

# Missing data lecture 9: Likelihood-based inference with incomplete data

## Likelihood inference with incomplete data

We said that the likelihood function is really a set of functions that are proportional the probability density such that the constant of proportionality doesn't depend the parameters.

Missing data methods distinguish themselves from other methods by modeling the joint distribution of  $Y$  and  $M$ . Let  $y$  and  $m$  represent the outcome measurements and  $m$  representing the missingness indicators for all  $i$  units:

$$f_{Y,M}(y, m \mid \theta, \phi) = f_Y(y \mid \theta) f_{M|Y}(m \mid y, \phi)$$

When we have missing data, we can partition the matrix  $y$  into  $y_{(1)}$  and  $y_{(0)}$ , representing the components of  $y$  that are missing and observed, respectively.

Let  $\mathcal{Y}$  be the sample space for  $y_i$  and let  $\mathcal{Y}_{(1)}$  and  $\mathcal{Y}_{(0)}$  be the sample space for the missing and observed components of  $y$ .

Then the distribution of the observed data is:

$$\int_{\mathcal{Y}_{(1)}} f_{Y,M}(y, m \mid \theta, \phi) dy_{(1)} = \int_{\mathcal{Y}_{(1)}} f_Y(y_{(0)}, y_{(1)} \mid \theta) f_{M|Y}(m \mid y_{(0)}, y_{(1)}, \phi) dy_{(1)}$$

This joint density is proportional to what we'll call the full-data likelihood:

$$L_{\text{full}}(\theta, \phi \mid y_{(0)}, m) = \int_{\mathcal{Y}_{(1)}} f_Y(y_{(0)}, y_{(1)} \mid \theta) f_{M|Y}(m \mid y_{(0)}, y_{(1)}, \phi) dy_{(1)}$$

We can also compute the likelihood *ignoring* the missingness process:

$$L_{\text{ign}}(\theta \mid y_{(0)}) = \int_{\mathcal{Y}_{(1)}} f_Y(y_{(0)}, y_{(1)} \mid \theta) dy_{(1)}$$

We'll say that the missingness mechanism is *ignorable* if inferences based on  $L_{\text{ign}}(\theta \mid y_{(0)})$  and  $L_{\text{full}}(\theta, \phi \mid y_{(0)}, m)$  are the same given  $m, y_{(0)}$ .

Formally, the missingness mechanism is ignorable for direct likelihood inference if the likelihood ratios for any two  $\theta, \theta^*$  given  $m, y_{(0)}$  are equal:

$$\frac{L_{\text{full}}(\theta, \phi \mid y_{(0)}, m)}{L_{\text{full}}(\theta^*, \phi \mid y_{(0)}, m_i)} = \frac{L_{\text{ign}}(\theta \mid y_{(0)})}{L_{\text{ign}}(\theta^* \mid y_{(0)})} \forall \theta, \theta^*, \phi$$

There are two sufficient conditions ensure ignorability:

1. Parameters  $\theta$  and  $\phi$  are *variationally independent*, i.e. the joint parameter space  $\Omega_{\theta, \phi} = \Omega_{\theta} \times \Omega_{\phi}$
2. The full likelihood factorizes as

$$L_{\text{full}}(\theta, \phi \mid y_{(0)}, m) = L_{\text{ign}}(\theta \mid y_{(0)}) L_{\text{rest}}(\phi \mid y_{(0)}, m)$$

The first condition is sufficient to ensure that the value of  $\phi$  doesn't lead to a different likelihood value for  $\theta$  vs.  $\theta^*$ .

If the data are MAR, then we will satisfy the second condition:

$$f_{M|Y}(m \mid y_{(0)}, y_{(1)}, \phi) = f_{M|Y}(m \mid y_{(0)}, y_{(1)}^*, \phi)$$

for all  $y_{(1)}, y_{(1)}^*, \phi$ . Then we can write the full-likelihood as:

$$f_{M|Y}(m \mid y_{(0)}, \phi) \int_{\mathcal{Y}_{(1)}} f_Y(y_{(0)}, y_{(1)} \mid \theta) dy_{(1)} = f_{M|Y}(m \mid y_{(0)}, \phi) f_Y(y_{(0)} \mid \theta)$$

Then by the above theorem, parameter distinctness and MAR are sufficient for ignorability.

When we do Bayesian inference we need to ensure that the posterior for  $\theta$  when using the ignorable likelihood is equal to the posterior for  $\theta$  when using the full likelihood. Under the full likelihood, the posterior for  $(\theta, \phi)$  is:

$$p(\theta, \phi \mid y_{(0)}, m) \propto p(\theta, \phi) L_{\text{full}}(\theta, \phi \mid y_{(0)}, m)$$

and under the ignorable likelihood we have:

$$p(\theta \mid y_{(0)}, m) \propto p(\theta) L_{\text{ign}}(\theta \mid y_{(0)})$$

Thus, sufficient conditions for the posteriors to be equal is that 1.  $p(\theta, \phi) = p(\theta)p(\phi)$ , or the prior independence of  $\theta$  and  $\phi$  2. the likelihood factorizes

$$L_{\text{full}}(\theta, \phi \mid y_{(0)}, m) = L_{\text{ign}}(\theta \mid y_{(0)}) L_{\text{rest}}(\phi \mid y_{(0)}, m)$$

### Example: Incomplete exponential sample

Let  $y_i \stackrel{\text{iid}}{\sim} \text{Exponential}(\theta)$  for  $i = 1, \dots, n$ . Let  $m_i$  be the missingness indicators, and suppose  $r = \sum_{i=1}^n m_i$ . The full likelihood is

$$f_Y(y \mid \theta) = \theta^{-n} \exp \left( - \sum_{i=1}^n y_i / \theta \right)$$

Let  $y_{(0)} = (y_1, \dots, y_r)$  and  $y_{(1)} = (y_{r+1}, \dots, y_n)$ . The likelihood that ignores the likelihood is

$$L_{\text{ign}}(\theta \mid y_{(0)}) = \theta^{-r} \exp \left( - \sum_{i=1}^r y_i / \theta \right)$$

Let  $m_i \stackrel{\text{iid}}{\sim} \text{Bernoulli}(\phi)$ , so

$$f_{M|Y}(m \mid y, \phi) = \phi^r (1 - \phi)^{n-r}$$

Then  $f(y_{(0)}, m \mid \theta, \phi) = \phi^r (1 - \phi)^{n-r} \theta^{-r} \exp \left( - \sum_{i=1}^r y_i / \theta \right)$ , which factorizes into a factor related to  $\theta$  and a factor related to  $\phi$ . This means we can base inferences on  $\theta$  on  $L_{\text{ign}}(\theta \mid y_{(0)})$  instead of the full likelihood. The MLE is  $\hat{\theta} = \sum_{i=1}^n y_i / r$ .

Now suppose we observe only observations for which  $y_i \leq c$ , so

$$f(m_i \mid y_i, \phi) = \mathbb{1}(y_i \geq c)^{m_i} \mathbb{1}(y_i < c)^{1-m_i}$$

Putting this together, the full likelihood is:

$$\prod_{i=1}^r f_Y(y_i \mid \theta) \mathbb{1}(y_i < c) \prod_{i=r+1}^n \int_{\mathbb{R}^+} \mathbb{1}(y_i \geq c) f(y \mid \theta) dy$$

Which of course simplifies to

$$\theta^{-r} \exp \left( - \sum_{i=1}^r y_i / \theta \right) \exp(-(n-r)c/\theta)$$

This shows that the missingness is nonignorable, because the full likelihood isn't equal to the ignorable likelihood we used in the first part of the problem.

The log-likelihood is:

$$\begin{aligned} \ell(\theta \mid y_{(0)}, m) &= -r \log \theta - \sum_{i=1}^r y_i / \theta - (n-r)c/\theta \\ \frac{\partial \ell(\theta \mid y_{(0)}, m)}{\partial \theta} &= -r/\theta + \sum_{i=1}^r y_i / \theta^2 + (n-r)c/\theta^2 \end{aligned}$$

Setting this equal to zero and solving for  $\theta$  gives the MLE:

$$\hat{\theta} = \frac{\sum_{i=1}^r y_i + (n-r)c}{r}$$

This of course doesn't equal the ignorable MLE,  $\bar{y}$  for the observed values.

## Missing data example: Parameter distinctness

Let the model be defined as

$$\begin{aligned} y_{ij} \mid \mu_i, \theta &\sim \text{Normal}(\alpha_i, \sigma^2) \\ \alpha_i \mid \theta &\sim \text{Normal}(\mu, \tau^2) \end{aligned}$$

Let the missingness mechanism be:

$$f_{M|Y}(m_{ij} \mid y, \alpha_i, \phi) = \pi(\alpha_i, \phi) = (1 + e^{-(\phi_0 + \phi_1 \alpha_i)})^{-1}$$

The joint density of the observations and parameters is:

$$\prod_{i=1}^I \prod_{j=1}^{n_i} \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(y_{ij}-\alpha_i)^2} \pi(\alpha_i, \phi)^{m_{ij}} (1 - \pi(\alpha_i, \phi))^{1-m_{ij}} \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(\alpha_i-\mu)^2}$$

Because the  $\alpha_i$  aren't observed, but do have a density, we need to integrate over them to compute the full likelihood:

$$\prod_{i=1}^I \int_{\mathbb{R}} \prod_{j=1}^{n_i} \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(y_{ij}-\alpha_i)^2} \pi(\alpha_i, \phi)^{m_{ij}} (1 - \pi(\alpha_i, \phi))^{1-m_{ij}} \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(\alpha_i-\mu)^2} d\alpha_i$$

This shows that the missingness process isn't ignorable here, even though we don't technically have the distribution of missingness depending on missing observable data, per se. This shows that in some sense,  $\alpha_i$  is missing data, and, indeed, this is what our textbook considers missing data; namely anything that has a distribution that is unobserved. This makes the problem MNAR.

Compare this to the ANOVA model with the same missingness mechanism:

$$y_{ij} \mid \alpha_i, \theta \sim \text{Normal}(\alpha_i, \sigma^2)$$

Then the joint likelihood is

$$\prod_{i=1}^I \prod_{j=1}^{n_i} \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(y_{ij}-\alpha_i)^2} \pi(\alpha_i, \phi)^{m_{ij}} (1 - \pi(\alpha_i, \phi))^{1-m_{ij}}$$

This shows that the data are MAR, because the missingness mechanism doesn't depend on missing data. However, the parameters for the missingness mechanism and the observations don't satisfy the distinctness condition, so the missingness is nonignorable.

## Partial MAR

Suppose we can partition  $\theta$  into two pieces,  $\theta_1$  and  $\theta_2$  so that the parameter of interest is  $\theta_1$ . The data are partially MAR for  $\theta_1$  if we can factorize the full likelihood:

$$L_{\text{full}}(\theta_1, \theta_2, \phi \mid y_{(0)}, m) = L_1(\theta_1 \mid y_{(0)}) L_{\text{rest}}(\theta_2, \phi \mid y_{(0)}, m)$$

### Example: Regression with missing data

An example of this is when we have covariates paired with each observation so that the complete data is  $(y_i, x_i), i = 1, \dots, n$  where  $y_i \in \mathbb{R}^d$  and  $x_i \in \mathbb{R}^p$ . let  $y_{(0)}, y_{(1)}$  be the observed and missing elements of  $y_i$  and  $x_{(0)}, x_{(1)}$  are the observed and missing elements of  $x_i$ . Let  $m_i^Z$  be the missingness indicators for the covariates,  $Z$ , and let  $m_i^Y$  be the missingness indicators for the observations  $y_i$ . Let  $m_i = (m_i^Y, m_i^Z)$  be the combined missingness indicators for unit  $i$ . Suppose that for  $i = 1, \dots, r$   $x_i$  is fully observed, while at least one component of  $y_i$  is observed, and for the remaining  $i = r + 1, \dots, n$   $y_i$  is completely missing and each  $z_i$  has at least one missing component. Let  $y_i, x_i, z_i$  be unit iid, so:

$$f_{Y,X,M}(y_i, x_i, m_i \mid \theta_1, \theta_2, \phi) = f_{Y|X}(y_i \mid x_i, \theta_1) f_X(x_i \mid \theta_2) f_{M|Y,X}(m_i \mid y_i, x_i, \phi)$$

We'll assume the missingness mechanism takes the following form:

$$f_{M|Y,X}(m_i \mid x_{i(1)}, x_{i(0)}, y_i, \phi) = f_{M|Y,X}(m_i \mid x_{i(1)}, x_{i(0)}, y_i^*, \phi)$$

for all  $y_i, y_i^*, x_{i(0)}, i = 1, \dots, n$ .

This missingness mechanism is MNAR because it depends on unobserved components of  $x_{i(1)}$ . Luckily we'll be able to factorize our likelihood so inference  $\theta_1$  is partially ignorable:

$$L_{\text{full}}(\theta_1, \theta_2, \phi \mid y_{(0)}, x_{(0)}, m) = L_{\text{p-ign}}(\theta_1 \mid y_{(0)}, x_{(0)}) L_{\text{rest}}(\theta_2, \phi \mid m, x_{(0)})$$

Let  $\mathcal{Y}_i$  be the sample space corresponding to the missing  $y_{i(1)}$ . Then we can write the ignorable part of the likelihood:

$$L_{\text{p-ign}}(\theta_1 \mid y_{(0)}, x_{(0)}) = \prod_{i=1}^r \int_{\mathcal{Y}_i} f_{Y|X}(y_{i(0)}, y_{i(1)} \mid x_i, \theta_1) dy_{i(1)}$$

Let  $\mathcal{X}_i$  be the sample space of the missing covariates for the  $i^{\text{th}}$  unit. Then the rest of the likelihood can be written as

$$L_{\text{rest}}(\theta_2, \phi \mid x_{i(0)}, m) = \prod_{i=1}^r f_X(z_i \mid \theta_2) f_{M|X}(m_i \mid x_i, \phi) \prod_{i=r+1}^n \int_{\mathcal{X}_i} f_X(x_{i(0)}, x_{i(1)} \mid \theta_2) f_{M|X}(m_i \mid x_{i(1)}, x_{i(0)}, \phi) dx_{i(1)}$$

Note the book has a typo here.

### Flawed approach to missing data

One generally flawed approach to inference in missing data problems is to treat the missing values as unknown parameters and to maximize the following function of parameters *and* missing values:

$$L_{\text{mispar}}(\theta, y_{(1)} \mid y_{(0)}) = f_Y(y_{(1)}, y_{(0)} \mid \theta)$$

The MLE using this distribution would require you to jointly maximize the likelihood with respect to  $\theta$  *and*  $y_{(1)}$ . Let's take this approach with the censored exponential samples and see what we get. We had that the likelihood was

$$\prod_{i=1}^r \theta^{-1} e^{-y_i/\theta} \mathbb{1}(y_i < c) \prod_{i=r+1}^n \theta^{-1} e^{-y_i/\theta} \mathbb{1}(y_i \geq c) f(y \mid \theta)$$

Because  $e^{-y_i/\theta}$  is monotonically decreasing in  $y_i$  for any  $\theta$ ,  $\hat{y}_i$  is  $c$  for the missing observations.

Plugging this into the log-likelihood gives

$$-r \log \theta - (n - r) \log \theta - \frac{\sum_{i=1}^r y_i + (n - r)c}{\theta}$$

Taking derivatives and setting this equal to zero gives:

$$\hat{\theta}_{\text{mispar}} = \frac{\sum_{i=1}^r y_i + (n - r)c}{n} = \frac{r}{n} \hat{\theta}$$

This understates the true value  $\theta$ , and one can show that this estimator isn't consistent for  $\theta$ . The only way this estimator is consistent for  $\theta$  is if  $r/n \rightarrow 1$ .

This example shows that the goal in missing data analysis isn't to predict missing values, it is to account for the uncertainty in missing values by integrating over the distribution of missing values.

## Coarsened data

Coarsened data is a generalization of missing data that includes other ways in which the resolution of data can be reduced. Examples include censoring, grouping, rounding, or heaping. Heaping is the phenomenon where there are varying levels of resolution reported in the same dataset. For example, on a questionnaire that asks for the the number of cigarettes smoked per day, some people will report exact numbers, and others will report multiples of packs. With rounded data, the coarsening is more deterministic, namely we know that an observation is exactly within the interval, say between  $[\text{floor}(y), \text{floor}(y) + 1]$ . With coarsened data, there is still the complete data matrix  $y = (y_{ij})$ , but there is now a coarsening variable  $c_{ij}$  that interacts with the true value to return the observed data.

Let  $y_{ij(0)}$  be the observed data, and let  $h_{ij}(y_{ij}, c_{ij})$  be the function of the true value and the coarsening variable that returns some subset of  $\mathcal{Y}_{ij}$  to which  $y_{ij}$  belongs. Thus  $y_{ij(0)} = h_{ij}(y_{ij}, c_{ij})$  with the requirement that  $y_{ij} \in h_{ij}(y_{ij}, c_{ij})$ . Let  $c_{ij(0)}$  be the observed part of the coarsening random variable, that is governed by a function  $g_{ij}(y_{ij}, c_{ij})$ .

The simplest nontrivial example is the censored exponential data from above, though we will modify the scenario so that each individual has a potentially different censoring time  $c_i$ . Let  $y_i$  be the true time to failure, while  $c_i$  is the censoring time.

$$y_{i(0)} = h(y_i, c_i) = \begin{cases} y_i & y_i \leq c_i \\ (c_i, \infty) & y_i > c_i \end{cases}$$

$$c_{i(0)} = g(y_i, c_i) = \begin{cases} (y_i, \infty) & y_i \leq c_i \\ c_i & y_i > c_i \end{cases}$$

Let  $c_i = c_{(0)}, c_{(1)}$  be the vector of coarsening values, where  $c_{(0)}$  is the observed value of the coarsening  $c_i$  is observed exactly, while  $c_{(1)}$  is the missing value of the coarsening. Let  $y_{(0)}$  be the set of values that we observe exactly, and  $y_{(1)}$  be the set of values that are censored. Furthermore, let the distribution of interest for  $y_i$  be  $f_Y(y_i | \theta)$ , while we let the coarsening distributuion be  $f_{C|Y}(c_i | y_i, \phi)$ . Then we can write:

$$L_{\text{full}}(\theta, \phi | y_{(0)}, c_{(0)}) = \prod_{i=1}^n \int \int f_{C|Y}(c | y, \phi) f(y | \theta) \mathbb{1}(y \in g(y_i, c_i)) \mathbb{1}(c \in h(y_i, c_i)) dy dc$$

This leads to a definition of coarsening at random, or CAR, that relates to conditions on the coarsening distribution:

$$f_{C|Y}(c_{(0)}, c_{(1)} | y_{(0)}, y_{(1)}, \phi) = f_{C|Y}(c_{(0)}, c_{(1)}^* | y_{(0)}, y_{(1)}^*, \phi)$$

For all  $c_{(1)}, c_{(1)}^*, y_{(1)}, y_{(1)}^*, \phi$ .

In the failure time example we have two contributions to the likelihood:

$$L_{\text{full}}(\theta, \phi | y_{(0)}, c_{(0)}) = \prod_{i|y_i \leq c_i} f(y_i | \theta) \int_{y_i}^{\infty} f_{C|Y}(c_i | y_i, \phi) dc \times \prod_{i|y_i > c_i} \int_{c_i}^{\infty} f_{C|Y}(c_i | y, \phi) f(y | \theta) dy$$

If we have that  $f(c_i | y_i, \phi) = f(c_i | \phi)$  for all  $i$ , and that  $\phi$  and  $\theta$  are variationally independent, we can write the likelihood as the product of  $L_{\text{ign}}(\theta | y_{(0)}, c_{(0)})$  and  $L_{\text{rest}}(\phi | y_{(0)}, c_{(0)})$

$$\prod_{i|y_i \leq c_i} f(y_i | \theta) \prod_{i|y_i > c_i} \int_{c_i}^{\infty} f(y | \theta) dy \times \prod_{i|y_i \leq c_i} \int_{y_i}^{\infty} f_C(c_i | \phi) dc \prod_{i|y_i > c_i} f(c_i | \phi)$$

In this case, the censoring mechanism is CAR, but not MAR, as we saw earlier.

In the cigarette smoking example, let  $y_i$  be the true number of cigarettes smoked per day, and let  $c_i$  be an indicator for the precision of reporting. Then define  $y_{ij}$  to be:

$$y_{ij(0)} = \begin{cases} [\text{floor}(y_i), \text{floor}(y_i) + 1] & c_i = 0 \\ [20\text{floor}(y_i/20), 20\text{floor}(y_i/20) + 20] & c_i = 1 \end{cases}$$

The main problem is that for  $y_i \leq 20$  we know that  $c_i = 0$ , but for  $y_i \geq 20$  we know that  $c_i \in \{0, 1\}$ .

Suppose that  $f_{C|Y}(c_i | y_i, \phi) = \Phi(\phi_1 + \phi_2 y_i)^{c_i} (1 - \Phi(\phi_1 + \phi_2 y_i))^{1-c_i}$

## Clarifying the distinction between MAR and MAAR

This previous discussion helps clarify the difference between MAR and MAAR.

Suppose we have measurements for  $y_{i1}, y_{i2}$  and the following model is assumed to hold for our missingness patterns:

$$\begin{aligned} f_{M|Y}(m_{i1} = 1, m_{i2} = 1 \mid y_i, \phi) &= g(y_i, \phi) \\ f_{M|Y}(m_{i1} = r, m_{i2} = s \mid y_{i1}, y_{i2}, (\{m_{i1} = 1\} \cap \{m_{i2} = 1\})^c, \phi) &= h(\pi) \\ (r, s) &\in \{(1, 0), (0, 1), (0, 0)\} \end{aligned}$$

The likelihood for an observed dataset is as follows:

$$\prod_i \left( \int_{\mathcal{Y}} f_Y(y_i \mid \theta) g(y_i, \phi) \right)^{\mathbb{1}(m_{i1}=1, m_{i2}=1)} \left( h(\pi) \int_{\mathcal{Y}_{(1)i}} f_Y(y_i \mid \theta) dy_{(1)i} \right)^{1-\mathbb{1}(m_{i1}=1, m_{i2}=1)}$$