

Missing values and imputation

data matrix D
 $n \times p$

rows: T_1, T_2, \dots, T_n i.i.d. $\sim \mathcal{N}_p(\mu, \Sigma)$

$$\theta = (\mu, \Sigma)$$

for every i : $T_i = (Y_i, Z_i)$

observed $d_1 \times 1$ missing $d_2 \times 1$

$$d_1 + d_2 = p$$

decompose for every i :

$$\mu_i = (\mu^{obs, i}, \mu^{mis, i})$$

$$\Sigma_i = \begin{pmatrix} \sum_{d_1 \times d_1}^{obs, i} & \sum_{d_1 \times d_2}^{cross, i} \\ \sum_{d_2 \times d_1}^{cross, i} & \sum_{d_2 \times d_2}^{mis, i} \end{pmatrix}$$

- log-likelihood:

$$- \sum_{i=1}^n \log f_{\theta}(Y_i)$$

complicated expression
depending on missingness pattern of i th row of D

complete likelihood is explicit

→ EM-algorithm is "suitable"

sketch of EM-algorithm:

1) start with $\hat{\mu}^{(0)}$, $\hat{\Sigma}^{(0)}$

2) E-step:

$$E[Z_i | Y_i, \hat{\mu}^{(m)}, \hat{\Sigma}^{(m)}]$$

$$\begin{matrix} \uparrow \\ \hat{\mu}^{(m)} \\ \text{Gaussian} \\ \text{assumption} \\ \Rightarrow \hat{\Sigma}_i \end{matrix} \quad \hat{\mu}_{\text{mis},i}^{(m)} + \left(\hat{\Sigma}_{\text{cross},i}^{(m)} \right)^T \left(\hat{\Sigma}_{\text{obs},i}^{(m)} \right)^{-1} \left(Y_i - \hat{\mu}_{\text{obs},i}^{(m)} \right)$$

Gaussian
assumption

$$\Rightarrow \hat{\Sigma}_i$$

$$3) \text{ M-step: } \hat{\mu}^{(m+1)} = E \left[\frac{1}{n} \sum_{i=1}^n T_i \mid Y, \hat{\theta}^{(m)} \right]$$

$$= \left(\frac{1}{n} \sum_{i=1}^n Y_i, \frac{1}{n} \sum_{i=1}^n \hat{\Sigma}_i \right)$$

(up to re-ordering)

$$\hat{\Sigma}^{(m+1)} = E \left[\frac{1}{n} \sum_{i=1}^n T_i T_i^T - \hat{\mu}^{(m+1)} \left(\hat{\mu}^{(m+1)} \right)^T \mid Y, \hat{\theta}^{(m)} \right]$$

→ have to compute

$$E[Y_i Y_i^T | Y, \hat{\theta}^{(m)}] = Y_i Y_i^T$$

$$E[Y_i Z_i^T | Y, \hat{\theta}^{(m)}] = Y_i \hat{Z}_i^T$$

$$E[Z_i Z_i^T | Y, \hat{\theta}^{(m)}] = \hat{C}_i + \hat{Z}_i \hat{Z}_i^T$$

$$\hat{C}_i = \hat{\Sigma}_{\text{mis},i}^{(m)} - (\hat{\Sigma}_{\text{cross},i}^{(m)})^T (\hat{\Sigma}_{\text{obs},i}^{(m)})^{-1} \hat{\Sigma}_{\text{cross},i}^{(m)}$$

R-package: `mvrnmle`