

The Multivariate Normal Model

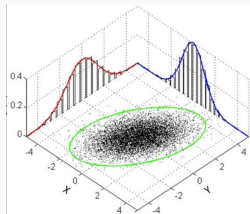
June 4, 2021

Table of contents

1. Overview
2. Conjugate inference
3. Semi-conjugate inference
4. Application: Reading Comprehension
5. Missing data and imputation

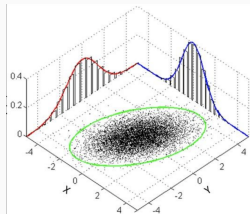
Overview

Some motivations for the normal



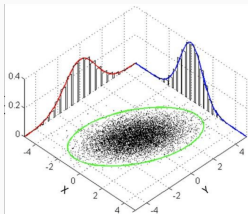
- *Maximum entropy* among all distributions with a given mean μ and variance Σ .

Some motivations for the normal



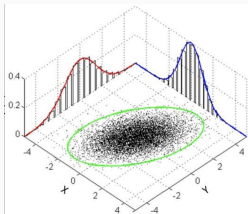
- *Maximum entropy* among all distributions with a given mean μ and variance Σ .
- Characterized by independence of sample mean and sample variance. (Bayesian take: ask if your beliefs about the sample mean are independent from those about the sample variance.)

Some motivations for the normal



- *Maximum entropy* among all distributions with a given mean μ and variance Σ .
- Characterized by independence of sample mean and sample variance. (Bayesian take: ask if your beliefs about the sample mean are independent from those about the sample variance.)
- Sample averages are generally approximately normally distributed due to the Central Limit Theorem.

Some motivations for the normal



- *Maximum entropy* among all distributions with a given mean μ and variance Σ .
- Characterized by independence of sample mean and sample variance. (Bayesian take: ask if your beliefs about the sample mean are independent from those about the sample variance.)
- Sample averages are generally approximately normally distributed due to the Central Limit Theorem.
- Sufficient statistics are sample mean and variance; so will consistently estimate population mean and variance even for non-normal distributions.

Why Bayesian normal?

- Prior information often exists and can be taken into account.
 - Population-level info: see the typing example, PIMA Indians example
 - Nature (e.g. support) of data: see the reading comprehension example

Why Bayesian normal?

- Prior information often exists and can be taken into account.
 - Population-level info: see the typing example, PIMA Indians example
 - Nature (e.g. support) of data: see the reading comprehension example
- ML estimates of covariance matrices have large variance.
 - Problem can be especially bad in certain contexts (e.g., small data, high-dimensions, missing data)
 - Spherical prior provides regularization
 - Posterior asymptotically concentrates around maximum likelihood (ML) solution

Why Bayesian normal?

- Prior information often exists and can be taken into account.
 - Population-level info: see the typing example, PIMA Indians example
 - Nature (e.g. support) of data: see the reading comprehension example
- ML estimates of covariance matrices have large variance.
 - Problem can be especially bad in certain contexts (e.g., small data, high-dimensions, missing data)
 - Spherical prior provides regularization
 - Posterior asymptotically concentrates around maximum likelihood (ML) solution
- Inference is no harder than for frequentist models
 - Easy, cheap updates (a conjugate prior exists)
 - Supports online learning
 - Fits nicely in more complex models
 - Nice hyperparameter interpretation

Conjugate inference

A conjugate prior

TODO: Fill in

Application: Modeling typing dynamics

See powerpoint slides.

Semi-conjugate inference

Semi-conjugate Bayesian MVN

Consider the following model with a normal sampling distribution and *semi-conjugate* prior

$$\boldsymbol{\mu} \sim \mathcal{N}_d(\boldsymbol{m}_0, \boldsymbol{V}_0)$$

$$\Sigma \sim \mathcal{W}^{-1}(\nu_0, \Psi_0)$$

$$\boldsymbol{x}_i \mid \boldsymbol{\mu}, \Sigma \stackrel{\text{iid}}{\sim} \mathcal{N}_d(\boldsymbol{\mu}, \Sigma), \quad i = 1, \dots, N$$

We define $\boldsymbol{x} := (\boldsymbol{x}_1, \dots, \boldsymbol{x}_N)$, where each $\boldsymbol{x}_i \in \mathbb{R}^d$.

Semi-conjugate Bayesian MVN

Consider the following model with a normal sampling distribution and *semi-conjugate* prior

$$\boldsymbol{\mu} \sim \mathcal{N}_d(\mathbf{m}_0, \mathbf{V}_0)$$

$$\Sigma \sim \mathcal{W}^{-1}(\nu_0, \Psi_0)$$

$$\mathbf{x}_i \mid \boldsymbol{\mu}, \Sigma \stackrel{\text{iid}}{\sim} \mathcal{N}_d(\boldsymbol{\mu}, \Sigma), \quad i = 1, \dots, N$$

We define $\mathbf{x} := (\mathbf{x}_1, \dots, \mathbf{x}_N)$, where each $\mathbf{x}_i \in \mathbb{R}^d$.

Fully conjugate vs semi-conjugate MVNs

This model is different than the model with the fully conjugate (Normal-Inverse-Wishart) prior on the pair $(\boldsymbol{\mu}, \Sigma)$. The conditionally conjugate prior lacks closed-form posterior updating, but is also more expressive. It is also easier to extend upwards.

Semi-conjugate models generally

Conjugate models

Conjugacy can be defined as follows (gelman2013bayesian). If \mathcal{F} is a class of sampling distributions and \mathcal{P} is a class of prior distributions for θ , then the class \mathcal{P} is *conjugate* for \mathcal{F} if

$$p(\theta \mid y) \in \mathcal{P} \text{ for all } p(\cdot \mid \theta) \in \mathcal{F} \text{ and } p(\cdot) \in \mathcal{P}$$

Semi-conjugate models

Conditional conjugacy (sometimes called semi-conjugacy) can be defined similarly (gelman2013bayesian). If \mathcal{F} is a class of sampling distributions and \mathcal{P} is a class of prior distributions for $\theta \mid \phi$, then the class \mathcal{P} is *conditionally conjugate* for \mathcal{F} if

$$p(\theta \mid \phi, y) \in \mathcal{P} \text{ for all } p(\cdot \mid \theta, \phi) \in \mathcal{F} \text{ and } p(\cdot \mid \phi) \in \mathcal{P}$$

Semi-conjugate models generally

Conjugate models

Conjugacy can be defined as follows (gelman2013bayesian). If \mathcal{F} is a class of sampling distributions and \mathcal{P} is a class of prior distributions for θ , then the class \mathcal{P} is *conjugate* for \mathcal{F} if

$$p(\theta \mid y) \in \mathcal{P} \text{ for all } p(\cdot \mid \theta) \in \mathcal{F} \text{ and } p(\cdot) \in \mathcal{P}$$

Semi-conjugate models

Conditional conjugacy (sometimes called semi-conjugacy) can be defined similarly (gelman2013bayesian). If \mathcal{F} is a class of sampling distributions and \mathcal{P} is a class of prior distributions for $\theta \mid \phi$, then the class \mathcal{P} is *conditionally conjugate* for \mathcal{F} if

$$p(\theta \mid \phi, y) \in \mathcal{P} \text{ for all } p(\cdot \mid \theta, \phi) \in \mathcal{F} \text{ and } p(\cdot \mid \phi) \in \mathcal{P}$$

In other words, a family of prior distributions for a parameter is called conditionally conjugate if the conditional posterior distribution (often called the *complete conditional*), given the data and all other parameters in the model, is also in that class.

Why are conditionally conjugate models of interest? The posterior distributions for conditionally conjugate models are easily approximated with Gibbs sampling or Mean Field Variational Inference – the former samples from the complete conditional, whereas the latter takes variational expectations with respect to the natural parameter of the complete conditional.

Complete conditionals for the Bayesian MVN

We sample from the posterior by iteratively sampling from the *complete conditionals*:

$$\boldsymbol{\mu} \mid \Sigma, \mathbf{x} \sim \mathcal{N}_d(\mathbf{m}, \mathbf{V})$$

where

$$\begin{aligned}\mathbf{m} &= \left(\mathbf{V}_0^{-1} + N\Sigma^{-1} \right)^{-1} \left(\mathbf{V}_0^{-1}\mathbf{m}_0 + N\Sigma^{-1}\bar{\mathbf{x}} \right) \\ \mathbf{V} &= \left(\mathbf{V}_0^{-1} + N\Sigma^{-1} \right)^{-1}\end{aligned}$$

and

$$\Sigma \mid \boldsymbol{\mu}, \mathbf{x} \sim \mathcal{W}^{-1}(\nu, \Psi)$$

where

$$\begin{aligned}\nu &= \nu_0 + N \\ \Psi &= \Psi_0 + \sum_{i=1}^N (\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^T\end{aligned}$$

Complete conditionals: Interpretation

These complete conditionals have nice interpretations:

- **Complete conditional for $(\mu \mid \Sigma, \mathbf{x})$:** On the precision scale, \mathbf{V} is the sum of the prior precision matrix \mathbf{V}_0^{-1} and N copies of the precision for each observation, Σ^{-1} . Similarly, \mathbf{m} is the precision-weighted convex combination of \mathbf{m}_0 , the prior mean, and the empirical average, $\bar{\mathbf{x}}$.
- **Complete conditional for $(\Sigma \mid \mu, \mathbf{x})$:** The covariance was estimated from ν observations with a sum of pairwise deviation products Ψ .

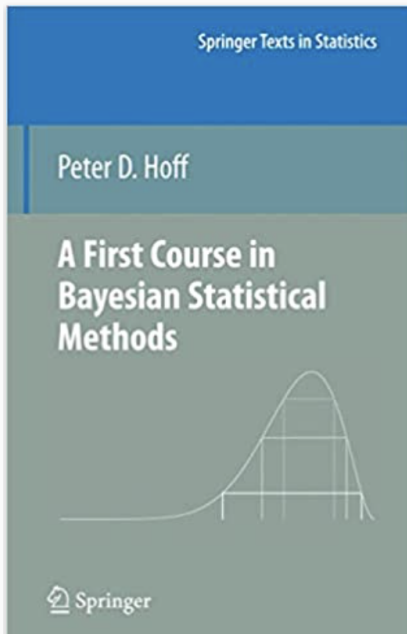
Complete conditionals: Proof

See exponential family notes.

Application: Reading Comprehension

See ipython notebook.

Missing data and imputation



Pima Dataset

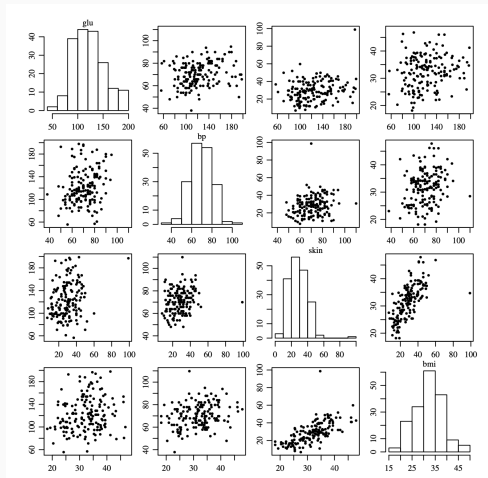


Figure 1: Univariate histograms and bivariate scatterplots for four variables taken from a dataset involving health-related measurements on 200 women of Pima Indian heritage living near Phoenix, Arizona. The four variables are `glu` (blood plasma glucose concentration), `bp` (diastolic blood pressure), `skin` (skin fold thickness), and `bmi` (body mass index).

Pima Dataset

	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
7	NA	58	31	34.3
8	193	50	16	25.9
9	142	80	15	NA
10	128	78	NA	43.3

Figure 2: Entries for the first ten subjects in the dataset. The NA's stand for "not available."

Description of problem

How to do parameter estimation in the presence of missing data?

We cannot do parameter estimation, because we cannot compute the likelihood

$$\prod_{i=1}^n p(\mathbf{y}_i \mid \boldsymbol{\theta}).$$

Two common approaches taken by software packages:

1. Throw away all subjects with missing data

Description of problem


How to do parameter estimation in the presence of missing data?

We cannot do parameter estimation, because we cannot compute the likelihood

$$\prod_{i=1}^n p(\mathbf{y}_i \mid \boldsymbol{\theta}).$$

Two common approaches taken by software packages:

1. Throw away all subjects with missing data

 Discards a potentially large amount of useful information.


Description of problem

How to do parameter estimation in the presence of missing data?

We cannot do parameter estimation, because we cannot compute the likelihood

$$\prod_{i=1}^n p(\mathbf{y}_i \mid \boldsymbol{\theta}).$$

Two common approaches taken by software packages:

1. Throw away all subjects with missing data
 Discards a potentially large amount of useful information.
2. Impute the population mean or some other fixed value.

Description of problem

How to do parameter estimation in the presence of missing data?

We cannot do parameter estimation, because we cannot compute the likelihood

$$\prod_{i=1}^n p(\mathbf{y}_i \mid \boldsymbol{\theta}).$$

Two common approaches taken by software packages:

1. Throw away all subjects with missing data

X Discards a potentially large amount of useful information.

2. Impute the population mean or some other fixed value.

X Assumes certainty about these values, when in fact we have not observed them.

Missing at random (MAR)

Let $\mathbf{O}_i = (O_1, \dots, O_p)^T$ be a binary vector such that

- $O_{ij} = 1 \implies Y_{ij}$ is observed
- $O_{ij} = 0 \implies Y_{ij}$ is missing

Definition

We say the missing data are *missing at random* if \mathbf{O}_i and \mathbf{Y}_i are conditionally independent given the model parameters θ and the distribution of \mathbf{O}_i does not depend on θ .

Missing at random (MAR)

Let $\mathbf{O}_i = (O_1, \dots, O_p)^T$ be a binary vector such that

- $O_{ij} = 1 \implies Y_{ij}$ is observed
- $O_{ij} = 0 \implies Y_{ij}$ is missing

Definition

We say the missing data are *missing at random* if \mathbf{O}_i and \mathbf{Y}_i are conditionally independent given the model parameters θ and the distribution of \mathbf{O}_i does not depend on θ .

Remark. This is one of the three types of missingness. In gist:

- Missing completely at random (MCAR) - missingness is independent of all data
- Missing at random (MAR) - missingness is independent of observed data
- Missing not at random (MNAR) - missingness depends on missing values (and perhaps observed data)

The likelihood in the presence of MAR data

When the data is missing at random, the sampling probability (density) for the data from observational unit i is given by

$$\begin{aligned} p(\mathbf{o}_i, \{y_{ij} : o_{ij} = 1\} \mid \boldsymbol{\theta}) &\stackrel{(1)}{=} p(\mathbf{o}_i) p(\{y_{ij} : o_{ij} = 1\} \mid \boldsymbol{\theta}) \\ &= p(\mathbf{o}_i) \int p(y_{i1}, \dots, y_{ip} \mid \boldsymbol{\theta}) \prod_{y_{ij}: o_{ij}=0} dy_{ij} \end{aligned}$$

where in (1) we applied the definition of MAR.

The likelihood in the presence of MAR data

When the data is missing at random, the sampling probability (density) for the data from observational unit i is given by

$$\begin{aligned} p(\mathbf{o}_i, \{y_{ij} : o_{ij} = 1\} \mid \boldsymbol{\theta}) &\stackrel{(1)}{=} p(\mathbf{o}_i) p(\{y_{ij} : o_{ij} = 1\} \mid \boldsymbol{\theta}) \\ &= p(\mathbf{o}_i) \int p(y_{i1}, \dots, y_{ip} \mid \boldsymbol{\theta}) \prod_{y_{ij}: o_{ij}=0} dy_{ij} \end{aligned}$$

where in (1) we applied the definition of MAR.

✓ So in the presence of MAR data, the correct thing to do is *integrate* over the missing data to obtain the marginal probability (density) of the observed data.

Utilization in multivariate normal models

In the case of multivariate normal models (so $\theta = (\mu, \Sigma)$), the integration is easy: Multivariate normals have normal marginals.

Example

Suppose $\mathbf{y}_i = (y_{i1}, \text{NA}, y_{i3}, \text{NA})^T$, so $\mathbf{o}_i = (1, 0, 1, 0)^T$.

Then

$$\begin{aligned} p(\mathbf{o}_i, y_{i1}, y_{i3} \mid \mu, \Sigma) &= p(\mathbf{o}_i) p(y_{i1}, y_{i3} \mid \mu, \Sigma) \\ &= p(\mathbf{o}_i) \int p(\mathbf{y}_i \mid \mu, \Sigma) dy_2 dy_4 \end{aligned}$$

The marginal density $p(y_{i1}, y_{i3} \mid \theta)$ is simply a bivariate normal density with mean $(\mu_1, \mu_3)^T$ and covariance matrix made up of $(\sigma_1^2, \sigma_{13}, \sigma_3^2)$.

Gibbs sampling with missing data

Complete data

If \mathbf{Y} is the $n \times p$ matrix in which $o_{i,j} = 1$ if $Y_{i,j}$ is observed and $o_{i,j} = 0$ if $Y_{i,j}$ is missing, then \mathbf{Y} has two parts

- $\mathbf{Y}_{\text{obs}} := \{y_{i,j} : o_{i,j} = 1\}$, the data that we observe, and
- $\mathbf{Y}_{\text{miss}} := \{y_{i,j} : o_{i,j} = 0\}$, the data that we do not observe.

Gibbs sampler

A Gibbs sampling scheme for approximating the posterior is given by:

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

Gibbs sampling with missing data

Complete data

If \mathbf{Y} is the $n \times p$ matrix in which $o_{i,j} = 1$ if $Y_{i,j}$ is observed and $o_{i,j} = 0$ if $Y_{i,j}$ is missing, then \mathbf{Y} has two parts

- $\mathbf{Y}_{\text{obs}} := \{y_{i,j} : o_{i,j} = 1\}$, the data that we observe, and
- $\mathbf{Y}_{\text{miss}} := \{y_{i,j} : o_{i,j} = 0\}$, the data that we do not observe.

Gibbs sampler

A Gibbs sampling scheme for approximating the posterior is given by:

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

The first two steps are the same as before! The third step is covered in the next slide. Any guesses?

Sampling the missing data

$$\begin{aligned} p(\mathbf{Y}_{\text{miss}} \mid \mathbf{Y}_{\text{obs}}, \boldsymbol{\mu}, \Sigma) &\propto p(\mathbf{Y}_{\text{miss}}, \mathbf{Y}_{\text{obs}} \mid \boldsymbol{\mu}, \Sigma) \\ &= \prod_{i=1}^n p(\mathbf{y}_{i,\text{miss}}, \mathbf{y}_{i,\text{obs}} \mid \boldsymbol{\mu}, \Sigma) \\ &\propto \prod_{i=1}^n p(\mathbf{y}_{i,\text{miss}} \mid \mathbf{y}_{i,\text{obs}}, \boldsymbol{\mu}, \Sigma) \end{aligned}$$

How to proceed?

Sampling the missing data

$$\begin{aligned} p(\mathbf{Y}_{\text{miss}} \mid \mathbf{Y}_{\text{obs}}, \boldsymbol{\mu}, \Sigma) &\propto p(\mathbf{Y}_{\text{miss}}, \mathbf{Y}_{\text{obs}} \mid \boldsymbol{\mu}, \Sigma) \\ &= \prod_{i=1}^n p(\mathbf{y}_{i,\text{miss}}, \mathbf{y}_{i,\text{obs}} \mid \boldsymbol{\mu}, \Sigma) \\ &\propto \prod_{i=1}^n p(\mathbf{y}_{i,\text{miss}} \mid \mathbf{y}_{i,\text{obs}}, \boldsymbol{\mu}, \Sigma) \end{aligned}$$

How to proceed? We apply standard results about conditional distributions formed from partitions of multivariate normals:

$$\begin{aligned} \mathbf{y}_{[b]} \mid \mathbf{y}_{[a]}, \boldsymbol{\mu}, \Sigma &\sim \mathcal{MVN}\left(\boldsymbol{\mu}_{b|a}, \Sigma_{b|a}\right), \quad \text{where} \\ \boldsymbol{\mu}_{b|a} &= \boldsymbol{\mu}_{[b]} + \Sigma_{[b,a]}(\Sigma_{[a,a]})^{-1}(\mathbf{y}_{[a]} - \boldsymbol{\mu}_{[a]}) \\ \Sigma_{b|a} &= \Sigma_{[b,b]} - \Sigma_{[b,a]}(\Sigma_{[a,a]})^{-1}\Sigma_{[a,b]} \end{aligned}$$

Some macroscopic properties:

Sampling the missing data

$$\begin{aligned} p(\mathbf{Y}_{\text{miss}} \mid \mathbf{Y}_{\text{obs}}, \boldsymbol{\mu}, \Sigma) &\propto p(\mathbf{Y}_{\text{miss}}, \mathbf{Y}_{\text{obs}} \mid \boldsymbol{\mu}, \Sigma) \\ &= \prod_{i=1}^n p(\mathbf{y}_{i,\text{miss}}, \mathbf{y}_{i,\text{obs}} \mid \boldsymbol{\mu}, \Sigma) \\ &\propto \prod_{i=1}^n p(\mathbf{y}_{i,\text{miss}} \mid \mathbf{y}_{i,\text{obs}}, \boldsymbol{\mu}, \Sigma) \end{aligned}$$

How to proceed? We apply standard results about conditional distributions formed from partitions of multivariate normals:

$$\begin{aligned} \mathbf{y}_{[b]} \mid \mathbf{y}_{[a]}, \boldsymbol{\mu}, \Sigma &\sim \mathcal{MVN}\left(\boldsymbol{\mu}_{b|a}, \Sigma_{b|a}\right), \quad \text{where} \\ \boldsymbol{\mu}_{b|a} &= \boldsymbol{\mu}_{[b]} + \Sigma_{[b,a]}(\Sigma_{[a,a]})^{-1}(\mathbf{y}_{[a]} - \boldsymbol{\mu}_{[a]}) \\ \Sigma_{b|a} &= \Sigma_{[b,b]} - \Sigma_{[b,a]}(\Sigma_{[a,a]})^{-1}\Sigma_{[a,b]} \end{aligned}$$

Some macroscopic properties:

- The conditional mean, $\boldsymbol{\mu}_{b|a}$, starts off at the unconditional mean, $\boldsymbol{\mu}_{[b]}$, but then is modified by $(\mathbf{y}_{[a]} - \boldsymbol{\mu}_{[a]})$ in a way that depends on the covariance $\Sigma_{[b,a]}$.
- The conditional variance $\Sigma_{b|a}$ is less than the unconditional variance $\Sigma_{[b,b]}$

Posterior Correlations

To each covariance matrix there corresponds a correlation matrix \mathbf{C} given by

$$\mathbf{C} := \left\{ c_{jk} : c_{jk} = \Sigma_{[j,k]} / \sqrt{\Sigma_{[j,j]} \Sigma_{[k,k]}} \right\}$$

Simply taking the mean across samples, we obtain the approximation

$$\mathbb{E}[\mathbf{C} \mid \mathbf{y}_1, \dots, \mathbf{y}_n] = \begin{bmatrix} 1.00 & 0.23 & 0.25 & 0.19 \\ 0.23 & 1.00 & 0.25 & 0.24 \\ 0.25 & 0.25 & 1.00 & 0.65 \\ 0.19 & 0.24 & 0.65 & 1.00 \end{bmatrix}$$

Notes:

- Bayesian paradigm again yielding unlimited access to posterior functionals of interest, without doing any extra inferential work!
- Correlations are generally of interest for multivariate normal models, but they are *especially* relevant to imputation.

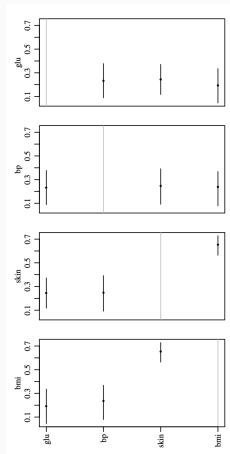


Figure 3: 95% posterior confidence intervals for correlations

Intelligent imputations

The posterior expectation gives a much better imputation than some flat fixed value.

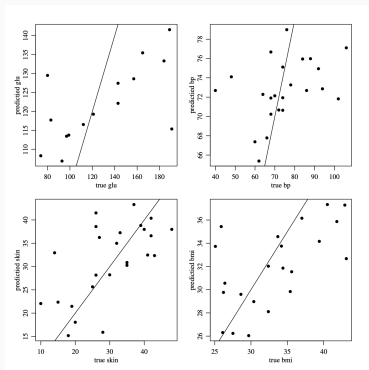


Figure 4: True values of the missing data vs. posterior expectations

Intelligent imputations

The posterior expectation gives a much better imputation than some flat fixed value.

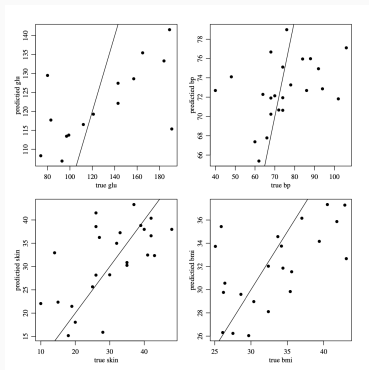


Figure 4: True values of the missing data vs. posterior expectations

Imputations are especially good for skin and bmi, due to their higher correlations.