8) = P(\text{loss} > 8) = \begin{cases} 1 - \frac{1}{10} = 0.9 & 0 \leq y < 2 \\ 0.9 \left(1 - \frac{1}{9}\right) = 0.7 & 2 \leq y < 3 \\ 0.7 \left(1 - \frac{1}{7}\right) = 0.6 & 3 \leq y < 5 \\ 0.6 \left(1 - \frac{1}{5}\right) = 0.48 & 5 \leq y < 6 \\ 0.48 \left(1 - \frac{1}{6}\right) = 0.36 & 6 \leq y < 7 \\ 0.36 \left(1 - \frac{1}{7}\right) = 0.36 & 7 \leq y < 9 \\ 0.36 & y > 9 \end{cases}$$

$$S_{10}(8) = 0.36$$

(6)

(2)

Nelson-Aalen method:

$$\hat{H}(t) = \begin{cases} 0 & t < \tau_1 \\ \sum_{i=1}^{j-1} \frac{d_i}{r_i} & \tau_{j-1} \leq t < \tau_j, j=2, \dots, k \\ \sum_{i=1}^k \frac{d_i}{r_i} & t \geq \tau_k \end{cases}$$

$$\hat{S}(t) = e^{-\hat{H}(t)}$$

Ex: Use the previous table to estimate the survival using Nelson-Aalen method.

$$\hat{H}(t) = \begin{cases} 0 & t < 2 \\ 0.1 & 2 \leq t < 3 \\ 0.1 + \frac{2}{9} = 0.32 & 3 \leq t < 5 \\ 0.32 + \frac{1}{7} = 0.46 & 5 \leq t < 6 \\ 0.46 + \frac{1}{5} = 0.66 & 6 \leq t < 7 \\ 0.66 + \frac{1}{4} = 0.91 & 7 \leq t < 9 \\ 0.91 + \frac{1}{2} = 1.41 & t \geq 9 \end{cases}$$

$$\hat{S}(7) = e^{-\hat{H}(7)} = e^{-0.91} = 0.4$$

III

Mean and variance for the empirical estimators:

(1) Complete data: • Individual data.

$$\hat{S}_n(t) = \frac{\sum_{i=1}^n 1_{(X_i > t)}}{n} = \frac{Y}{n}$$

$$1_{(X_i > t)} \sim \text{Bernoulli}(S(t)).$$

$$Y \sim \text{Binomial}(n, S(t)).$$

$$E(Y) = n S(t); \text{var}(Y) = n S(t) (1 - S(t)).$$

$$E(\hat{S}_n(t)) = S(t); \text{var}(\hat{S}_n(t)) = \frac{\text{var}(Y)}{n^2} = \frac{S(t)(1 - S(t))}{n} \rightarrow 0$$

(7)

estimated variance:

$$\hat{\text{Var}}(\hat{S}_n^1(\theta)) = \frac{\hat{S}_n^1(\theta)(1 - \hat{S}_n^1(\theta))}{n}$$

Ex: 4.9, 5.0, 5.0, 5.0, 6.0, 7.5, 8.0, 12.0, 13.0.
Find the estimated variance of the estimate $\hat{S}_q^1(6.0)$.

$$\hat{S}_q^1(6.0) = \frac{4}{9}$$

$$\hat{\text{Var}}(\hat{S}_q^1(6.0)) = \frac{\frac{4}{9}(1 - \frac{4}{9})}{9} = \frac{20}{9^2} = 0.024.$$

• Grouped data:

$$\hat{S}_n^1(x) = 1 - \frac{n_1 + \dots + n_j}{n} = \frac{n_{j+1} + \dots + n_k}{n}$$

$$c_{j-1} \leq x \leq c_j: \hat{S}_n^1(x) = \frac{c_j - x}{c_j - c_{j-1}} \hat{S}_n^1(c_{j-1}) + \frac{x - c_{j-1}}{c_j - c_{j-1}} \hat{S}_n^1(c_j)$$

$$\rightarrow E(\hat{S}_n^1(x)) = \frac{c_j - x}{c_j - c_{j-1}} S(c_{j-1}) + \frac{x - c_{j-1}}{c_j - c_{j-1}} S(c_j)$$

$$\rightarrow \text{Var}(\hat{S}_n^1(x)) = \frac{(c_j - c_{j-1})^2 \text{Var}(Y) + (x - c_{j-1})^2 \text{Var}(Z) + 2(c_j - c_{j-1})(x - c_{j-1}) \text{Cov}(Y, Z)}{n^2 (c_j - c_{j-1})^2}$$

$$\text{Var}(Y) = n S(c_{j-1})(1 - S(c_{j-1}))$$

$$\text{Var}(Z) = n (S(c_{j-1}) - S(c_j))(1 - \frac{S(c_{j-1}) + S(c_j)}{2})$$

$$\text{Cov}(Y, Z) = -n (1 - S(c_{j-1}))(S(c_{j-1}) - S(c_j))$$

$$\left[x = c_{j-1}: \text{Var}(\hat{S}_n^1(c_{j-1})) = \frac{S(c_{j-1})(1 - S(c_{j-1}))}{n} \right]$$

(8)

Ex:

interval	#
(0, 2]	25
(2, 10]	10
(10, 100]	10
(100, 1000]	5

Find the estimated variance of f the estimate of $\frac{1}{S_0}(100)$?

$$\frac{1}{S_0}(100) = \frac{5}{50} = 0.1$$
$$\hat{\text{Var}}\left(\frac{1}{S_0}(100)\right) = \frac{0.1(1-0.1)}{50} = \frac{0.1 \times 0.9}{50}$$
$$= \frac{0.09}{50} = 0.0018$$

(9)

4) Estimation of variance of the Kaplan-Meier estimator :

The Greenwood's approximation of the variance estimate for the Kaplan-Meier estimator is :

$$\text{Var}(S_n(y_j)) \approx (S_n(y_j))^2 \sum_{i=1}^j \frac{s_i}{r_i(r_i - s_i)}$$

Ex: Let consider 10 payments :

4, 4, 5⁺, 5⁺, 5⁺, 8, 10⁺, 10⁺, 12, 15.

Find the Greenwood's approximation to the variance of the Kaplan-Meier estimate $S_{10}(11)$?

j	y _j	s _j	r _j	S ₁₀ (y _j)
1	4	2	10	$1 - \frac{2}{10}$
2	8	1	5	$(1 - \frac{2}{10})(1 - \frac{1}{5})$
3	12	1	2	$(1 - \frac{1}{2})$
4	15	1	1	0

i	d _i	x _i	u _i
1	0	4	5
2	0	4	5
3	0	1	5
4	0	1	5
5	1	8	10
6	1	8	10
7	1	12	1
8	0	15	1

8 ≤ 11 < 12.

$$S_{10}(11) = S_{10}(8) = \frac{16}{25}$$

$$\text{Var}(S_{10}(8)) = \left(\frac{16}{25}\right)^2 \left[\frac{2}{10(10-2)} + \frac{1}{5(5-1)} \right]$$

$$= 0.03072$$

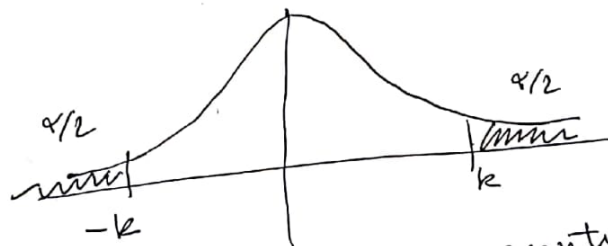
⑤ Variance estimate of the Nelson-Åalen estimator:

$$\text{Var}(\hat{H}(y_j)) = \sum_{i=1}^j \frac{s_i}{r_i^2}$$

Ex: (Previous example). Find the variance estimate of the Nelson-Åalen estimator at $y = 4$.

$$\text{Var}(\hat{H}(4)) = \text{Var}(\hat{H}(8)) = \frac{2}{10^2} + \frac{1}{5^2} = 0.06$$

⑥ Confidence interval:
 Let X a random variable with mean μ and variance σ^2 .
 Under the assumption that $X \sim N(\mu, \sigma^2)$,
 the confidence interval of μ is:
 $(\mu - k\sigma, \mu + k\sigma)$



$$k = (1 - \frac{\alpha}{2})_{100} \text{-percentile} = z_{\alpha/2}$$

→ $(1 - \alpha)$ -confidence interval:
 $(\mu - z_{\alpha/2} \sigma, \mu + z_{\alpha/2} \sigma)$

(6.1) Kaplan-Meier estimator:

The log-transformed $100(1-\alpha)\%$ confidence interval for $S_n(t)$:

$$\left(S_n(t)^{\frac{1}{u}}; S_n(t)^u \right),$$

$$u = \exp \left[z_{\alpha/2} \frac{\sqrt{\widehat{\text{Var}}(\hat{S}_n(t))}}{S_n(t) \ln S_n(t)} \right].$$

Ex: (Previous example) Find a 95% confidence interval for $S_{10}(11)$?

$$\hat{S}_{10}(11) = \frac{16}{25}; \quad \widehat{\text{Var}}(\hat{S}_{10}(11)) = 0.03072.$$

$$z_{0.025} = 1.96$$

$$u = \exp \left(1.96 \frac{\sqrt{0.03072}}{\frac{16}{25} \ln \left(\frac{16}{25} \right)} \right) = 0.3$$

Confidence interval:

$$\left(\left(\frac{16}{25} \right)^{\frac{1}{0.3}}; \left(\frac{16}{25} \right)^{0.3} \right) = (0.23; 0.87).$$

(6.2) Nelson-Aalen estimator:

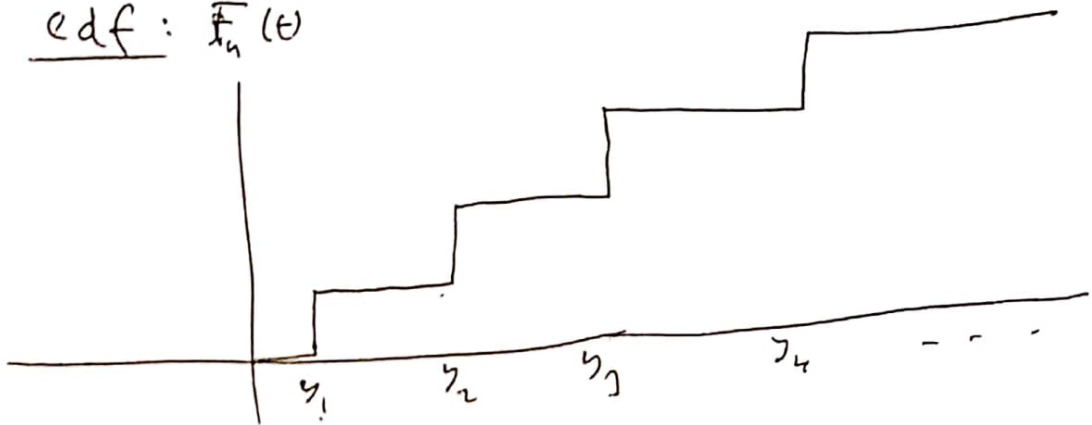
log-transformed $100(1-\alpha)\%$ confidence interval for $\hat{H}(t)$ is:

$$\left(\frac{\hat{H}(t)}{u}; \hat{H}(t)u \right);$$

$$u = \exp \left(z_{\alpha/2} \frac{\sqrt{\widehat{\text{Var}}(\hat{H}(t))}}{\hat{H}(t)} \right).$$

(7) Kernel density estimation:

edf: $F_n(\theta)$



let $y_1 < y_2 < \dots < y_k$ and

$$p(y_j) = F_n(y_j) - F_n(y_{j-1}) \quad \left(\frac{s_j}{n}\right)$$

let the kernel density estimator $K_{y_j}(x)$.
we define the kernel smoothed estimate of the cdf by:

$$\hat{F}(x) = \sum_{j=1}^k p(y_j) K_{y_j}(x)$$

and the density function:

$$\hat{f}(x) = \sum_{j=1}^k p(y_j) k_{y_j}(x)$$

(7.1)

Uniform kernel with bandwidth b :

$$k_y(x) = \begin{cases} 0 & ; \text{ elsewhere} \\ \frac{1}{2b} & ; y-b \leq x \leq y+b \end{cases}$$

is the density function of $\text{Unif}(y-b, y+b)$.

$$k_y(x) = \begin{cases} 0 & x < y-b \\ \frac{x-y+b}{2b} & y-b \leq x < y+b \\ 1 & x \geq y+b \end{cases}$$

Ex: 25, 30, 35, 37, 39, 45, 47, 49, 55.

Use a uniform kernel of bandwidth h to estimate the density function at $x=40$.

$$\begin{aligned} \hat{f}(40) &= \sum_{j=1}^9 p(y_j) k_{y_j}(40) \\ &= \frac{1}{10} k_{30}(40) + \frac{2}{10} k_{35}(40) + \frac{1}{10} k_{37}(40) \\ &\quad + \frac{1}{10} k_{39}(40) + \frac{1}{10} k_{45}(40) + \frac{1}{10} k_{47}(40) \\ &\quad + \frac{1}{10} k_{49}(40) \\ &= \frac{8}{10} \cdot \frac{1}{20} = 0.04. \end{aligned}$$

7.2

Triangular kernel with bandwidth h .

$$k_y(x) = \begin{cases} \frac{x-y+b}{b^2} & y-b \leq x \leq y \\ \frac{y+b-x}{b^2} & y \leq x \leq y+b \\ 0 & \text{elsewhere} \\ 0 & x \leq y-b \\ 0 & y-b \leq x < y \\ 0 & y \leq x < y+b \\ 1 & x \geq y+b \end{cases}$$

$$K_y(x) = \begin{cases} \frac{(x-y+b)^2}{2b^2} & y-b \leq x < y \\ 1 - \frac{(y+b-x)^2}{2b^2} & y \leq x < y+b \\ 1 & x \geq y+b \end{cases}$$

Ex: data: 25, 30, 35, 37, 39, 45, 47, 49, 55.

Use triangular kernel of bandwidth h to estimate the pdf at $x=40$.

$$\begin{aligned} \hat{f}(40) &= \frac{1}{10} k_{37}(40) + \frac{1}{10} k_{39}(40) \\ &= \frac{1}{10} \frac{1}{16} + \frac{1}{10} \frac{3}{16} = \frac{1}{40} = 0.025. \end{aligned}$$

7.3

Gamma kernel:

The Gamma kernel $k_y(x)$ is the density function of Gamma $(\alpha, \frac{x}{y})$.

$$f_y(x) = \frac{\left(\frac{x}{y}\right)^{\alpha-1} e^{-x/y}}{\Gamma(\alpha)}$$

Ex: Same data as before:

Use Gamma kernel with $\alpha = 1$ to estimate the pdf at $x = 40$.

$$\begin{aligned} \hat{f}(40) &= \frac{1}{10} \frac{1}{25} + \frac{1}{10} \frac{1}{30} + \dots \\ &= \frac{1}{10} \frac{1}{25} e^{-40/25} + \dots \\ &= 0.049 \end{aligned}$$