

as $h \rightarrow 0$, $\hat{F}_h(y) \rightarrow \frac{1}{n} \sum_{i=1}^n \mathbb{I}(Y_i \leq y) \leftarrow$ this is the empirical cdf! \Rightarrow MLE!

1.2.7 Censored Data

Censored data occur when the value is only partially known. This is different from *truncation*, in which the data does not include any values below (or above) a certain limit.

truncation:

For example, we might sample only households that have an income above a limit, L_0 . If all incomes have distribution $F(x; \theta)$, then for $y > L_0$,

$$P(Y_1 \leq y | Y_1 > L_0) = \frac{P(Y_1 \leq y, Y_1 > L_0)}{P(Y_1 > L_0)} = \frac{F(y; \theta) - F(L_0; \theta)}{1 - F(L_0; \theta)}$$

The likelihood is then

$$L(\theta | \mathcal{Y}) = \prod_{i=1}^n \left\{ \frac{f(Y_i; \theta)}{1 - F(L_0; \theta)} \right\}$$

\nearrow
this is just an iid likelihood w/ densities adjusted to take into account that $Y_i > L_0$.

1.2.7.1 Type I Censoring

Suppose a random variable X is normally distributed with mean μ and variance σ^2 , but whenever $X \leq 0$, all we observe is that it is less than or equal to 0. If the sample is set to 0 in the censored cases, then define censoring \neq truncation.

$$Y = \begin{cases} 0 & \text{if } X \leq 0 \\ X & \text{if } X > 0. \end{cases} \quad \begin{array}{l} Y: \text{observed} \\ X: \text{latent/hidden} \end{array}$$

The distribution function of Y is

$$X \sim N(\mu, \sigma^2).$$

$$\begin{aligned} \text{at } y=0: F_Y(0) &= \overset{P(Y=0)}{P(Y \leq 0)} = P(X \leq 0) \\ &= P(\sigma Z + \mu \leq 0) \\ &= P(Z \leq -\frac{\mu}{\sigma}) = \Phi\left(-\frac{\mu}{\sigma}\right). \quad \text{where } Z \sim N(0,1) \\ &\quad \Phi \text{ st. normal cdf.} \end{aligned}$$

$$\text{for } y > 0: F_Y(y) = P(Y \leq y) = \Phi\left(\frac{y - \mu}{\sigma}\right)$$

$$\text{for } y < 0: F_Y(y) = P(Y \leq y) = 0$$

Suppose we have a sample Y_1, \dots, Y_n and let n_0 be the number of sample values that are 0. Then $m = n - n_0$ and

$$\begin{aligned} L_h(\underline{\theta} | \underline{Y}) &= \left(\frac{1}{2h}\right)^m \prod_{i=1}^n \{F_Y(Y_i + h; \underline{\theta}) - F_Y(Y_i - h; \underline{\theta})\} \\ &= \left\{ \Phi\left(\frac{h - \mu}{\sigma}\right) - 0 \right\}^{n_0} \times \prod_{Y_i > 0} \left\{ \frac{\Phi((Y_i + h - \mu)/\sigma) - \Phi((Y_i - h - \mu)/\sigma)}{2h} \right\} \end{aligned}$$

$$L(\underline{\theta} | \underline{Y}) = \lim_{h \rightarrow 0^+} L_h(\underline{\theta} | \underline{Y})$$

$$= \underbrace{\left\{ \Phi\left(-\frac{\mu}{\sigma}\right) \right\}^{n_0}}_{P_\theta(Y \text{ censored})^{n_0}} \underbrace{\prod_{Y_i > 0} \left\{ \frac{1}{\sigma} \phi\left(\frac{Y_i - \mu}{\sigma}\right) \right\}}_{\prod_{Y_i \text{ not censored}} f(Y_i | \underline{\theta})}$$

essentially same as our mixture!

More generally, for left censoring @ L_0

$$L(\underline{\theta} | \underline{Y}) = \{F(L_0; \underline{\theta})\}^{n_0} \prod_{Y_i > L_0} f_Y(Y_i | \underline{\theta})$$

L_0 fixed = "type I censoring". Type II censoring = "let's observe first r failures"

We might have censoring on the left at L_0 and censoring on the right at R_0 , but observe all values of X between L_0 and R_0 . Suppose X has density $f(x; \theta)$ and distribution function $F(x; \theta)$ and

$$Y_i = \begin{cases} L_0 & \text{if } X_i \leq L_0 \\ X_i & \text{if } L_0 < X_i < R_0 \\ R_0 & \text{if } X_i \geq R_0 \end{cases}$$

If we let n_L and n_R be the number of X_i values $\leq L_0$ and $\geq R_0$ then the likelihood of the observed data Y_1, \dots, Y_n is

$$L(\underline{\theta} | \underline{Y}) = \underbrace{\{F(L_0; \underline{\theta})\}^{n_L}}_{\text{left censored part}} \underbrace{\left\{ \prod_{L_0 < Y_i < R_0} f_Y(Y_i; \underline{\theta}) \right\}}_{\text{observed}} \underbrace{\{1 - F(R_0; \underline{\theta})\}^{n_R}}_{\text{right censored part}}$$

not surprising.

We could also let each X_i be subject to its own censoring values L_i and R_i . For the special case of right censoring, define $Y_i = \min(X_i, R_i)$. In addition, define $\delta_i = \mathbb{I}(X_i \leq R_i)$. Then the likelihood can be written as

$$L(\underline{\theta}|\underline{y}) = \prod_{i=1}^n f(y_i; \underline{\theta})^{\delta_i} \underbrace{[1 - F(R_i; \underline{\theta})]^{1-\delta_i}}_{\text{right censored}}$$

$= \begin{cases} 1 & \text{if obs. } X_i \\ 0 & \text{if } X_i \text{ censored} \end{cases}$

Example (Equipment failure times): Pieces of equipment are regularly checked for failure (but started at different times). By a fixed date (when the study ended), three of the items had not failed and therefore were censored.

y	2	72	51	50	33	27	14	24	4	21
delta	1	0	1	0	1	1	1	1	1	0

$R_2 = 72 = Y_2$ $R_4 = 50 = Y_4$ $R_{10} = 21 = Y_{10}$
← censored!

Suppose failure times follow an exponential distribution $F(x; \sigma) = 1 - \exp(-x/\sigma)$, $x \geq 0$. Then

$$L(\sigma|\mathbf{Y}) = \prod_{i=1}^n \left[\frac{1}{\sigma} \exp\left(-\frac{y_i}{\sigma}\right) \right]^{\delta_i} \exp\left(-\frac{y_i}{\sigma}\right)^{1-\delta_i}$$

let n_R be # obs. were right censored.

$$= \left(\frac{1}{\sigma}\right)^{n-n_R} \exp\left(-\frac{n\bar{y}}{\sigma}\right)$$

$$\ell(\sigma) = -(n-n_R) \log \sigma - n\bar{y}/\sigma$$

$$\frac{d\ell}{d\sigma} = -\frac{(n-n_R)}{\sigma} + \frac{n\bar{y}}{\sigma^2} \stackrel{\text{set}}{=} 0$$

$$n\bar{y} = \sigma(n-n_R) \Rightarrow \hat{\sigma}_{MLE} = \frac{n\bar{y}}{n-n_R}$$

$$= \frac{10 \cdot 30.8}{7} = 44.0$$

1.2.7.2 Random Censoring

So far we have considered censoring times to be fixed. This is not required.

e.g. in medical studies patients enter at different times, modeled as random variables (w/ fixed end date)

How is this different than previous example?

↳ don't know when patients will enter!

This leads to random censoring times, e.g. R_i , where we assume that the censoring times are independent of X_1, \dots, X_n and iid with distribution function $G(t)$ and density $g(t)$.

Again $Y_i = \min(X_i, R_i)$ and $\delta_i = \mathbb{I}(X_i \leq R_i)$

Let's consider the contributions to the likelihood:

$$\begin{aligned}
 \text{due to } (Y_i, \delta_i=1): \quad & \frac{P(Y_i \in (y-h, y+h], \delta_i=1)}{2h} = \frac{P(X_i \in (y-h, y+h], X_i \leq R_i)}{2h} \\
 & = \frac{1}{2h} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \mathbb{I}(y-h < t \leq y+h, t \leq r) f(t; \theta) g(r) dt dr \\
 & = \frac{1}{2h} \int_{y-h}^{y+h} \left[\int_{-\infty}^{\infty} \mathbb{I}(t \leq r) g(r) dr \right] f(t; \theta) dt \\
 \text{FTC} \quad & = \frac{1}{2h} \int_{y-h}^{y+h} [1 - G(t)] f(t; \theta) dt \\
 & \xrightarrow{h \rightarrow 0^+} [1 - G(y)] f(y; \theta)
 \end{aligned}$$

$$\begin{aligned}
 \text{due to } (Y_i, \delta_i=0): \quad & \frac{P(Y_i \in (y-h, y+h], \delta_i=0)}{2h} = \frac{P(R_i \in (y-h, y+h], X_i > R_i)}{2h} \\
 & = \frac{1}{2h} \int_{y-h}^{y+h} \left[\int_{-\infty}^{\infty} \mathbb{I}(t > r) f(t; \theta) dt \right] g(r) dr \\
 \text{FTC} \quad & = \frac{1}{2h} \int_{y-h}^{y+h} [1 - F(r; \theta)] g(r) dr \\
 & \xrightarrow{h \rightarrow 0^+} [1 - F(y; \theta)] g(y)
 \end{aligned}$$

which results in

$$\begin{aligned}
 L(\theta | \mathbf{Y}, \delta) &= \left\{ \prod_{i=1}^n f(Y_i; \theta)^{\delta_i} [1 - G(Y_i)]^{\delta_i} \right\} \left\{ \prod_{i=1}^n g(Y_i)^{1-\delta_i} [1 - F(Y_i; \theta)]^{1-\delta_i} \right\} \\
 &= \underbrace{\prod_{i=1}^n f(Y_i; \theta)^{\delta_i} [1 - F(Y_i; \theta)]^{1-\delta_i}}_{\text{Same as for fixed!}} \underbrace{g(Y_i)^{1-\delta_i} [1 - G(Y_i)]^{\delta_i}}_{\text{not necessary for MLE}}
 \end{aligned}$$

1.3 Likelihoods for Regression Models

We will start with linear regression and then talk about more general models.

1.3.1 Linear Model

Consider the familiar linear model

$$Y_i = \mathbf{x}_i^\top \boldsymbol{\beta} + \epsilon_i, \quad i = 1, \dots, n, \quad \text{independent observations.}$$

where $\mathbf{x}_1, \dots, \mathbf{x}_n$ are known nonrandom vectors.

$$E[\epsilon_i] = 0$$

$$\text{Var}[\epsilon_i] = \sigma^2$$

often we estimate $\boldsymbol{\beta}$ by $\hat{\boldsymbol{\beta}}_{OLS}$, which does not require a distribution for ϵ_i .

For likelihood-based estimation, we need a distribution for ϵ_i ! Namely start w/ $\epsilon_i \sim N(0, \sigma^2)$.

$$\begin{aligned} \Rightarrow L(\boldsymbol{\beta}, \sigma | \{Y_i, \mathbf{x}_i\}_{i=1}^n) &= \prod_{i=1}^n \left(\frac{1}{\sqrt{2\pi}\sigma} \right) \exp\left(-\frac{(Y_i - \mathbf{x}_i^\top \boldsymbol{\beta})^2}{2\sigma^2}\right) \\ &= \left(\frac{1}{\sqrt{2\pi}\sigma} \right)^n \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (Y_i - \mathbf{x}_i^\top \boldsymbol{\beta})^2\right) \end{aligned}$$

take logs,

take derivatives,
set = 0, solve



$$\hat{\boldsymbol{\beta}}_{MLE} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$$

$$\hat{\sigma}_{MLE}^2 = \frac{1}{n} \sum_{i=1}^n (Y_i - \mathbf{x}_i^\top \hat{\boldsymbol{\beta}}_{MLE})^2$$

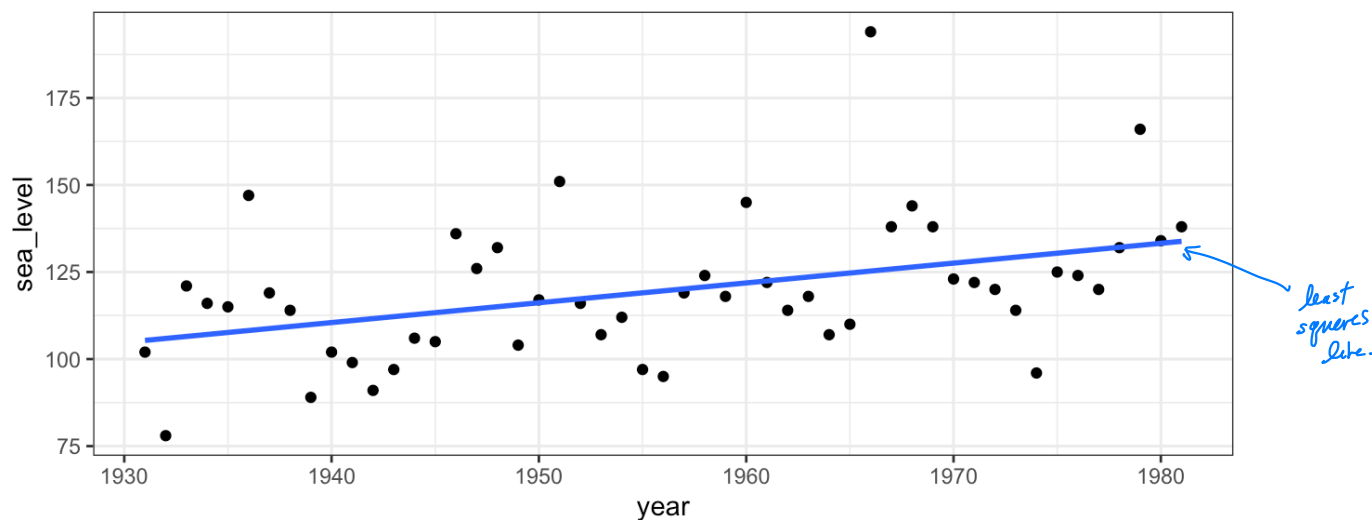
same as $\hat{\boldsymbol{\beta}}_{OLS}$

only asymptotically unbiased.

What do you do when ϵ_i are not Gaussian?

- transform data so $\hat{\epsilon}_i$ look gaussian
- Use a different distribution for ϵ_i !

Example (Venice sea levels): The annual maximum sea levels in Venice for 1931–1981 are :



We know maxima not gaussian!

Approach 1: OLS $E[\epsilon_i] = 0$, $\text{Var}[\epsilon_i] = \sigma^2$ no distributional assumption.

Approach 2: Assume $\epsilon_i \sim \text{Gumbel}$ extreme value distribution, use ML

$$f_{\epsilon}(y) = \frac{1}{\sigma} \exp\left(-\frac{y}{\sigma}\right) \exp\left(-\exp\left(-\frac{y}{\sigma}\right)\right)$$

$$\Rightarrow L(\beta_0, \sigma | \{y_i, x_i\}_{i=1}^n) = \prod_{i=1}^n f_{\epsilon}\left(y_i - x_i^T \beta\right) = \prod_{i=1}^n \frac{1}{\sigma} \exp\left(-\frac{y_i - x_i^T \beta}{\sigma}\right) \exp\left(-\exp\left(-\frac{y_i - x_i^T \beta}{\sigma}\right)\right)$$

Your Turn: Fit both approaches to Venice data!

OLS:

$$\hat{\beta}_0 = 104.8, \hat{\beta}_1 = .567 \quad \left(\begin{smallmatrix} \text{se} \\ .177 \end{smallmatrix} \right)$$

ML, GUMBEL

$$\hat{\beta}_0 = 96.8, \hat{\beta}_1 = .564 \quad \left(\begin{smallmatrix} \text{se} \\ .136 \end{smallmatrix} \right)$$

Difference in $\hat{\beta}_0$? $E[\epsilon_i] = 0.577 \sigma = 0.577 \hat{\sigma}_{\text{MLE}} \neq 0$ (14.5) $(96.8 + .577(14.5) = 105.1)$

OLS or ML? If EV model is correct, more efficient (note: st. errors).