

MISSING DATA

Ex PIMA INDIANS ; PRECISION MEDICINE

We record:

glu : blood glucose concentration

bp : diastolic blood pressure

skin : skin fold thickness

bmi : body mass index

Questions : How do these measurements in the PIMA pop'n compare to national averages?
How do these measurements covary in PIMA pop'n?

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

↑
vector of
measurements
for ith indiv.

Complication : some data are missing

How to handle?

- throw out missing data? No! Lose a lot of info.
- impute w/ mean of column? No! Lose all cov. structure.

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

Let $O_i = \begin{bmatrix} O_{i1} \\ \vdots \\ O_{in} \end{bmatrix}$ be an observation indicator
 vector.

$$O_{ij} = 1 \quad \text{if } Y_{ij} \text{ obs.} \\
 = 0 \quad \text{if } Y_{ij} \text{ missing}$$

COMPLETE DATA LIKELIHOOD: Let $Y = Y_1, \dots, Y_n$
 $O = O_1, \dots, O_n$

$$P(Y, O | \theta, \Sigma, \phi) \\
 = P(O | Y, \theta, \Sigma, \phi) \cdot P(Y | \theta, \Sigma)$$

$$\text{Let } Y = [Y_{\text{obs}}, Y_{\text{mis}}] \\
 Y_{\text{obs}} = Y[O = 1] \\
 Y_{\text{mis}} = Y[O = 0]$$

obs. matrix

OBSERVED DATA LIKELIHOOD

$$P(Y_{\text{obs}}, O | \theta, \Sigma, \phi) = \int P(Y_{\text{obs}}, Y_{\text{mis}}, O | \theta, \Sigma, \phi) dY_{\text{mis}}$$

ASSUMPTION : DATA are MAR
 "missing at random"

$$p(O|y, \theta, \Sigma, \phi) = p(O|\phi)$$

where ϕ does not depend on θ, Σ ,
 or y_{mis} .

$$O \perp y_{\text{mis}}$$

MAR MAR missingness does not depend
 on missing data whatsoever, but may
 depend on observed data.

MCAR : "missing completely at random"
 $\Rightarrow p(O|\phi) \& \phi$ does not even depend
 on y_{obs} .

missingness has
no dependency on data.

MNAR : missing NOT at random
 missingness depends on missing data.

We are interested in, as Bayesians ⁽⁴⁾

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

$$p(\Theta, \Sigma, Y_{\text{mis}} | O, Y_{\text{obs}}) \propto$$

$$p(\Theta, \Sigma, Y_{\text{mis}}, Y_{\text{obs}}, O) \propto$$

$$\underbrace{p(Y_{\text{mis}}, Y_{\text{obs}} | \Theta, \Sigma, O)}_{\text{complete data likelihood}} \cdot \underbrace{\frac{p(O, \Theta, \Sigma)}{p(\Theta, \Sigma | O) p(O | \Phi)}}_{\text{prior}}$$

I want to approx. this posterior.
What priors would enable Gibbs sampling?

$$\Theta \sim \text{MVN}(\mu_0, \Sigma_0)$$

$$\Sigma \sim \text{inv-Wishart}(\nu_0, S_0)$$

assumption:

$$p(\Theta, \Sigma | O) = p(\Theta) p(\Sigma)$$

Gibbs sampling proceeds for each unknown:

$$p(\theta | \cdot) = \text{LMVN}(\mu_n, \Sigma_n)$$

$\downarrow \quad \downarrow$
 function of complete data y_{mis} & y_{obs} and μ_0, Σ_0

$$p(\Sigma | \cdot) = \text{inv-wishart}(\nu_n, S_n)$$

$\downarrow \quad \downarrow$
 function of complete data & ν_0, S_0

" " means everything else

$$p(y_{\text{mis}} | \cdot) \propto \frac{p(y_{\text{obs}}, y_{\text{mis}} | \theta, \Sigma, \mu)}{p(y_{\text{mis}} | y_{\text{obs}}, \theta, \Sigma, \mu)} \cdot \cancel{p(y_{\text{obs}} | \cdot)}$$

$\underbrace{\hspace{10em}}$
 cond'l normal