

# The Multivariate Normal Model

---

June 9, 2021

# Table of contents

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

# Overview

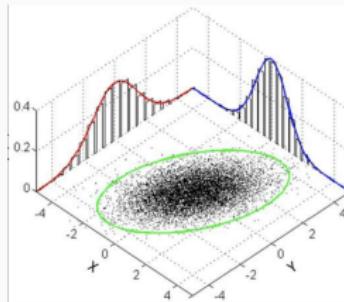
---

# Overview

---

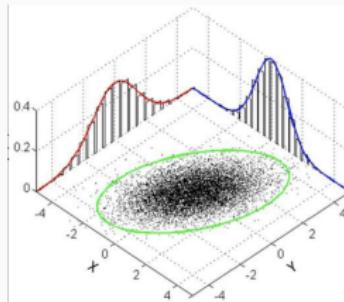
## Motivations

## Some motivations for the normal



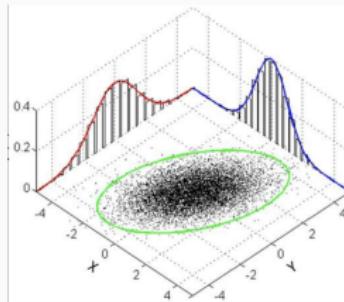
- *Maximum entropy* among all distributions with a given mean  $\mu$  and variance  $\Sigma$ .

## Some motivations for the normal



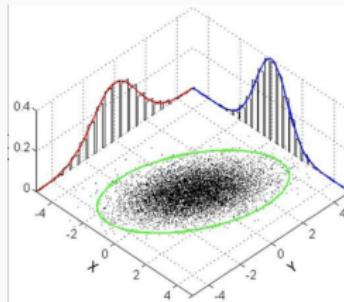
- *Maximum entropy* among all distributions with a given mean  $\mu$  and variance  $\Sigma$ .
- Characterized by independence of sample mean and sample variance.  
(Bayesian take: ask if your beliefs about the sample mean are independent from those about the sample variance.)

## Some motivations for the normal



- *Maximum entropy* among all distributions with a given mean  $\mu$  and variance  $\Sigma$ .
- 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  $\mu$  and variance  $\Sigma$ .
- 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 (Previous example: disease prevalence. Forthcoming example: biometrics and PIMA Indians)
  - Nature (e.g. support) of data (Forthcoming example: reading comprehension)

# Why Bayesian normal?

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

# Why Bayesian normal?

- Prior information often exists and can be taken into account.
  - Population-level info (Previous example: disease prevalence. Forthcoming example: biometrics and PIMA Indians)
  - Nature (e.g. support) of data (Forthcoming example: reading comprehension)
- ML estimates of covariance matrices can have large variance.
  - Problem can be especially bad in certain contexts ( e.g., small data, high-dimensions, missing data)
  - Spherical prior provides regularization
  - Posterior still asymptotically concentrates around maximum likelihood (ML) solution
- Inference can be as easy or easier than in frequentist models
  - Easy, cheap updates (esp. when using a conjugate prior)
  - Supports online learning
  - Fits nicely in more complex models

# **Overview**

---

## **Building blocks**

## Review: Exponential families

We have observed that the Bernoulli and the Poisson distributions are both **exponential families**.

### Exponential family

An *exponential family* is a set of probability distributions whose probability density functions have the following form

$$p(x \mid \theta) = h(x) \exp\{\eta(\theta)^T t(x) - a(\theta)\} \quad (1)$$

where we refer to  $h$  as the base measure,  $\eta$  as the natural parameter,  $t$  as the sufficient statistics, and  $a$  as the log normalizer.

## Multivariate normal is an exponential family

We can write the density of a multivariate normal  $\mathcal{N}(\mu, \Sigma)$  distribution in exponential form. (*Try it!*)

# Multivariate normal is an exponential family

We can write the density of a multivariate normal  $\mathcal{N}(\mu, \Sigma)$  distribution in exponential form. (*Try it!*)

$$\begin{aligned} p(\mathbf{x} \mid \boldsymbol{\mu}, \boldsymbol{\Sigma}) &= (2\pi)^{-d/2} |\boldsymbol{\Sigma}|^{-1/2} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\} \\ &= (2\pi)^{-d/2} \exp \left\{ -\frac{1}{2} \underbrace{\mathbf{x}^T \boldsymbol{\Sigma}^{-1} \mathbf{x}}_{-\frac{1}{2} \text{vec}(\boldsymbol{\Sigma}^{-1})^T \text{vec}(\mathbf{x}\mathbf{x}^T)} + \underbrace{(\boldsymbol{\Sigma}^{-1} \boldsymbol{\mu})^T}_{\text{vec}(\boldsymbol{\Sigma}^{-1} \boldsymbol{\mu})} \mathbf{x} - \frac{1}{2} \boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} + \frac{1}{2} \log |\boldsymbol{\Sigma}^{-1}| \right\} \end{aligned} \tag{2}$$

Note: the underbrace representation is given by  $\mathbf{x}^T \boldsymbol{\Sigma}^{-1} \mathbf{x} = \text{tr}(\mathbf{x}^T \boldsymbol{\Sigma}^{-1} \mathbf{x}) = \text{tr}(\boldsymbol{\Sigma}^{-1} \mathbf{x} \mathbf{x}^T) = \text{vec}(\boldsymbol{\Sigma}^{-1})^T \text{vec}(\mathbf{x} \mathbf{x}^T)$ .

# Multivariate normal is an exponential family

We can write the density of a multivariate normal  $\mathcal{N}(\mu, \Sigma)$  distribution in exponential form. (*Try it!*)

$$\begin{aligned} p(x | \mu, \Sigma) &= (2\pi)^{-d/2} |\Sigma|^{-1/2} \exp \left\{ -\frac{1}{2}(x - \mu)^T \Sigma^{-1} (x - \mu) \right\} \\ &= (2\pi)^{-d/2} \exp \left\{ -\frac{1}{2} \underbrace{x^T \Sigma^{-1} x}_{-\frac{1}{2} \text{vec}(\Sigma^{-1})^T \text{vec}(xx^T)} + (\Sigma^{-1} \mu)^T x - \frac{1}{2} \mu^T \Sigma^{-1} \mu + \frac{1}{2} \log |\Sigma^{-1}| \right\} \end{aligned} \quad (2)$$

Note: the underbrace representation is given by  $x^T \Sigma^{-1} x = \text{tr}(x^T \Sigma^{-1} x) = \text{tr}(\Sigma^{-1} x x^T) = \text{vec}(\Sigma^{-1})^T \text{vec}(xx^T)$ .

- **Natural parameter:**  $\eta(\mu, \Sigma) = (-\frac{1}{2} \text{vec}(\Sigma^{-1}), \Sigma^{-1} \mu)$ ,
- **Sufficient statistics:**  $t(x) = (\text{vec}(xx^T), x)$ ,
- **Log normalizer:**  $a(\mu, \Sigma) = -\frac{1}{2} \mu^T \Sigma^{-1} \mu + \frac{1}{2} \log |\Sigma^{-1}|$ .
- **Base measure:**  $h(x) = (2\pi)^{-d/2}$

# Qualitative points

## The natural parametrization of the MVN

From Equation (2), we see that the natural parameters of the MVN are

- the *precision*  $\Sigma^{-1}$ , and
- the *precision-weighted mean*  $\Sigma^{-1}\mu$ .

## Are you exponentiating a quadratic?

We also see that if a random vector  $\mathbf{x}$  has a density on  $\mathbb{R}^d$  that satisfies

$$p(\mathbf{x}) \propto \exp\left\{-\frac{1}{2}\mathbf{x}^T \mathbf{A}\mathbf{x} + \mathbf{x}^T \mathbf{b}\right\}$$

for some matrix  $\mathbf{A}$  and vector  $\mathbf{b}$ , then  $\mathbf{x}$  must be multivariate normal.

Moreover its natural parameters are precision  $\mathbf{A}$  and precision-weighted mean  $\mathbf{b}$ . In other words, the covariance is  $\mathbf{A}^{-1}$  and the mean is  $\mathbf{A}^{-1}\mathbf{b}$ .

# Self-help guide for deriving Gibbs samplers



orientating a quadratic?

## Inverse Wishart Distribution

The Inverse Wishart is a distribution on symmetric, positive definite matrices. The Inverse Wishart distribution, denoted  $\mathcal{W}^{-1}(\nu, \Psi)$ , has density

$$p(\Sigma) \propto |\Sigma|^{-(\nu+d+1)/2} \exp \left[ -\frac{1}{2} \text{tr}(\Sigma^{-1} \Psi) \right] \quad (3)$$

where  $\Sigma \succ 0$  and  $\nu > d - 1$  to have a proper prior. The expected value of an Inverse Wishart random variable parametrized as in (3) is given by  $\mathbb{E}[\Sigma] = \frac{\Psi}{\nu-d-1}$ .

### Interpreting the parameters of the Inverse Wishart

Note that the parameters of the Inverse Wishart can be interpreted (as per conjugacy; we will see this below) in the following way: the covariance was estimated from  $\nu$  observations with a residual sum of squares (a.k.a. sum of pairwise deviation products)  $\Psi$ .

# Inverse Wishart Distribution

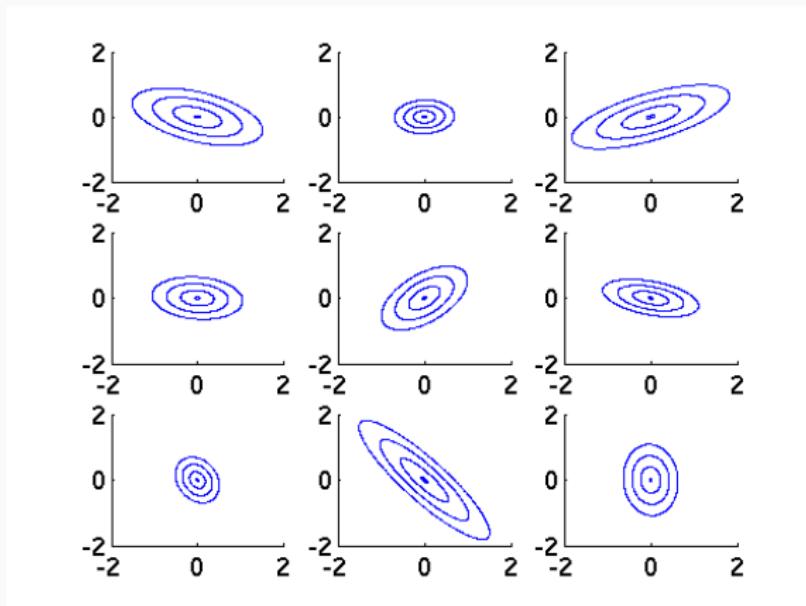


Image Credit: Michael Hughes, Tufts University

## Sampling from the Inverse Wishart

A sample  $\Sigma$  from the  $\mathcal{W}^{-1}(\nu, \Psi)$  distribution can be obtained by the following scheme:

1. Sample  $\mathbf{z}_1, \dots, \mathbf{z}_\nu \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \Psi^{-1})$
2. Calculate  $\mathbf{Z}^T \mathbf{Z} = \sum_{i=1}^\nu \mathbf{z}_i \mathbf{z}_i^T$ .
3. Set  $\Sigma = (\mathbf{Z}^T \mathbf{Z})^{-1}$ .

The intuition is that the Inverse Wishart models covariance matrices as an inverse sum of squares.

## **Conjugate inference**

---

## A fully conjugate formulation

Since the MVN is an EF, we can form a fully conjugate model.

**Fully conjugate Bayesian MVN** (has **closed-form** posteriors and predictive posteriors!)

Given observations  $\mathbf{x} := (\mathbf{x}_1, \dots, \mathbf{x}_N)$ , where each  $\mathbf{x}_i \in \mathbb{R}^d$ , we take

$$\boldsymbol{\mu}, \boldsymbol{\Sigma} \sim \text{NIW}_d(\boldsymbol{\alpha}_0, \boldsymbol{\mu}_0, \nu_0, \boldsymbol{\Psi}_0)$$

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

# A fully conjugate formulation

Since the MVN is an EF, we can form a fully conjugate model.

## Fully conjugate Bayesian MVN (has closed-form posteriors and predictive posteriors!)

Given observations  $\mathbf{x} := (\mathbf{x}_1, \dots, \mathbf{x}_N)$ , where each  $\mathbf{x}_i \in \mathbb{R}^d$ , we take

$$\boldsymbol{\mu}, \Sigma \sim \text{NIW}_d(\alpha_0, \boldsymbol{\mu}_0, \nu_0, \boldsymbol{\Psi}_0)$$

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

## The Normal-Inverse-Wishart prior

$$\boldsymbol{\mu}, \Sigma \sim \text{NIW}_d(\alpha_0, \boldsymbol{\mu}_0, \nu_0, \boldsymbol{\Psi}_0)$$

which means

$$\Sigma \sim \mathcal{W}_d^{-1}(\nu_0, \boldsymbol{\Psi}_0)$$

$$\boldsymbol{\mu} \mid \Sigma \sim \mathcal{N}_d(\boldsymbol{\mu}_0, \frac{1}{\alpha_0} \Sigma)$$

## Application: Modeling typing dynamics

See powerpoint slides.

## **Semi-conjugate inference**

---

# Semi-conjugate Bayesian normal

## Semi-conjugate Bayesian MVN

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

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

$$\boldsymbol{\Sigma} \sim \mathcal{W}^{-1}(\nu_0, \boldsymbol{\Psi}_0)$$

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

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

# Semi-conjugate models: Gist

**Main idea:** A family of prior distributions for a parameter is called *semi-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 family.

## Comparison

Fully conjugate model

- $(\mu, \Sigma) | \text{rest is}$   
Normal-Inverse-Wishart

Semi-conjugate model

- $\mu | \text{rest is Normal}$
- $\Sigma | \text{rest is Inverse Wishart}$

## Semi-conjugate models: Gist

**Main idea:** A family of prior distributions for a parameter is called *semi-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 family.

### Comparison

Fully conjugate model

- $(\mu, \Sigma) | \text{rest is}$   
Normal-Inverse-Wishart

Semi-conjugate model

- $\mu | \text{rest is Normal}$
- $\Sigma | \text{rest is Inverse Wishart}$

### Relative evaluation of semi-conjugate model

# Semi-conjugate models: Gist

**Main idea:** A family of prior distributions for a parameter is called *semi-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 family.

## Comparison

Fully conjugate model

- $(\mu, \Sigma) | \text{rest is}$   
Normal-Inverse-Wishart

Semi-conjugate model

- $\mu | \text{rest is Normal}$
- $\Sigma | \text{rest is Inverse Wishart}$

## Relative evaluation of semi-conjugate model

✗ lacks closed-form posterior updating

# Semi-conjugate models: Gist

**Main idea:** A family of prior distributions for a parameter is called *semi-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 family.

## Comparison

Fully conjugate model

- $(\mu, \Sigma) | \text{rest is}$   
Normal-Inverse-Wishart

Semi-conjugate model

- $\mu | \text{rest is Normal}$
- $\Sigma | \text{rest is Inverse Wishart}$

## Relative evaluation of semi-conjugate model

- ✗ lacks closed-form posterior updating
- ✓ more expressive

# Semi-conjugate models: Gist

**Main idea:** A family of prior distributions for a parameter is called *semi-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 family.

## Comparison

### Fully conjugate model

- $(\mu, \Sigma) \mid \text{rest is}$   
Normal-Inverse-Wishart

### Semi-conjugate model

- $\mu \mid \text{rest is Normal}$
- $\Sigma \mid \text{rest is Inverse Wishart}$

## Relative evaluation of semi-conjugate model

- ✗ lacks closed-form posterior updating
- ✓ more expressive
- ✓ easier to embed in more complex models

# Semi-conjugate models: Definition

## Review: Conjugate models

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

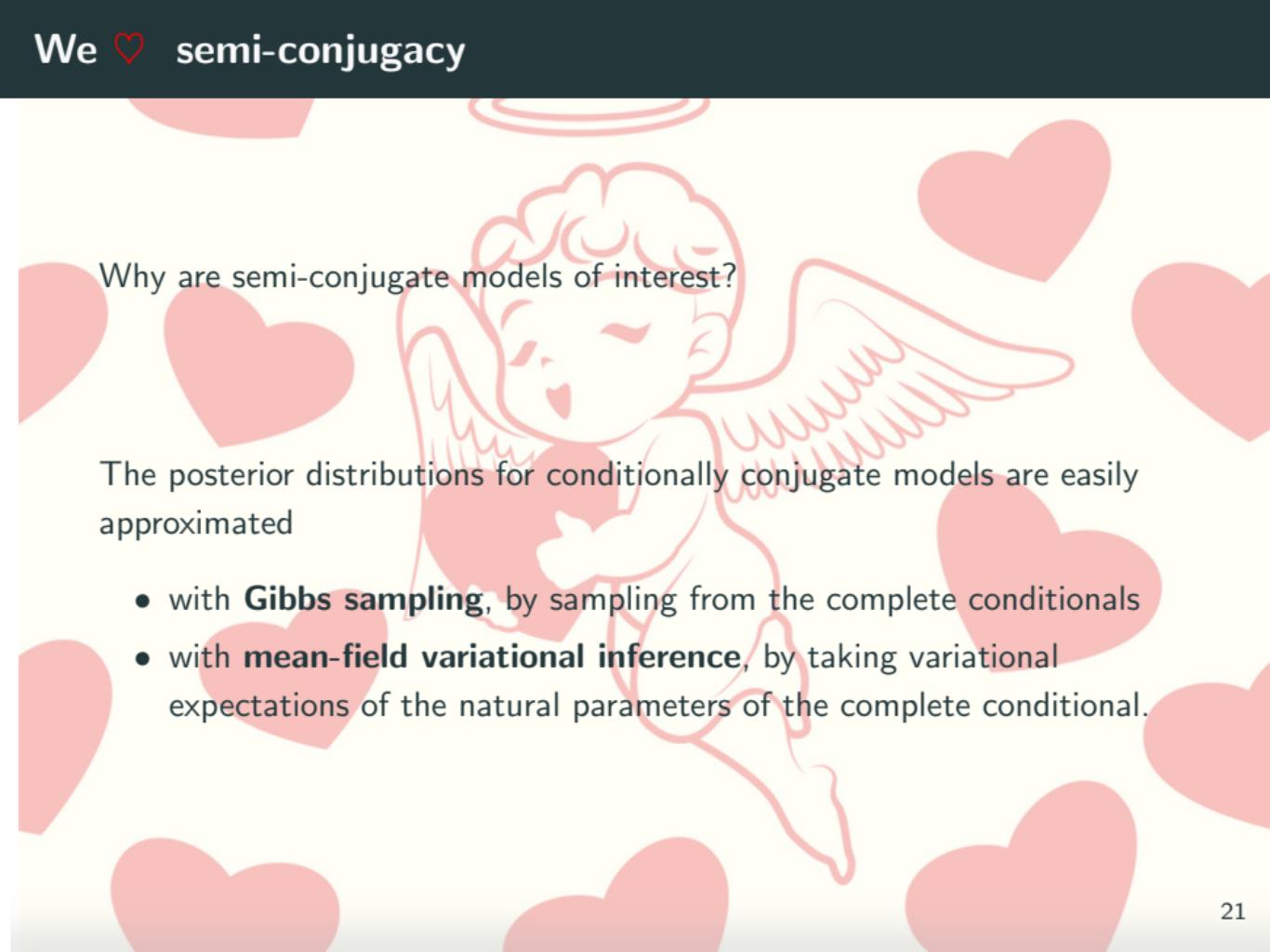
$$p(\theta | y) \in \mathcal{P} \text{ for all } p(\cdot | \theta) \in \mathcal{F} \text{ and } p(\cdot) \in \mathcal{P}$$

## Semi-conjugate models

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

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

# We ❤️ semi-conjugacy



Why are semi-conjugate models of interest?

The posterior distributions for conditionally conjugate models are easily approximated

- with **Gibbs sampling**, by sampling from the complete conditionals
- with **mean-field variational inference**, by taking variational expectations of the natural parameters of the complete conditional.

# The complete conditional for $\mu$

We use the exponential family representation of the MVN to represent the prior in terms of its natural parameters

$$p(\mu) \propto \exp \left\{ -\frac{1}{2} \mu^T \left( \underbrace{\boldsymbol{V}_0^{-1}}_{\text{prior precision}} \right) \mu + \mu^T \left( \underbrace{\boldsymbol{V}_0^{-1} \mathbf{m}_0}_{\text{prior precision-weighted mean}} \right) \right\} \quad (4)$$

And similarly, we write the likelihood  $L(\mu) = p(\mathbf{y} | \mu) = \prod_{i=1}^n p(\mathbf{y}_i | \mu)$  as

$$\begin{aligned} L(\mu) &\propto \exp \left\{ -\frac{1}{2} \sum_{i=1}^n (\mathbf{x}_i - \mu)^T \Sigma^{-1} (\mathbf{x}_i - \mu) \right\} \\ &= \exp \left\{ -\frac{1}{2} \mu^T \left( \underbrace{n \Sigma^{-1}}_{\text{data precision}} \right) \mu + \mu^T \left( \underbrace{\Sigma^{-1} n \bar{\mathbf{x}}}_{\text{data precision-weighted mean}} \right) \right\} \end{aligned} \quad (5)$$

Can you finish the derivation? (Hint: Remember Jimi!)

# The complete conditional for $\mu$

We use the exponential family representation of the MVN to represent the prior in terms of its natural parameters

$$p(\mu) \propto \exp \left\{ -\frac{1}{2} \mu^T \left( \underbrace{\mathbf{V}_0^{-1}}_{\text{prior precision}} \right) \mu + \mu^T \left( \underbrace{\mathbf{V}_0^{-1} \mathbf{m}_0}_{\text{prior precision-weighted mean}} \right) \right\} \quad (4)$$

And similarly, we write the likelihood  $L(\mu) = p(\mathbf{y} | \mu) = \prod_{i=1}^n p(\mathbf{y}_i | \mu)$  as

$$\begin{aligned} L(\mu) &\propto \exp \left\{ -\frac{1}{2} \sum_{i=1}^n (\mathbf{x}_i - \mu)^T \Sigma^{-1} (\mathbf{x}_i - \mu) \right\} \\ &= \exp \left\{ -\frac{1}{2} \mu^T \left( \underbrace{n \Sigma^{-1}}_{\text{data precision}} \right) \mu + \mu^T \left( \underbrace{\Sigma^{-1} n \bar{\mathbf{x}}}_{\text{data precision-weighted mean}} \right) \right\} \end{aligned} \quad (5)$$

Can you finish the derivation? (Hint: Remember Jimi!) By Bayes' law, and combining like terms,

$$\begin{aligned} p(\mu | \mathbf{x}, \Sigma) &\propto \underbrace{p(\mu)}_{\text{prior}} \underbrace{p(\mathbf{x} | \mu, \Sigma)}_{\text{likelihood}} \\ &= \exp \left\{ -\frac{1}{2} \mu^T \left( \underbrace{\mathbf{V}_0^{-1} + n \Sigma^{-1}}_{\text{posterior precision}} \right) \mu + \mu^T \left( \underbrace{\mathbf{V}_0 \mathbf{m}_0 + \Sigma^{-1} n \bar{\mathbf{x}}}_{\text{posterior precision-weighted mean}} \right) \right\} \end{aligned}$$

which reveals that the posterior is normal, along with the particular form of its parameter (covariance  $\mathbf{A}^{-1}$  and mean  $\mathbf{A}^{-1}\mathbf{b}$ ).

# Complete conditionals for the Bayesian MVN

The complete conditions under natural parametrization

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

where

$$\underbrace{\mathbf{V}^{-1}}_{\text{posterior precision}} = \underbrace{\mathbf{V}_0^{-1}}_{\text{prior precision}} + \underbrace{N\Sigma^{-1}}_{\text{data precision}}$$

$$\underbrace{\mathbf{V}^{-1}\mathbf{m}}_{\text{posterior precision-weighted mean}} = \underbrace{\mathbf{V}_0^{-1}\mathbf{m}_0}_{\text{prior precision-weighted mean}} + \underbrace{N\Sigma^{-1}\bar{\mathbf{x}}}_{\text{data precision-weighted mean}}$$

and

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

where

$$\underbrace{\nu}_{\text{posterior sample size}} = \underbrace{\nu_0}_{\text{prior sample size}} + \underbrace{N}_{\text{sample size}}$$

$$\underbrace{\Psi}_{\text{posterior RSS}} = \underbrace{\Psi_0}_{\text{prior RSS}} + \underbrace{\sum_{i=1}^N (\mathbf{x}_i - \mu)(\mathbf{x}_i - \mu)^T}_{\text{data RSS}}$$

Note: RSS = Residual Sum of Squares

## Complete conditionals: Interpretation

These complete conditionals have nice interpretations:

- **Complete conditional for  $(\mu | \Sigma, \mathbf{x})$ :** On the precision scale,  $\mathbf{V}$  is the sum of the prior precision matrix  $\mathbf{V}_0^{-1}$  and the data precision matrix (which is  $N$  copies of the precision for each observation,  $\Sigma^{-1}$ ). 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  $(\Sigma | \mu, \mathbf{x})$ :** The covariance was estimated from  $\nu$  observations with a sum of pairwise deviation products  $\Psi$ .

## Gibbs Sampler (in standard parametrization)

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

$$\mathbf{V} = \left( \underbrace{\mathbf{V}_0^{-1}}_{\text{prior precision}} + \underbrace{N\Sigma^{-1}}_{\text{data precision}} \right)^{-1}$$
$$\mathbf{m} = \mathbf{V} \left( \underbrace{\mathbf{V}_0^{-1} \mathbf{m}_0}_{\text{prior precision-weighted mean}} + \underbrace{N\Sigma^{-1} \bar{\mathbf{x}}}_{\text{data precision-weighted mean}} \right)$$

and

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

where

$$\nu = \nu_0 + N$$

$$\Psi = \Psi_0 + \sum_{i=1}^N (\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^T$$

## **Application: Reading Comprehension**

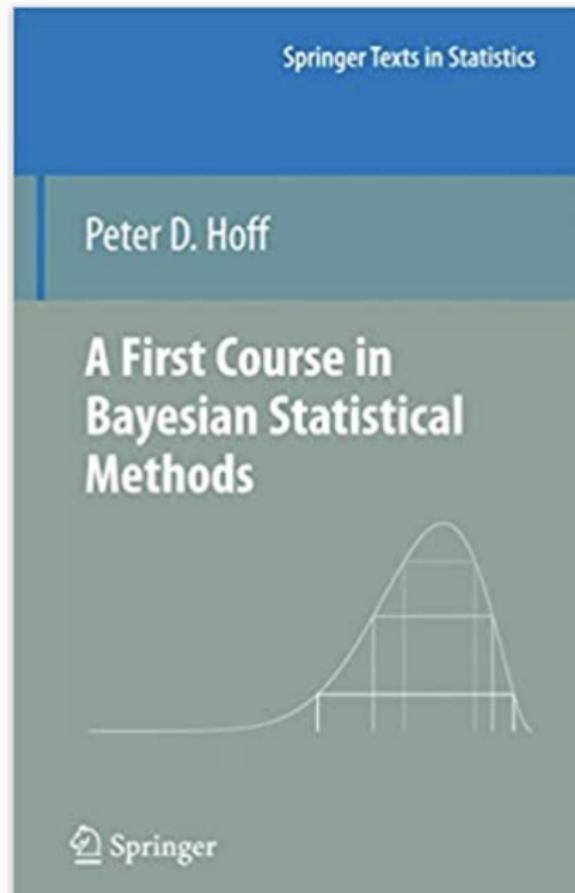
---

See ipython notebook.

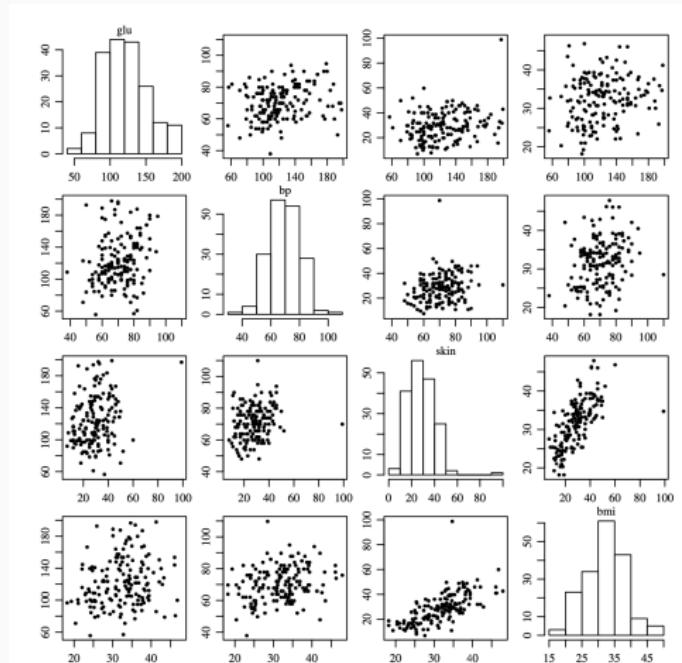
## **Missing data and imputation**

---

## References



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

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

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

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

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

## Missing at random (MAR)

Let  $\mathbf{O}_i = (O_{i1}, \dots, O_{ip})^T$  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  $\theta$  and the distribution of  $\mathbf{O}_i$  does not depend on  $\theta$ .

# Missing at random (MAR)

Let  $\mathbf{O}_i = (O_{i1}, \dots, O_{ip})^T$  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  $\theta$  and the distribution of  $\mathbf{O}_i$  does not depend on  $\theta$ .

**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

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

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

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

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

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

The marginal density  $p(y_{i1}, y_{i3} | \boldsymbol{\theta})$  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  $p(\boldsymbol{\mu} | \mathbf{Y}_{\text{obs}}, \mathbf{Y}_{\text{miss}}^{(s)}, \boldsymbol{\Sigma}^{(s)})$ ;
2. Sampling  $\boldsymbol{\Sigma}^{(s+1)}$  from  $p(\boldsymbol{\Sigma} | \mathbf{Y}_{\text{obs}}, \mathbf{Y}_{\text{miss}}^{(s)}, \boldsymbol{\mu}^{(s+1)})$ ;
3. Sampling  $\mathbf{Y}_{\text{miss}}^{(s+1)}$  from  $p(\mathbf{Y}_{\text{miss}} | \mathbf{Y}_{\text{obs}}, \boldsymbol{\mu}^{(s+1)}, \boldsymbol{\Sigma}^{(s+1)})$ ;

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

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

# Sampling the missing data

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

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

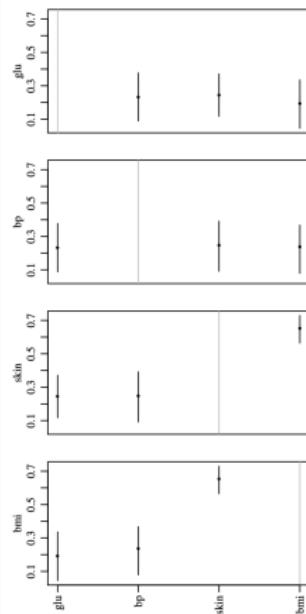
$$\mathbf{C} := \left\{ c_{jk} : c_{jk} = \Sigma_{[j,k]} / \sqrt{\Sigma_{[j,j]} \Sigma_{[k,k]}} \right\}$$

Simply taking the mean across samples, we obtain the approximation

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

Notes:

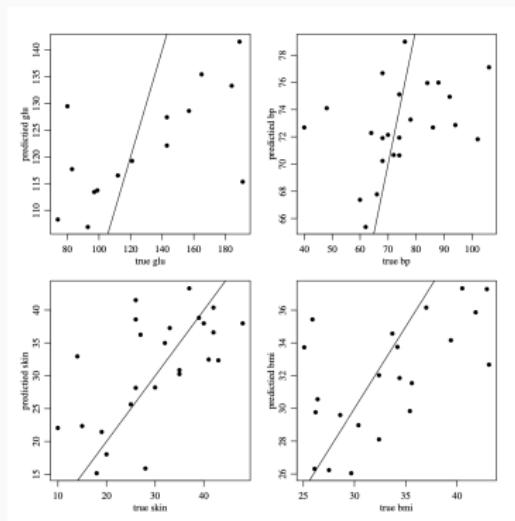
- Bayesian paradigm again yielding unlimited access to posterior functionals of interest, without doing any extra inferential work!
- Correlations are generally of interest for multivariate normal models, but they are *especially* relevant to imputation.



**Figure 3:** 95% posterior confidence intervals for correlations

# Intelligent imputations

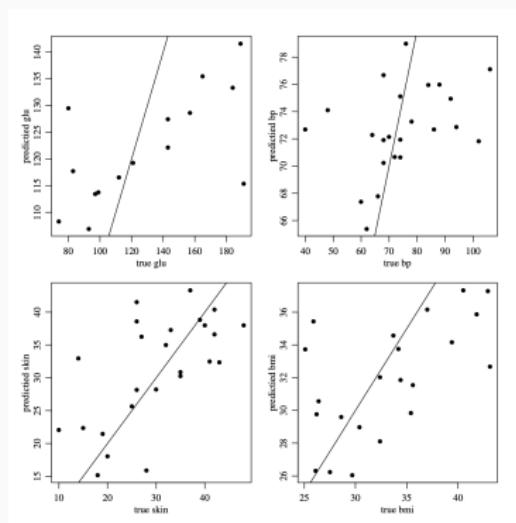
The posterior expectation gives a much better imputation than some flat fixed value.



**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.