

Missing Data and Imputation

Rebecca C. Steorts

Bayesian Methods and Modern Statistics: STA 360/601

Lecture

- ▶ Missing data for real applications
- ▶ Common approaches
- ▶ Bayesian approach
- ▶ Lab: Example

Health related measurements on women from Pima Indian heritage village

- ▶ glu: glucose concentration
- ▶ bp: diastolic blood pressure
- ▶ skin: skin fold thickness
- ▶ bmi: body mass index

	glu	bp	skin	bmi
1	86	68	28	30.2
2	195	70	33	NA
3	77	82	NA	35.8
4	NA	76	43	47.9
5	107	60	NA	NA
6	97	76	27	NA

How would we do parameter estimate with missing values?

Remove the missing values (throwing away data).

Impute these with the population mean (statistically incorrect).

Notation and Missing at Random

- ▶ \mathbf{Y} are observed covariates.
- ▶ $\mathbf{O}_i = (O_1, \dots, O_p)^T$
- ▶ $O_i = 1$ implies Y_{ij} is not missing, 0 otherwise.
- ▶ We assume data is missing data random, meaning that \mathbf{O}_i are \mathbf{Y}_i are statistically independent
- ▶ We assume also that \mathbf{O}_i does not depend on $\boldsymbol{\theta}$ or Σ .
- ▶ For when it's not missing at random, see Gelman and Rubin, Chapter 21.

Missing at Random

$$p(\{y_{ij} : o_{ij} = 1\} \mid \boldsymbol{\theta}, \Sigma) \quad (1)$$

$$= p(\mathbf{o}_i) \times p(\{y_{ij} : o_{ij} = 1\} \mid \boldsymbol{\theta}, \Sigma) \quad (2)$$

$$= p(\mathbf{o}_i) \times \int \left\{ p(y_{i,1}, \dots, y_{i,p} \mid \boldsymbol{\theta}, \Sigma) \prod_{y_{ij}: o_{ij}=0} dy_{ij} \right\} \quad (3)$$

Main point: integrate out all the missing y's.

Simple Example

Let $\mathbf{y}_i = (y_{i1}, \text{NA}, y_{i3}, \text{NA})^T$.

Then $\mathbf{o}_i = (1, 0, 1, 0)^T$.

$$p(\{y_{ij} : o_{ij} = 1\} \mid \boldsymbol{\theta}, \Sigma) \tag{4}$$

$$= p(\mathbf{o}_i, y_{i1}, y_{i3} \mid \boldsymbol{\theta}, \Sigma) \tag{5}$$

$$= p(\mathbf{o}_i) \times \int \{p(\mathbf{y}_i \mid \boldsymbol{\theta}, \Sigma) dy_2 dy_4\} \tag{6}$$

Missing data with Gibbs Sampling

- ▶ Let \mathbf{Y} be the observed data.
- ▶ Let \mathbf{Y}_{obs} be the data we do observe (not missing).
- ▶ Let \mathbf{Y}_{miss} be the data we do not observe (missing).

For any observed data, we want to estimate

$$p(\boldsymbol{\theta}, \Sigma, \mathbf{Y}_{miss} \mid \mathbf{Y}_{obs}).$$

We do this via Gibbs sampling.

Suppose starting values $\Sigma^{(o)}, \mathbf{Y}_{miss}^{(o)}$.

We generate

$$\boldsymbol{\theta}^{(s+1)}, \Sigma^{(s+1)}, \mathbf{Y}_{miss}^{(s+1)}$$

and

$$\boldsymbol{\theta}^{(s)}, \Sigma^{(s)}, \mathbf{Y}_{miss}^{(s)}$$

by

1. Sampling $\boldsymbol{\theta}^{(s)}$ from $p(\boldsymbol{\theta} \mid \mathbf{Y}_{obs}, \mathbf{Y}_{miss}^{(s)}, \Sigma^{(s)})$
2. Sampling $\Sigma^{(s+1)}$ from $p(\Sigma \mid \mathbf{Y}_{obs}, \mathbf{Y}_{miss}^{(s)}, \boldsymbol{\theta}^{(s+1)})$
3. Sampling $\mathbf{Y}_{miss}^{(s+1)}$ from $p(\mathbf{Y}_{miss} \mid \mathbf{Y}_{obs}, \boldsymbol{\theta}^{(s+1)}, \Sigma^{(s)})$

Using steps 1 and 2, we obtain a full matrix \mathbf{Y} .

Then $\boldsymbol{\theta} \mid \mathbf{Y}, \Sigma \sim \text{MVN}(\mu_n, \Sigma_n)$

(Slide 10, MVN and Wishart Lecture, Hoff eqn 7.6)

Also, $\Sigma \mid \mathbf{Y}, \boldsymbol{\theta} \sim \text{inverseWishart}(\nu_o + n, S_n)$

(Slide 15, MVN and Wishart Lecture, Hoff eqn 7.9).

Consider

$$p(\mathbf{Y}_{\text{miss}} \mid \mathbf{Y}_{\text{obs}}, \boldsymbol{\theta}, \Sigma) \propto p(\mathbf{Y}_{\text{miss}}, \mathbf{Y}_{\text{obs}} \mid \boldsymbol{\theta}, \Sigma) \quad (7)$$

$$\propto \prod_i p(\mathbf{y}_{i,\text{miss}}, \mathbf{y}_{i,\text{obs}} \mid \boldsymbol{\theta}, \Sigma) \quad (8)$$

$$\propto \prod_i p(\mathbf{y}_{i,\text{miss}} \mid \mathbf{y}_{i,\text{obs}}, \boldsymbol{\theta}, \Sigma) \quad (9)$$

How can we compute this? There's clever little result (eqn 7.11, Hoff) that makes it possible. (It's just another MVN result).

Illustration from Hoff's example, p. 120

The prior mean is set at $(120, 64, 26, 26)^T$

We use Hoff's code given in the book and run 1,000 iterations of the GS.

Posterior mean is $(123.4671.0329.3532.18)^T$

We convert the covariance matrix into a correlation matrix (see Hoff for formula).

Then we look at the marginal posterior of 95 percent quantile based confidence intervals.

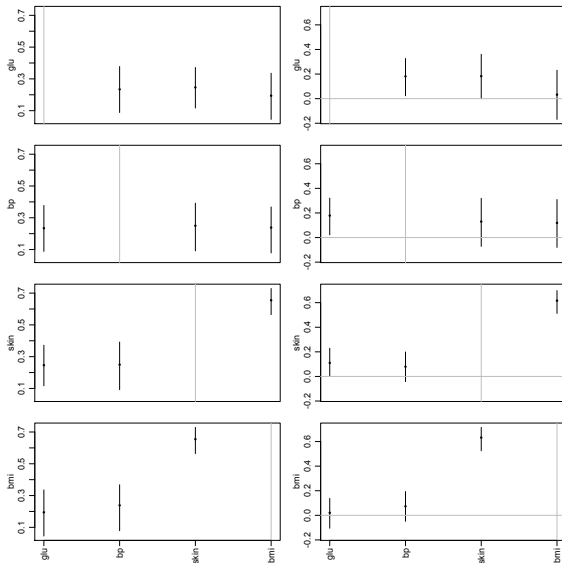


Figure 1: 95 percent posterior confidence intervals for correlations and regression coefficients.

Then do the prediction example from Hoff