

MT4531/5731: (Advanced) Bayesian Inference

Data augmentation - Auxiliary variables

Nicolò Margaritella

School of Mathematics and Statistics, University of St Andrews



University
of
St Andrews

Outline

- 1 Introduction
- 2 Example - mixture of normals
- 3 NIMBLE analysis

Outline

- 1 Introduction
- 2 Example - mixture of normals
- 3 NIMBLE analysis

Data augmentation (1)

- Data augmentation occurs when...
 - we have missing data, and we include them in the analysis by treating them as unknown random variables.
 - (Remember also that prediction can be put into a missing data framework, as we have seen in the Exponential-Gamma example of Lecture 12.)

Data augmentation (2)

- Data augmentation also occurs when...
 - The likelihood of the data is intractable (or difficult to calculate), but, ...
 - conditionally on a collection of unobserved data or parameters (auxiliary quantities), likelihood becomes tractable (or simpler) and,...
 - posterior calculations become simpler and quicker. We can define auxiliary variables in such a way that the resulting posterior conditional distributions are easier to sample from and/or result in better mixing.

Data augmentation (3)

- So, in simple words, data augmentation occurs when we include missing data in the analysis, or auxiliary quantities (usually model parameters) that are not an essential part of the model.
- Let \mathbf{x} denote the observed data and \mathbf{y} the auxiliary variables.
- Bayes' Theorem states that,

$$\pi(\boldsymbol{\theta}, \mathbf{y} | \mathbf{x}) \propto f(\mathbf{y}, \mathbf{x} | \boldsymbol{\theta}) p(\boldsymbol{\theta}).$$

- We use MCMC to sample from the posterior $\pi(\boldsymbol{\theta}, \mathbf{y} | \mathbf{x})$.
- However, we are interested in the posterior $\pi(\boldsymbol{\theta} | \mathbf{x})$. This is the marginal posterior,

$$\pi(\boldsymbol{\theta} | \mathbf{x}) = \int \pi(\boldsymbol{\theta}, \mathbf{y} | \mathbf{x}) d\mathbf{y}.$$

- We simply obtain a sample from the joint posterior distribution and only use the realisations of $\boldsymbol{\theta}$.

Outline

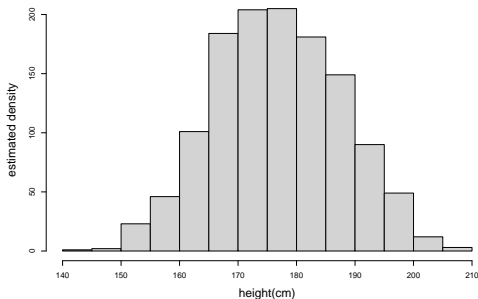
- 1 Introduction
- 2 Example - mixture of normals
- 3 NIMBLE analysis

Example - mixture of normals (1)

- Consider a data set on the height of 700 women and 550 men.
- Suppose we have the list of heights, but we don't know which data points are from women and which are from men.

Example - mixture of normals (2)

- From the plot can we still infer the distribution of female and male heights?



- The answer is yes!

Example - mixture of normals (3)

- We will assume that both mixture components (female and male) have the same precision s , which is fixed and known.
- The model is

$$\begin{aligned}X_1, \dots, X_n | \boldsymbol{\mu}, \pi &\sim g(x | \boldsymbol{\mu}, \pi); \\ \boldsymbol{\mu} \equiv (\mu_1, \mu_2) &\sim N(\boldsymbol{\phi}, \lambda^{-1}); \\ \pi &\sim \text{Beta}(a, b),\end{aligned}\tag{1}$$

where the pdf of $g(x | \boldsymbol{\mu}, \pi)$ is $\pi N(\mu_1, s^{-1}) + (1 - \pi)N(\mu_2, s^{-1})$

Example - mixture of normals (4)

- The likelihood is

$$\begin{aligned}f(x_1, \dots, x_n | \boldsymbol{\mu}, \pi) &= \prod_{i=1}^n g(x_i | \boldsymbol{\mu}, \pi); \\ &= \prod_{i=1}^n [\pi N(\mu_1, s^{-1}) + (1 - \pi)N(\mu_2, s^{-1})],\end{aligned}$$

which is a complicated function of $\boldsymbol{\mu}$ and π and it will produce a difficult posterior to sample from!

Example - mixture of normals (5)

- ⇒ Define an allocation variable $Z_{i=1,\dots,n}$ which identifies to what mixture component each data point comes from; i.e. Z_i identifies whether observation i is male or female.
- Using this auxiliary variable, we can rewrite our model as follows:

$$\begin{aligned}X_1, \dots, X_n | \mu, \pi &\sim N(\mu_{Z_i}, s^{-1}); \\Z_1, \dots, Z_n &\sim \text{Bernoulli}(\pi); \\(\mu_1, \mu_2) &\sim N(\phi, \lambda^{-1}); \\\pi &\sim \text{Beta}(a, b),\end{aligned}\tag{2}$$

which is the same model as (1)!

Example - mixture of normals (5)

- In fact we can rewrite (using what probability law?)

$$p(x_i|\mu, \pi) = \sum_{j=1}^2 p(x_i|Z_i = j, \mu, \pi)p(Z_i = j|\mu, \pi)$$

and then

$$\begin{aligned} p(x_i|\mu, \pi) &= p(x_i|Z_i = 1, \mu, \pi)p(Z_i = 1|\mu, \pi) + \\ &\quad + p(x_i|Z_i = 2, \mu, \pi)p(Z_i = 2|\mu, \pi), \\ &= N(\mu_1, s^{-1})\pi + N(\mu_2, s^{-1})(1 - \pi), \\ &= g(x_i|\boldsymbol{\mu}, \pi). \end{aligned}$$

Example - mixture of normals (6)

Let's look at the full conditionals now:

- (1) $p(\pi|\boldsymbol{\mu}, \mathbf{Z}, x)$ - Conditional on \mathbf{Z} , π is independent of the other parameters and the conditional distribution reduces to a Beta-Bernoulli model and we obtain

$$p(\pi|\boldsymbol{\mu}, \mathbf{Z}, x) = p(\pi|\mathbf{Z}) = \text{Beta}(a + n_1, b + n_2),$$

where $n_j = \sum_i \mathbb{1}(Z_i = j)$ with $j = \{1, 2\}$

Example - mixture of normals (7)

Let's look at the full conditionals now:

- (2) $p(\boldsymbol{\mu}|\mathbf{Z}, \pi, \mathbf{x})$ - Conditional on \mathbf{Z} , we know from which component each observation comes from! Therefore, conditionally on \mathbf{z} , the conditional distribution reduces to two independent Normal-Normal models.

$$\begin{aligned}\mu_1|\mu_2, \mathbf{Z}, \pi, \mathbf{x} &\sim \text{N}(M_1, L_1^{-1}); \\ \mu_2|\mu_1, \mathbf{Z}, \pi, \mathbf{x} &\sim \text{N}(M_2, L_2^{-1})\end{aligned}$$

where $n_j = \sum_i \mathbb{1}(Z_i = j)$, $L_j = \lambda + n_j s$,
 $M_j = \frac{\lambda \phi + s \sum_{i: Z_i = j} x_i}{\lambda + n_j s}$ and $j = \{1, 2\}$.

Example - mixture of normals (8)

Let's look at the full conditionals now:

(3) $p(\mathbf{Z}|\boldsymbol{\mu}, \pi, \mathbf{x})$ - First observe that

$$p(\mathbf{Z}|\boldsymbol{\mu}, \pi, \mathbf{x}) \propto p(\mathbf{Z}, \boldsymbol{\mu}, \pi, \mathbf{x}) \propto p(\mathbf{x}|\mathbf{Z}, \boldsymbol{\mu})p(\mathbf{Z}|\pi),$$

hence,

$$\begin{aligned} p(\mathbf{Z}|\boldsymbol{\mu}, \pi, \mathbf{x}) &= \prod_{i=1}^n \text{N}(x_i|\mu_{z_i}, s^{-1}) \text{Bernoulli}(z_i|\pi), \\ &= \prod_{i=1}^n (\text{N}(x_i|\mu_1, s^{-1})\pi)^{z_i} (\text{N}(x_i|\mu_2, s^{-1})(1-\pi))^{1-z_i}, \\ &\propto \prod_{i=1}^n \text{Bernoulli}\left(z_i \middle| \frac{\theta_{i,1}}{\theta_{i,1} + \theta_{i,2}}\right) \end{aligned}$$

with $\theta_{i,1} = \text{N}(x_i|\mu_1, s^{-1})\pi$ and $\theta_{i,2} = \text{N}(x_i|\mu_2, s^{-1})(1-\pi)$

Outline

- 1 Introduction
- 2 Example - mixture of normals
- 3 NIMBLE analysis**

Example - NIMBLE analysis (1)

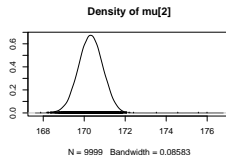
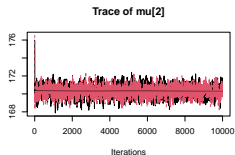
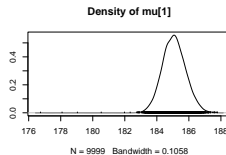
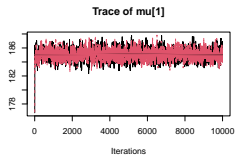
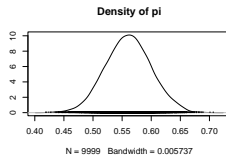
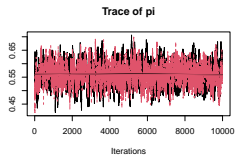
- Model settings:
 - $s = 1/\sigma^2$ where $\sigma = 8\text{cm}$ is the SD of subjects' height within each component.
 - $a = 1, b = 1$ for the Beta prior correspond to a Uniform prior on $[0,1]$.
 - $\phi = 177\text{cm}$ is the mean of the prior on the component means.
 - $1/\sqrt{\lambda} = 15\text{cm}$ is the SD of the prior on the component means.

Example - NIMBLE analysis (2)

- Let's take a look at the BUGS code for the 2-component mixture model:

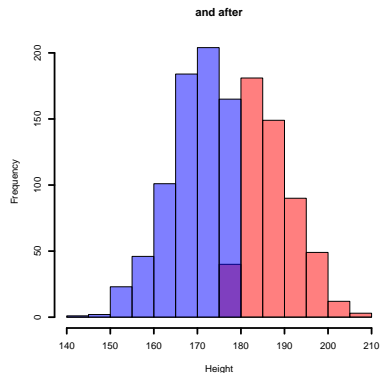
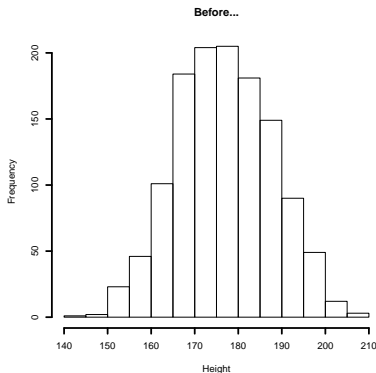
```
MixtureCode <- nimbleCode({  
  # Specify the likelihood:  
  for (i in 1:N){  
    x[i] ~ dnorm(mu[t[i]],s)  
    t[i] <- z[i]+1  
    z[i] ~ dbin(pi,1)# auxiliary variable  
  }  
  # Prior specification:  
  mu[1] ~ dnorm(phi,lambda)  
  mu[2] ~ dnorm(phi,lambda)  
  pi ~ dbeta(1,1)  
})
```

Example - NIMBLE analysis (3)



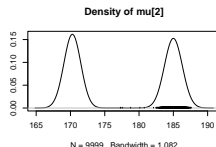
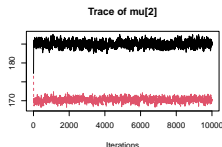
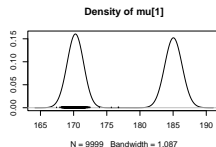
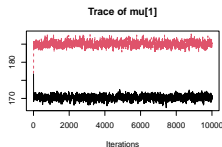
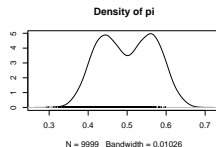
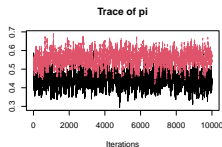
Example - NIMBLE analysis (4)

- Histograms of the heights of subjects assigned to each component, according to the posterior mode of z_j (right panel)



Watch out for multiple modes!

- If we run the sampler multiple times starting from the same initial values, sometimes it will settle on [females = 1, males = 2] and sometimes on [females = 2, males = 1].
- If the sampler were behaving properly, it would move back and forth between these two modes. But it gets stuck in one mode and stays there.
- Nevertheless, all parameters show good rate of convergence and mixing and we could just relabel the groups and then calculate the posterior summaries of interest



Try this at home

- **Task:** There is another interesting example of data augmentation in Section 2.5 of the lecture notes (Genetic Linkage) for you to read. This will also be useful to answer question 2 of tutorial 7.