

We formalize the concepts of missing-data mechanisms by following the notation in Little and Rubin (2002).

## Notation in missing data mechanisms

- Data:  $\mathbf{Y} = (y_{ij})_{n \times K}$  where  $i$ th row  $\mathbf{y}_i = (y_{i1}, \dots, y_{iK})$ ,  $y_{ij}$  is the value of variable  $\mathbf{Y}_j$  for observation  $i$ , and  $\mathbf{Y}$  is characterized by  $\theta$
- Missing data indicator matrix:  $\mathbf{M} = (m_{ij})$ ,  $m_{ij} = 1$  if  $y_{ij}$  is missing and  $m_{ij} = 0$  if  $y_{ij}$  is present; and  $\mathbf{M}$  is characterized by unknown parameter  $\psi$ .
- The missing-data mechanism is characterized in terms of the conditional distribution of  $\mathbf{M}$  given  $\mathbf{Y}$ ,  $f(\mathbf{M}|\mathbf{Y}, \psi)$ .

- **Missing completely at random (MCAR)**

This means there is no relationship between the missingness of the data and any values, observed or missing. Those missing data points are a random subset of the data.

$$f(\mathbf{M}|\mathbf{Y}, \psi) = f(\mathbf{M}|\psi) \text{ for all } \mathbf{Y}, \psi.$$

- **Missing at random (MAR)**

Let  $\mathbf{Y} = (\mathbf{Y}_{obs}, \mathbf{Y}_{mis})$  where  $\mathbf{Y}_{obs}$  is defined as the observed component of  $\mathbf{Y}$  and  $\mathbf{Y}_{mis}$  is the missing component. When the missingness depends on  $\mathbf{Y}_{obs}$ , the missing-data mechanism is said to be MAR if

$$f(\mathbf{M}|\mathbf{Y}, \psi) = f(\mathbf{M}|\mathbf{Y}_{obs}, \psi) \text{ for all } \mathbf{Y}_{mis}, \psi.$$

Whether an observation is missing has nothing to do with the missing values, but it does have to do with the values of observed variables.

- **Missing not at random (MNAR)**

Under MNAR, the distribution of  $\mathbf{M}$  depends the data  $\mathbf{Y}$ . This means there is a relationship between the propensity of a value to be missing and its values. For instance, people with high incomes tend to not providing their salaries in surveys.

- **Complete case analysis:** the complete-case analysis only analyzes the complete observations in a dataset. It is easy to implement but only valid under MCAR.
- **Single imputation**
  - **Mean imputation:** mean calculated based on the complete observations in the dataset.
  - **Regression imputation:** predictions from the regression model built on complete observations.
  - **Stochastic regression imputation:** predictions from the regression model **plus** a random draw from the estimated distribution of residuals.
  - **Logistic regression imputation:** predictions from the logistic regression model built on complete observations.
  - **“Worst-rank” method:** Lachin (1999) proposed “worst-rank” analysis by assigning more extreme values (values indicating “worst” treatment effects) than observed values as the imputed values for missing data. All missing values share the same values (ranks) if a worst-rank analysis applies.
  - **“Best-worst and worst-best” method:** suppose the treatment  $A$  is the beneficial group and  $B$  is the placebo. If the “best - worst” method is adopted, the imputed values for missing responses in treatment  $A$  will be values representing harmful outcomes (i.e., the “worst” values among observed values). While in treatment  $B$ , the imputed values will be values representing beneficial outcomes (i.e., the “best” values among observed values).

- **Direct maximum likelihood method** is a method handling missing data without imputing missing values, and it is valid under MAR assumption.
- Under MAR, we have

$$f(\mathbf{Y}_{obs}, \mathbf{M}|\boldsymbol{\theta}, \boldsymbol{\psi}) = f(\mathbf{M}|\mathbf{Y}_{obs}, \boldsymbol{\psi})f(\mathbf{Y}_{obs}|\boldsymbol{\theta}),$$

where we are interested in the inference for  $\boldsymbol{\theta}$ . If  $\boldsymbol{\theta}, \boldsymbol{\psi}$  are distinct, as the distribution  $f(\mathbf{M}|\mathbf{Y}_{obs}, \boldsymbol{\psi})$  does not depend on  $\boldsymbol{\theta}$ , to estimate the ML of  $\boldsymbol{\theta}$  is equivalent to maximize the likelihood  $f(\mathbf{Y}_{obs}|\boldsymbol{\theta})$ ; i.e, to ignore the missing data.

- Consider the linear regression model with the following format  $\mathbf{y} = \tilde{\mathbf{X}}^\top \boldsymbol{\beta} + \boldsymbol{\epsilon}$
- $\tilde{\mathbf{X}} = (1, \mathbf{X}^\top)^\top$  and the predictors  $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_p)^\top \sim N_p(\boldsymbol{\mu}_X, \boldsymbol{\Sigma}_X)$
- $\epsilon_i \sim N(0, \sigma^2)$ .
- Some values are missing within  $n$  independent observations  $(y_i, \mathbf{X}_i^\top)^\top$ ,  $i = 1, \dots, n$ .
- Based on the normality and independence assumptions of  $\mathbf{y}$  and  $\mathbf{X}$ , we have

$$(\mathbf{y}, \mathbf{X}) \sim N(\boldsymbol{\mu}_{\mathbf{y}, \mathbf{X}}, \boldsymbol{\Sigma}_{\mathbf{y}, \mathbf{X}})$$

where  $\boldsymbol{\mu}_{\mathbf{y}, \mathbf{X}} = \begin{pmatrix} \boldsymbol{\mu}_y \\ \boldsymbol{\mu}_X \end{pmatrix}$  and  $\boldsymbol{\Sigma}_{\mathbf{y}, \mathbf{X}} = \begin{pmatrix} \boldsymbol{\Sigma}_y & \boldsymbol{\Sigma}_{y, \mathbf{X}} \\ \boldsymbol{\Sigma}_{\mathbf{X}, y} & \boldsymbol{\Sigma}_X \end{pmatrix}$ .

- The parameters  $\boldsymbol{\theta} = (\boldsymbol{\mu}_{\mathbf{y}, \mathbf{X}}, \boldsymbol{\Sigma}_{\mathbf{y}, \mathbf{X}})$  are estimated through the expectation maximization algorithm (EM) algorithm.

## Steps of EM algorithm

- **Initiate**  $\theta^{(0)}$ ;  $\theta^{(t)}$  is the estimate of  $\theta$  at the  $t$ th iteration.
- **E step:** compute the expectation of complete-data log-likelihood with respect to the conditional distribution of  $\mathbf{Y}_{mis} | \mathbf{Y}_{obs}$  with  $\theta^{(t)}$ , i.e.:

$$Q(\theta | \theta^{(t)}) = E[l(\mathbf{Y}; \theta) | \mathbf{Y}_{obs}; \theta^{(t)}] = \int l(\mathbf{Y}; \theta) f(\mathbf{Y}_{mis} | \mathbf{Y}_{obs}; \theta^{(t)}) d\mathbf{Y}_{mis}.$$

- **M step:** maximize the  $Q$  function to obtain  $\theta^{(t+1)}$  :

$$\theta^{(t+1)} = \arg \max_{\theta} Q(\theta | \theta^{(t)}).$$

- Iterate between E step and M step until the change in function  $Q$  is very small.

After the  $\theta$  has been estimated by EM algorithm, we can have the estimate of coefficient  $\beta = (\mu_y - \Sigma_{y,x} \Sigma_x^{-1} \mu_x, \Sigma_{y,x} \Sigma_x^{-1})^\top$ . Also, the standard deviations can be estimated by  $\mathbb{V}[\beta] = \text{diag}(\mathbf{C})$ , with

$$\mathbf{C} = (\Sigma_y - \beta^\top \Sigma_x \beta) \left( (0_{p+1}, (0_p, \Sigma_x)^\top)^\top + (1, \mu_x^\top)^\top (1, \mu_x^\top) \right)^{-1} / n.$$

**Multiple imputation:** an alternative method for dealing with missing data under MAR.

## Steps of multiple imputation:

- **Imputation step:**  $m$  imputations are conducted. Hence,  $m$  completed datasets are generated by replacing the missing values with the imputed values  $m$  times. Usually, setting  $m = 50$  or higher is acceptable to reduce the sampling uncertainty from the imputation process.
- **The complete-data analysis step:** a desirable statistical analysis is conducted individually on each complete dataset generated from the previous step.
- **Pooling step:** collect  $m$  statistical inference results (e.g., parameter estimates and their standard errors) from the previous step. Based on the Rubin (2004)'s rules, overall parameter estimates and their standard errors, confidence intervals and  $p$ -values can be generated by combining separate results.