

MISSING DATA

①

Ex: PIMA INDIANS

; PRECISION MEDICINE

We record:

- glu : blood glucose concentration (mg/deciliter)
- bp : diastolic blood pressure (bottom #; pressure in arteries when heart rests btwn beats)
- skin : skin fold thickness (mm)
- bmi : body mass index (kg/m²)

Questions : - How do these measurements compare (in PIMA pop'n) to national average?
- How do these measurements relate (covary) w/in the PIMA pop'n?
- What kind of model could help answer these questions?

Proposed Model :
$$Y_i = \begin{bmatrix} Y_{i1} \\ Y_{i2} \\ Y_{i3} \\ Y_{i4} \end{bmatrix} \sim \text{MVN} \left(\begin{matrix} \underline{\theta} \\ \underline{\Sigma} \end{matrix}, \begin{matrix} 4 \times 1 \\ 4 \times 4 \end{matrix} \right)$$

4x1 vector of measurements for ith individual

Complication: some data are missing.

How can we handle this?

- throw out missing data? No! Lose a lot of info.
- impute (w/ mean, w/ nearest neighbor)
 - issue: we have uncertainty about the value & throw that out.
 - another issue: we might hurt cov. structure estimation.

To handle this in a principled Bayesian way, we need to account for the missingness. ②

Let $\underline{\theta}_i = \begin{bmatrix} \theta_{i1} \\ \vdots \\ \theta_{i4} \end{bmatrix}$ be an observation indicator vector.

$\theta_{ij} = 1$ if Y_{ij} obs.

$\theta_{ij} = 0$ if Y_{ij} missing

Let $\underline{Y} = \underline{Y}_1, \dots, \underline{Y}_n$

$\underline{\theta} = \underline{\theta}_1, \dots, \underline{\theta}_n$

This notation lets us write the COMPLETE DATA Likelihood:

$$P(\underline{Y}, \underline{\theta} \mid \underline{\theta}, \underline{\Sigma}, \underline{\phi})$$

Let $\underline{Y} = [\underline{Y}_{\text{obs}}, \underline{Y}_{\text{mis}}]$

$$\underline{Y}_{\text{obs}} = \underline{Y} [\underline{\theta} == 1]$$

$$\underline{Y}_{\text{mis}} = \underline{Y} [\underline{\theta} == 0]$$

observed data likelihood:

$$P(\underline{Y}_{\text{obs}}, \underline{\theta} \mid \underline{\theta}, \underline{\Sigma}, \underline{\phi}) = \int P(\underline{Y}_{\text{obs}}, \underline{Y}_{\text{mis}}, \underline{\theta} \mid \underline{\theta}, \underline{\Sigma}, \underline{\phi}) d\underline{Y}_{\text{mis}}$$

3 WAYS TO CONSIDER MISSING DATA MECHANISM: ③

1) MCAR : Missing Completely At RANDOM
 - no dependence on obs or missing data whatsoever.

2) MAR : Missing at Random
 - missingness only depends on obs data (but not missing data).
 i.e. $\underline{\theta}_i \perp \underline{Y}_{mis,i} \mid \underline{Y}_{obs,i}$

3) MNAR : missing not at random
 - missingness may depend on missing data.

Assumption we'll make: Data are MAR

Implication:

complete data likelihood:

$$p(\underline{Y}, \underline{\theta} \mid \underline{\theta}, \underline{\Sigma}, \underline{\phi}) = \underbrace{p(\underline{\theta} \mid \underline{Y}_{obs}, \underline{Y}_{mis}, \underline{\theta}, \underline{\Sigma}, \underline{\phi})}_{p(\underline{\theta} \mid \underline{Y}_{obs}, \underline{\phi})} p(\underline{Y} \mid \underline{\theta}, \underline{\Sigma}, \underline{\phi})$$

We are interested in, as Bayesians,

$p(\text{unknowns} \mid \text{knowns})$

$$p(\underline{\theta}, \underline{\Sigma}, \underline{Y}_{mis} \mid \underline{Y}_{obs}, \underline{\phi}) \propto p(\underline{\theta}, \underline{\Sigma}, \underline{Y}_{obs}, \underline{Y}_{mis}, \underline{\phi})$$

$$\propto \underbrace{p(\underline{Y}_{mis}, \underline{Y}_{obs} \mid \underline{\theta}, \underline{\Sigma}, \underline{\phi})}_{\text{complete data likelihood}} \cdot \underbrace{p(\underline{\theta}, \underline{\Sigma})}_{p(\underline{\theta}, \underline{\Sigma} \mid \underline{\phi})} \cdot p(\underline{\phi})$$

$$\text{complete data likelihood} \quad p(\underline{\theta}, \underline{\Sigma} \mid \underline{\phi}) \cdot p(\underline{\phi})$$

I want to approximate this posterior.

What priors could enable Gibbs sampling?

$$\underline{\theta} \sim \text{MVN}(\underline{\mu}_0, \underline{V}_0)$$

$$\underline{\Sigma} \sim \text{inv-Wishart}(\underline{n}_0, \underline{S}_0)$$

assumption:

$$p(\underline{\theta}, \underline{\Sigma} | \underline{y}) = p(\underline{\theta}) p(\underline{\Sigma})$$

Gibbs sampler:

$$p(\underline{\theta} | \cdot) = \text{dMVN}(\underline{\mu}_n, \underline{V}_n)$$

$$p(\underline{\Sigma} | \cdot) = \text{dinv-Wishart}(\underline{n}_n, \underline{S}_n)$$

↓ ↓
 functions of complete data \underline{y} .
 ↑ ↑

$$p(\underline{y}_{\text{mis}} | \cdot) = \text{conditional Normal}$$

↓
 can see this easily from
 data generative model or
 write out

$$p(\underline{y}_{\text{mis}} | \cdot) \propto \underbrace{p(\underline{y}_{\text{mis}}, \underline{y}_{\text{obs}} | \underline{\theta}, \underline{\Sigma}, \underline{\theta})}_{\text{MVN}}$$