

SPATIO-TEMPORAL METHODS IN ENVIRONMENTAL EPIDEMIOLOGY

Alexandra M. Schmidt

Session 2 - An introduction to Bayesian Inference

Department of Epidemiology, Biostatistics and Occupational Health
McGill University

alexandra.schmidt@mcgill.ca

Applied Bayesian Statistics School 2025

University of Genova, Department of Architecture and Design
03-06, June 2025

Bayesian inference

From <http://www.bayesian.org/>

"Scientific inquiry is an iterative process of integrating accumulating information. Investigators assess the current state of knowledge regarding the issue of interest, gather new data to address remaining questions, and then update and refine their understanding to incorporate both new and old data. Bayesian inference provides a logical, quantitative framework for this process. It has been applied in a multitude of scientific, technological, and policy settings."

Bayes' Theorem and methodology

Summarized presentation of methodology

Bayes Theorem

Observations y : described by density $f(y|\theta)$

Likelihood: $l(\theta) = f(y|\theta)$

θ : index of f (parameter)

Canonical situation: random sample

$y = (y_1, \dots, y_n)$ taken from $f(y|\theta)$.

Example

Measurements of a physical quantity θ with errors $e_i \sim N(0, \sigma^2)$, σ^2 known.

$y_i = \theta + e_i$, $i = 1, \dots, n$ and $f(y|\theta) =$

$$\prod_{i=1}^n f_N(y_i; \theta, \sigma^2) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{1}{2} \frac{(y_i - \theta)^2}{\sigma^2} \right\}$$

θ is more than simple index

- Situation repeats in more general cases
- Very likely that researcher has prior information about θ
- This may be modelled through density $p(\theta)$
- Lots of controversy in the past

Bayes' Theorem and methodology

Inference process based on distribution of θ after observing $y \rightarrow$
posterior distribution (as **opposed** to **prior**)

Obtained through Bayes theorem as

$$p(\theta|y) = \frac{p(\theta)f(y|\theta)}{f(y)} \quad \text{or}$$
$$\pi(\theta) \propto p(\theta)l(\theta)$$

$$\text{posterior dist.} \propto \text{prior dist.} \times \text{likelihood}$$

And

$$f(y) = \int f(y|\theta)p(\theta)d\theta$$

It is not generally necessary to compute the denominator.

Bayes' Theorem and methodology

Predictive Distribution

Prediction (or forecast) of a future observation y^* after observing y based on the distribution of $(y^*|y)$

$$f(y^*|y) = \int f(y^*, \theta|y) d\theta = \int f(y^*|\theta) p(\theta | y) d\theta$$

if y and y^* are conditionally independent given $\theta \rightarrow$ eg. random sample

The main conjugate families

Theorem 1

Let $\theta \sim N(\mu, \tau^2)$ and $X | \theta \sim N(\theta, \sigma^2)$, with σ^2 known.

The posterior distribution of θ is $(\theta | X = x) \sim N(\mu_1, \tau_1^2)$ where

$$\mu_1 = \frac{\tau^{-2}\mu + \sigma^{-2}x}{\tau^{-2} + \sigma^{-2}} \quad \text{e} \quad \tau_1^{-2} = \tau^{-2} + \sigma^{-2} \quad (1)$$

- The precision is the inverse of the variance
- From the result above, the posterior precision of μ is the sum between the prior precision and the likelihood, and it does not depend on x

Sampling from the normal distribution with known variance

- Looking at the precision as a measure on the amount of information, and defining $w = \tau^{-2}/(\tau^{-2} + \sigma^{-2}) \in (0, 1)$
- w measures the relative information contained in the prior distribution with respect to the total information (prior + likelihood) Then we can write

$$\mu_1 = w\mu + (1 - w)x$$

which is a weighted average between the prior mean and the likelihood mean

Proving the previous theorem

Proof : From Bayes theorem

$$\begin{aligned} p(\theta | x) &\propto l(\theta; x)p(\theta) \\ &\propto \exp \left\{ -\frac{1}{2\sigma^2}(x - \theta)^2 - \frac{1}{2\tau^2}(\theta - \mu)^2 \right\} \\ &\propto \exp \left\{ -\frac{\theta^2}{2\sigma^2} - \frac{\theta^2}{2\tau^2} + \frac{x\theta}{\sigma^2} + \frac{\mu\theta}{\tau^2} \right\} \\ &= \exp \left\{ -\frac{\theta^2}{2} \left(\frac{1}{\sigma^2} + \frac{1}{\tau^2} \right) + \theta \left(\frac{x}{\sigma^2} + \frac{\mu}{\tau^2} \right) \right\} \end{aligned}$$

where, in the first step, all the constants were included in the proportionality constant

Now let $\tau_1^2 = (\tau^{-2} + \sigma^{-2})^{-1}$ and $\mu_1 = (\sigma^{-2}x + \mu\tau^{-2})\tau_1^2$

Proving the previous theorem (cont.)

Substituting the expressions in the previous slide, we have

$$\begin{aligned} p(\theta | x) &\propto \exp \left\{ -\frac{\theta^2}{2\tau_1^2} + \frac{\theta\mu_1}{\tau_1^2} \right\} \\ &\propto \exp \left\{ -\frac{1}{2\tau_1^2}(\theta - \mu_1)^2 \right\} \\ &\propto \frac{1}{\sqrt{2\pi\tau_1^2}} \exp \left\{ -\frac{1}{2\tau_1^2}(\theta - \mu_1)^2 \right\} \end{aligned}$$

Note that the last term in the expression above corresponds to that of a normal density

Therefore, the constant of proportionality is equal to 1 and $(\theta | x) \sim N(\mu_1, \tau_1^2)$.

The main conjugate families

Binomial Distribution

The family of beta distributions is conjugate to the binomial (or Bernoulli) model.

Normal with known variance

Theorem 1 stated that the normal family is conjugate to the normal model. For a sample of size n we have

$$l(\theta; \mathbf{x}) \propto \exp \left\{ -\frac{n}{2\sigma^2} (\bar{X} - \theta)^2 \right\}$$

Note that $p(\mathbf{x} | \theta) \propto p(\bar{X} | \theta)$. Assuming $\theta \sim N(\mu, \tau^2)$ we have that $\theta | \mathbf{x} \sim N(\mu_1, \tau_1^2)$ where

$$\mu_1 = \frac{n\sigma^{-2}\bar{X} + \tau^{-2}\mu}{n\sigma^{-2} + \tau^{-2}} \text{ and } \tau_1^{-2} = n\sigma^{-2} + \tau^{-2}.$$

The main conjugate families

Poisson distribution

Suppose that $\mathbf{X} = (X_1, \dots, X_n)$ is a random sample from the Poisson distribution with parameter θ , denoted $Pois(\theta)$. Its joint probability function is

$$p(\mathbf{x} | \theta) = \prod_{i=1}^n p(x_i | \theta) = \prod_{i=1}^n \frac{e^{-\theta} \theta^{x_i}}{x_i!}.$$

The Gamma family of distributions is conjugate to the Poisson model.

Exponential distribution

Suppose that $\mathbf{X} = (X_1, \dots, X_n)$ is a random sample from the Exponential distribution with parameter θ , denoted $Exp(\theta)$. Its joint probability function is

$$p(\mathbf{x} | \theta) = \prod_{i=1}^n p(x_i | \theta) = \prod_{i=1}^n \theta e^{-\theta x_i}.$$

The Gamma family of distributions is conjugate to the Exponential model.

The main conjugate families

Multinomial distribution

Let $\mathbf{X} = (X_1, \dots, X_p)$ and $\theta = (\theta_1, \dots, \theta_p)$ be, respectively, the number of observed cases and the probabilities associated with each of p categories in a sample of size n . Assume that $\sum_{i=1}^n X_i = n$ and $\sum_{i=1}^p \theta_i = 1$. \mathbf{X} is said to have a multinomial distribution with parameters n and $(\theta_1, \dots, \theta_p)$. And

$$p(\mathbf{x} | \theta) = \frac{n!}{\prod_{i=1}^p x_i!} \prod_{i=1}^p \theta_i^{x_i}.$$

It also belongs to the exponential family and $l(\theta) \propto \prod \theta_i^{x_i}$. Its kernel is the same as the kernel of the density of a Dirichlet distribution. The Dirichlet family with integer parameters a_1, \dots, a_p is natural conjugate with respect to the multinomial sampling distribution. Again, little is lost by extending natural conjugacy over all Dirichlet distributions.

The main conjugate families

The posterior will then be

$$p(\theta \mid \mathbf{x}) \propto \left[\prod_{i=1}^p \theta_i^{x_i} \right] \left[\prod_{i=1}^p \theta_i^{a_i-1} \right] = \prod_{i=1}^p \theta_i^{x_i+a_i-1},$$

which is a Dirichlet with parameters $a_1 + x_1, \dots, a_p + x_p$ and is denoted by $(\theta \mid \mathbf{x}) \sim D(a_1 + x_1, \dots, a_p + x_p)$.

The proportionality constant is given by

$$\frac{\Gamma(a+n)}{\prod_{i=1}^p \Gamma(a_i + x_i)}$$

This conjugate analysis generalizes the analysis for binomial samples with beta priors.

The main conjugate families

Normal with known mean and unknown variance

Let $\mathbf{X} = (X_1, \dots, X_p)$ be a random sample from the $N(\theta, \sigma^2)$, θ known, and $\phi = \sigma^{-2}$. In this case we have that

$$l(\phi; \mathbf{x}) = p(\mathbf{x} | \theta, \phi) \propto \phi^{n/2} \exp \left\{ -\frac{\phi}{2} n s_0^2 \right\} \text{ where}$$

$$s_0^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \theta)^2.$$

The conjugate prior may have the kernel of $l(\phi; \mathbf{x})$, which is in the gamma distribution form. As the gamma family is closed under sampling we can consider $\phi \sim Ga(n_0/2, n_0 \sigma_0^2/2)$ or, equivalently, $n_0 \sigma_0^2 \phi \sim \chi_{n_0}^2$. Then

$$\phi | \mathbf{x} \sim Ga \left(\frac{n_0 + n}{2}, \frac{n_0 \sigma_0^2 + n s_0^2}{2} \right),$$

or, equivalently, $(n_0 \sigma_0^2 + n s_0^2) \phi | \mathbf{x} \sim \chi_{n_0 + n}^2$.

The main conjugate families

Normal with unknown mean and variance

Assume that (θ, ϕ) be such that $X_i \sim N(\theta, \sigma^2)$, where $\phi = \sigma^{-2}$. We will consider the prior for (θ, ϕ) in two stages:

$$(\theta \mid \phi) \sim N(\mu_0, (c_0\phi)^{-1}) \text{ and} \\ n_0\sigma_0^2\phi \sim \chi_{n_0}^2 \text{ or } \phi \sim Ga\left(n_0/2, n_0\sigma_0^2/2\right),$$

where $(\mu_0, c_0, n_0, \sigma_0^2)$ are obtained from the initial information H . This distribution is usually called normal-gamma or normal- χ^2 with parameters $(\mu_0, c_0, n_0, \sigma_0^2)$ and joint density given by:

$$p(\theta, \phi) \propto \phi^{(n_0+1)/2-1} \exp\left\{-\frac{\phi}{2} \left[n_0\sigma_0^2 + c_0(\theta - \mu_0)^2\right]\right\}$$

Note that

$$\begin{aligned} p(\theta) &= \int_0^\infty \phi^{a-1} e^{-b\phi} d\phi = \frac{\Gamma(a)}{b^a} \\ &\propto [n_0\sigma_0^2 + c_0(\theta - \mu_0)^2]^{-(n_0+1)/2} \end{aligned}$$

The main conjugate families

Rearranging the terms of the last equation gives

$$p(\theta) \propto \left[1 + \frac{(\theta - \mu_0)^2}{n_0(\sigma_0^2/c_0)} \right]^{-(n_0+1)/2}$$

implying that $\theta \sim t_{n_0}(\mu_0, \sigma_0^2/c_0)$.

The conditional distribution of $\phi \mid \theta$ can be obtained from the joint distribution of (θ, ϕ) and

$$\phi \mid \theta \sim Ga\left(\frac{(n_0+1)}{2}, \frac{[n_0\sigma_0^2 + c_0(\theta - \mu_0)^2]}{2}\right)$$

The joint distribution of a random sample $\mathbf{X} = (X_1, \dots, X_n)$ is

$$\begin{aligned} p(\mathbf{x} \mid \theta, \phi) &= \prod_{i=1}^n \phi^{1/2} \exp\left\{-\frac{\phi}{2}(x_i - \theta)^2\right\} \\ &\propto \phi^{n/2} \exp\left\{-\frac{\phi}{2}\left[ns^2 + n(\bar{x} - \theta)^2\right]\right\} \end{aligned}$$

where $s^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$.

It has the same kernel as the normal-gamma density for (θ, ϕ) . The normal-gamma is closed under sampling.

The main conjugate families

The posterior distribution will then be

$$\begin{aligned} p(\theta, \phi \mid \mathbf{x}) &\propto p(\mathbf{x} \mid \theta, \phi) p(\theta, \phi) \\ &\propto \phi^{[(n+n_0+1)/2]-1} \\ &\quad \times \exp \left\{ -\frac{\phi}{2} [n_0 \sigma_0^2 + ns^2 + c_0(\theta - \mu_0)^2 + n(\bar{x} - \theta)^2] \right\}. \end{aligned}$$

But

$$c_0(\theta - \mu_0)^2 + n(\theta - \bar{x})^2 = (c_0 + n)(\theta - \mu_1)^2 + \frac{c_0 n}{c_0 + n}(\mu_0 - \bar{x})^2$$

where $\mu_1 = (c_0 \mu_0 + n \bar{x}) / (c_0 + n)$. Then we have that

$$\begin{aligned} p(\theta, \phi) &\propto \phi^{[(n+n_0+1)/2]-1} \\ &\exp \left\{ -\frac{\phi}{2} \left[n_0 \sigma_0^2 + ns^2 + \frac{c_0 n}{c_0 + n}(\mu_0 - \bar{x})^2 + (c_0 + n)(\theta - \mu_1)^2 \right] \right\} \end{aligned}$$

The main conjugate families

The joint posterior for $(\theta, \phi \mid \mathbf{x})$ is normal-gamma with parameters $(\mu_1, c_1, n_1, \sigma_1^2)$ given by

$$\begin{aligned}\mu_1 &= \frac{c_0 \mu_0 + n \bar{x}}{c_0 + n} & c_1 &= c_0 + n \\ n_1 &= n_0 + n & n_1 \sigma_1^2 &= n_0 \sigma_0^2 + n s^2 + \frac{c_0 n}{c_0 + n} (\mu_0 - \bar{x})^2.\end{aligned}$$

The normal-gamma family is conjugate with respect to the normal sampling model when θ and σ^2 are both unknown.

The main conjugate families

Summary of the Distributions in the normal case

	Prior	Posterior
$\theta \mid \phi$	$N(\mu_0, (c_0 \phi)^{-1})$	$N(\mu_1, (c_1 \phi)^{-1})$
ϕ	$n_0 \sigma_0^2 \phi \sim \chi_{n_0}^2$	$n_1 \sigma_1^2 \phi \sim \chi_{n_1}^2$
θ	$t_{n_0}(\mu_0, \sigma_0^2/c_0)$	$t_{n_1}(\mu_1, \sigma_1^2/c_1)$
$\phi \mid \theta$	$[n_0 \sigma_0^2 + c_0(\theta - \mu_0)^2] \phi \sim \chi_{n_0+1}^2$	$[n_1 \sigma_1^2 + c_1(\theta - \mu_1)^2] \phi \sim \chi_{n_1+1}^2$

Example (1)

Normal data, $Y_i \sim N(\theta, \sigma^2)$, σ^2 known, then

$$\begin{aligned} l(\theta) &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{1}{2} \frac{(y_i - \theta)^2}{\sigma^2} \right\} \\ &\propto \exp \left\{ -\frac{n}{2\sigma^2} (\bar{y} - \theta)^2 \right\} \end{aligned}$$

where \bar{y} is the average of the y_i 's.

Model completed with prior for θ ,

$p(\theta) = N(\mu, \tau^2)$, μ and τ^2 known.

Then, we have that

$$\pi(\theta) \sim N(\mu_1, \tau_1^2)$$

where

$\tau_1^{-2} = n\sigma^{-2} + \tau^{-2}$ and $\mu_1 = \tau_1^2(n\sigma^{-2}\bar{y} + \tau^{-2}\mu)$

$\tau^2 \rightarrow \infty$: non-informative prior $p(\theta) \propto c$ and

$\pi(\theta) = N(\bar{y}, \sigma^2/n)$.

[See ConjugatePriors.r](#)

Bayes Estimators

A Bayes estimator is an estimator that is chosen to minimize the posterior mean of some measure of how far the estimator is from the parameter.

Loss Function A loss function is a real-valued function of two variables $L(\theta, a)$, where $\theta \in \Omega$ and a is a real number. The interpretation is that the statistician loses $L(\theta, a)$ if the parameter equals θ and the estimate equals a

Definition of a Bayes Estimator

Suppose that one can observe the value \mathbf{y} of random vector \mathbf{Y} before estimating θ , and let $p(\theta | \mathbf{x})$ denote the posterior pdf of θ on Ω . For each estimate a that the statistician might use, her expected loss in this case will be

$$E[L(\theta, a) | \mathbf{x}] = \int_{\Omega} L(\theta, a) p(\theta | \mathbf{x}) d\theta \quad (2)$$

You can choose an estimate a for which the expectation above is a minimum.

Bayes Estimators

Definition Let $L(\theta, a)$ be a loss function. For each possible value \mathbf{x} of \mathbf{X} , let $\delta^*(\mathbf{x})$ be a value of a such that $E[L(\theta, a) | \mathbf{x}]$ is minimized. Then δ^* is called a *Bayes estimator* of θ . Once $\mathbf{X} = \mathbf{x}$ is observed, $\delta^*(\mathbf{x})$ is called a *Bayes estimate* of θ .

Then, for each possible value of \mathbf{x} of \mathbf{X} , the value $\delta^*(\mathbf{x})$ is chosen so that

$$E[L(\theta, \delta^*(\mathbf{x})) | \mathbf{x}] = \min_{\text{All } a} E[L(\theta, a) | \mathbf{x}]$$

Corollary Let θ be a real-valued parameter. Suppose that the squared error loss function, $L(\theta, a) = (\theta - a)^2$ is used, and that $E(\theta | \mathbf{x}) < \infty$. Then, a Bayes estimator of θ is $\delta^*(\mathbf{x}) = E(\theta | \mathbf{X})$.

Corollary Let θ be a real-valued parameter. Suppose that the absolute error loss function, $L(\theta, a) = |\theta - a|$ is used, then a Bayes estimator of θ is equal to the median of the posterior distribution of θ .

Bayes Estimators

Another form to reduce the effect of large estimation errors is to consider loss functions that remain constant whenever $|\delta - \theta| > k$ for some k arbitrary.

⇒ The most common choice is the limiting value as $k \rightarrow 0$.

This loss function associates a fixed loss when an error is committed, irrespective of its magnitude. This loss is usually known as the **0-1 loss**.

Lemma: Let $L_3(\delta, \theta) = \lim_{\varepsilon \rightarrow 0} I_{|\theta - \delta|}([\varepsilon, \infty))$. The estimator of θ is $\delta_3 = \text{mode}(\theta)$, the mode of the posterior distribution of θ .

→ also known as the generalized maximum likelihood estimator (GMLE).

Consistency of Bayes Estimators

Under fairly general conditions, and for a wide class of loss functions, the Bayes estimators of some parameters θ will form a consistent sequence of estimators as the sample size $n \rightarrow \infty$.

Linear Models

Model specification

- $\mathbf{Y} \mid \beta, \sigma^2 \sim N_n(\mathbf{X}\beta, \sigma^2 \mathbf{I}_n)$
 - $\mathbf{Y}_{n \times 1}$ vector of observations
 - $\mathbf{X}_{n \times p}$ design matrix
 - $\beta_{p \times 1}$ parameter vector
- $E(\mathbf{Y} \mid \mathbf{X}) = \mathbf{X}\beta$.
- Kernel of the likelihood function

$$l(\beta, \sigma^2; \mathbf{y}) \propto \sigma^{-n} \exp \left\{ -\frac{S(\beta)}{2\sigma^2} \right\},$$
$$S(\beta) = (\mathbf{y} - \mathbf{X}\beta)'(\mathbf{y} - \mathbf{X}\beta).$$

$$\text{MLE} : \hat{\beta} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}.$$

Bayesian linear models

Bayesian inference

- Model: $\mathbf{y} \mid \beta, \sigma^2 \sim N_n(\mathbf{X}\beta, \sigma^2 \mathbf{I}_n)$.
- Prior specification for β and $\phi = \sigma^{-2}$.

Conjugate Prior

- $$\begin{cases} \beta \mid \phi \sim N_p(\mu_0, \phi^{-1} \mathbf{C}_0^{-1}), \phi \mathbf{C}_0 \text{ precision} \\ \phi \sim Ga\left(\frac{n_0}{2}; \frac{n_0 \sigma_0^2}{2}\right) \end{cases}$$
- Joint prior: $p(\beta, \phi) \sim \text{Normal} \times \text{Gamma}$;
- Marginal prior for β : $p(\beta) = \frac{p(\beta, \phi)}{p(\phi \mid \beta)}(*)$
- From the joint prior $p(\beta, \phi)$ it follows that
 - ▶ $\phi \mid \beta \sim Ga\left(\frac{n_0+p}{2}, \frac{(n_0+p)[n_0 \sigma_0^2 + (\beta - \mu_0)' \mathbf{C}_0 (\beta - \mu_0)]}{2}\right)$
 - ▶ Using (*) or integrating $p(\beta, \phi)$ we have that
$$p(\beta) \propto [n_0 \sigma_0^2 + (\beta - \mu_0)' \mathbf{C}_0 (\beta - \mu_0)]^{-(n_0+p)/2},$$
$$\Rightarrow \beta \sim tM_{n_0}(\mu_0, \sigma_0^2 \mathbf{C}_0^{-1}).$$

Bayesian linear models

- On the other hand, the kernel of $l(\beta, \phi; \mathbf{y})$, evaluated at the MLE $\hat{\beta}$, given by

$$\phi^{n/2} \exp \left\{ -\frac{\phi}{2} [\mathbf{S}_e + (\beta - \hat{\beta})' \mathbf{X}' \mathbf{X} (\beta - \hat{\beta})] \right\}$$

has the same form of the prior distribution

- It can be shown that

$$\begin{aligned} (\beta - \mu_0)' \mathbf{C}_0 (\beta - \mu_0) + (\beta - \hat{\beta})' \mathbf{X}' \mathbf{X} (\beta - \hat{\beta}) \\ = (\beta - \mu_1)' \mathbf{C}_1 (\beta - \mu_1) + \mu_0' \mathbf{C}_0 \mu_0 + \hat{\beta}' \mathbf{X}' \mathbf{X} \hat{\beta} + \mu_1' \mathbf{C}_1 \mu_1 \end{aligned}$$

$$\mu_1 = \mathbf{C}_1^{-1} (\mathbf{C}_0 \mu_0 + \mathbf{X}' \mathbf{y}) \text{ e } \mathbf{C}_1 = \mathbf{C}_0 + \mathbf{X}' \mathbf{X}.$$

- Also

$$\begin{aligned} \mathbf{S}_e &+ \mu_0' \mathbf{C}_0 \mu_0 + \hat{\beta}' \mathbf{X}' \mathbf{X} \hat{\beta} + \mu_1' \mathbf{C}_1 \mu_1 \\ &= \mathbf{y}' \mathbf{y} - \hat{\beta}' \mathbf{X}' \mathbf{X} \hat{\beta} + \mu_0' \mathbf{C}_0 \mu_0 + \hat{\beta}' \mathbf{X}' \mathbf{X} \hat{\beta} + \mu_1' (\mathbf{C}_0 \mu_0 + \mathbf{X}' \mathbf{y}) \\ &= (\mathbf{y} - \mathbf{X} \mu_1)' \mathbf{y} + (\mu_0 - \mu_1)' \mathbf{C}_0 \mu_0. \end{aligned}$$

Bayesian linear models

Conjugate posterior distribution

- $$p(\beta, \phi | \mathbf{y}) \propto \phi^{p/2} \exp \left\{ -\frac{\phi}{2} (\beta - \mu_1)' \mathbf{C}_1 (\beta - \mu_1) \right\} \\ \times \phi^{(n_1/2)-1} \exp \left\{ -\frac{\phi}{2} n_1 \sigma_1^2 \right\},$$

$$n_1 = n + n_0$$

$$n_1 \sigma_1^2 = n_0 \sigma_0^2 + (\mathbf{y} - \mathbf{X} \mu_1)' \mathbf{y} - (\mu_1 - \mu_0)' \mathbf{C}_0 \mu_0.$$

- Then $p(\beta, \phi | \mathbf{y}) \sim \text{Normal} \times \text{Gamma}$

$$\begin{cases} \beta | \phi, \mathbf{y} \sim N_p(\mu_1, (\phi \mathbf{C}_1)^{-1}), \\ \phi | \mathbf{y} \sim Ga\left(\frac{n_1}{2}; \frac{n_1 \sigma_1^2}{2}\right) \end{cases}$$

- Marginal posterior: $\beta | \mathbf{y} \sim tM_{n_1}(\mu_1, \sigma_1^2 \mathbf{C}_1^{-1})$,

$$\begin{cases} E(\beta | \mathbf{y}) = \mu_1 = \mathbf{C}_1^{-1} (\mathbf{C}_0 \mu_0 + \mathbf{X}' \mathbf{y}), \\ \text{Var}(\beta | \mathbf{y}) = \frac{n_1}{n_1 - 2} \sigma_1^2 \mathbf{C}_1^{-1}, n_1 > 2. \end{cases}$$

Bayesian linear models

- Marginal Posterior : $\phi|\mathbf{y} \sim Ga\left(\frac{n_1}{2}, \frac{n_1\sigma_1^2}{2}\right)$,

$$\left\{ \begin{array}{l} E(\phi|\mathbf{y}) = \phi_1, \\ Var(\phi|\mathbf{y}) = \frac{2\phi_1^2}{n_1} \end{array} \right.$$

- Credible intervals for β_j and ϕ : are obtained through the percentiles of the t_{n_1} and $Ga\left(\frac{n_1}{2}, \frac{n_1\sigma_1^2}{2}\right)$, respectively.
- Inference about β :

$$(\beta - \mu_1)' \mathbf{C}_1 (\beta - \mu_1) \sigma_1^2 \mid \mathbf{y} \sim F(p, n_1 - p).$$

Bayesian linear models

Jeffreys prior:

$$p(\beta, \phi) \propto \phi^{-1} \Leftrightarrow \mu_0 \equiv 0, \sigma_0 \rightarrow 0 \text{ e } C_0 \rightarrow 0.$$

- Posterior density Normal \times Gamma: $p(\beta, \phi \mid \mathbf{y})$

$$\propto \phi^{(n/2)-1} \exp \left\{ -\frac{\phi}{2} [S_e + (\beta - \hat{\beta})' \mathbf{X}' \mathbf{X} (\beta - \hat{\beta})] \right\}$$

$$\propto \phi^{p/2} \exp \left\{ -\frac{\phi}{2} (\beta - \hat{\beta}) \mathbf{X}' \mathbf{X} (\beta - \hat{\beta}) \right\} \\ \times \phi^{((n-p)/2)-1} \exp \left\{ -\frac{\phi}{2} (n-p) s^2 \right\}.$$

- Then, $\beta \mid \mathbf{y} \sim tM_{n-p}(\hat{\beta}, s^2(\mathbf{X}'\mathbf{X})^{-1})$,

$$\phi \mid \mathbf{y} \sim Ga \left(\frac{n-p}{2}, \frac{(n-p)s^2}{2} \right) \text{ and}$$

$$(\beta - \hat{\beta})' \mathbf{X}' \mathbf{X} (\beta - \hat{\beta}) / s^2 \sim F(p, n-p).$$

- Bayesian and classical results are similar

Bayesian linear models - Example

Multiple linear regression: analysis of heart data

See files `Heart_data.pdf` and `HeartBayesian.r`

Bayesian Estimation Methods

- We can usually write down the functional form of the posterior distributions required for Bayesian inference.
- For most realistic models, these functions are complex and high dimensional: Difficult to evaluate them analytically.
- Instead, it is often straightforward to *simulate* realisations from the required posterior distributions.
- Posterior summaries (e.g. posterior mean) are easily obtained by simple data summaries of the simulated values.
- Markov Chain Monte Carlo (MCMC) methods are a convenient class of simulation algorithms for this purpose.

Heuristic View of Simulation Methods for Bayesian Inference

- Imagine generating a random sample of values from a probability distribution (e.g.: normal);
- Construct a histogram from the sample;
- If the sample is large enough, histogram can provide virtually complete information about the distribution from which these samples were drawn:
 - ▶ Mean, variance, percentiles of sample
 \approx mean, variance, percentiles original distribution.
- MCMC methods enable us to generate large samples from the posterior distributions of model parameters:
 - ▶ These samples can be summarised to estimate properties (e.g. mean, variance, percentiles) of the posterior distribution.

Markov chain Monte Carlo Methods

A Markov chain Monte Carlo algorithm to simulate from $\pi(\cdot)$ is any method which produces a homogeneous, ergodic and irreducible Markov chain, whose stationary distribution is $\pi(\cdot)$;

A chain is **ergodic** if it is (i) aperiodic and (ii) positively recurrent;

Periodicity: a chain is aperiodic if none of its states is visited after d steps with probability one, for any d integer and $d > 0$.

Positive Recurrence: a chain is positively recurrent when the mean number of steps for the chain to return to any state is finite.

Markov chain Monte Carlo Methods

A chain is irreducible if, with positive probability, it moves from one (any) point to another in a finite number of iterations.

Results:

- If the Markov chain is homogeneous, irreducible, positively recurrent and aperiodic, then the limit distribution exists and the states of this chain are, approximately, realizations of this stationary distribution.
- Ergodic Means:

$$\bar{h}_m = \frac{1}{m} \sum_{i=1}^m h(\theta_i) \rightarrow E_{\pi}[h(\theta)]$$

for $m \rightarrow \infty$.

This result is similar to Monte Carlo integration.

Gibbs Sampling

Geman and Geman (1984), Gelfand and Smith (1990)

It is an algorithm that generates a sequence

$$\{\theta^{(0)}, \theta^{(1)}, \theta^{(2)}, \dots\},$$

from a Markov chain, whose limit/equilibrium distribution is $\pi(\theta)$ and whose transition kernel is given by the product of the full conditional distributions.

Algorithm:

1. $\theta^{(0)} = (\theta_1^{(0)}, \theta_2^{(0)}, \dots, \theta_p^{(0)})$
2. $\theta^{(j)}$ obtained from $\theta^{(j-1)}$

$$\begin{aligned}\theta_1^{(j)} &\sim \pi(\theta_1 | \theta_2^{(j-1)}, \dots, \theta_p^{(j-1)}) \\ \theta_2^{(j)} &\sim \pi(\theta_2 | \theta_1^{(j)}, \theta_3^{(j-1)}, \dots, \theta_p^{(j-1)}) \\ \theta_3^{(j)} &\sim \pi(\theta_3 | \theta_1^{(j)}, \theta_2^{(j)}, \theta_4^{(j-1)}, \dots, \theta_p^{(j-1)}) \\ &\vdots \\ \theta_p^{(j)} &\sim \pi(\theta_p | \theta_1^{(j)}, \theta_2^{(j)}, \dots, \theta_{p-1}^{(j)})\end{aligned}$$

Example

Assume that $\mathbf{X} \sim N_2(\mu, \Sigma)$, where

$$\mu = \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix} \Sigma = \begin{pmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{21} & \sigma_2^2 \end{pmatrix}$$

It is also known that for $i \neq j = 1, 2$,

$$\begin{aligned} (X_i | X_j = x_j) &\sim N(\mu_{i|j}, \sigma_{i|j}^2) \\ \mu_{i|j} &= \mu_i + \sigma_{ij} \sigma_j^{-2} (x_j - \mu_j) \\ \sigma_{i|j}^2 &= \sigma_i^2 - \sigma_{ij}^2 \sigma_j^{-2} \end{aligned}$$

Let

$$\mu = \begin{pmatrix} 2 \\ 1 \end{pmatrix} \Sigma = \begin{pmatrix} 1 & 0.7 \\ 0.7 & 1 \end{pmatrix}$$

See file `GibbsExample.r`

Change of mean in the Poisson

Assume that

$$y_i | \lambda, \phi, k \sim \begin{cases} Po(\lambda) & \text{para } i = 1, 2, \dots, k \\ Po(\phi) & \text{para } i = k+1, k+2, \dots, n \end{cases}$$

such that,

$$p(\mathbf{y} | \lambda, \phi, k) = \left[\prod_{i=1}^k \frac{\lambda^{y_i} e^{-\lambda}}{y_i!} \right] \left[\prod_{i=k+1}^n \frac{\phi^{y_i} e^{-\phi}}{y_i!} \right]$$

and likelihood function given by:

$$l(\lambda, \phi, k | \mathbf{y}) \propto \lambda^{t_1(\mathbf{y}, k)} e^{-k\lambda} \phi^{t_2(\mathbf{y}, k)} e^{-(n-k)\phi}$$

where

$$t_1(\mathbf{y}, k) = \sum_{i=1}^k y_i \quad t_2(\mathbf{y}, k) = \sum_{i=k+1}^n y_i$$

Carlin, Gelfand and Smith (1992) apply this model to the British coalmining disaster data.

Prior Distributions:

$$\pi(\lambda|\alpha, \beta) \sim \text{Ga}(\alpha, \beta)$$

$$\pi(\phi|\gamma, \delta) \sim \text{Ga}(\gamma, \delta)$$

$$\text{Pr}(k = i) = 1/n$$

Posterior Full Conditional Distributions

$$\pi(\lambda|\mathbf{y}, k, \phi) \sim \text{Ga}(\alpha + t_1(\mathbf{y}, k), \beta + k)$$

$$\pi(\phi|\mathbf{y}, k, \lambda) \sim \text{Ga}(\gamma + t_2(\mathbf{y}, k), \delta + n - k)$$

$$\text{Pr}(k = i|\mathbf{y}, \phi, \lambda) \propto \lambda^{t_1(\mathbf{y}, i)} \phi^{t_2(\mathbf{y}, i)} e^{-(i\lambda + (n-i)\phi)}$$

Metropolis-Hastings

Metropolis et. al. (1953), Hastings (1970)

Like the Gibbs sampling, it is an algorithm that generates a sequence

$$\{\theta^{(0)}, \theta^{(1)}, \theta^{(2)}, \dots\},$$

from a Markov chain whose stationary distribution is $\pi(\theta)$.

1. initial value $\theta^{(0)}$.
2. proposed value $\xi \sim q(\xi|\theta^{(i-1)})$
3. accepted value:

$$\theta^{(i)} = \begin{cases} \xi & \text{with probability } \alpha \\ \theta^{(i-1)} & \text{with probability } 1 - \alpha \end{cases}$$

where

$$\alpha = \min \left\{ 1, \frac{\frac{\pi(\xi)}{q(\xi|\theta^{(i-1)})}}{\frac{\pi(\theta^{(i-1)})}{q(\theta^{(i-1)}|\xi)}} \right\}$$

Special Cases

1. Symmetric Chains: $q(\theta|\xi) = q(\xi|\theta)$

$$\alpha = \min \left\{ 1, \frac{\pi(\xi)}{\pi(\theta)} \right\}$$

2. Random Walk: $q(\theta|\xi) = q(|\theta - \xi|)$

$$\alpha = \min \left\{ 1, \frac{\pi(\xi)}{\pi(\theta)} \right\}$$

3. Independent Chains: $q(\theta|\xi) = q(\theta)$

$$\alpha = \min \left\{ 1, \frac{\omega(\xi)}{\omega(\theta)} \right\}$$

where $\omega(\xi) = \pi(\xi)/q(\xi)$.

Gibbs \subset Metropolis

The Gibbs sampling is equivalent to a composition of p Metropolis-Hastings algorithms, whose acceptance probabilities are always equals 1.

$$\frac{\pi(\theta_1^{(j)}, \dots, \theta_{i-1}^{(j)}, \xi, \theta_{i+1}^{(j-1)}, \dots, \theta_p^{(j-1)})}{\pi(\xi | \theta_1^{(j)}, \dots, \theta_{i-1}^{(j)}, \theta_{i+1}^{(j-1)}, \dots, \theta_p^{(j-1)})}$$

equals

$$\pi(\theta_1^{(j)}, \dots, \theta_{i-1}^{(j)}, \theta_{i+1}^{(j-1)}, \dots, \theta_p^{(j-1)})$$

Example: Mixture of Normals

Let,

$$\mathbf{X} \sim 0.7N(\mu_1, \Sigma_1) + 0.3N(\mu_2, \Sigma_2)$$

where

$$\mu_1 = \begin{pmatrix} 4 \\ 5 \end{pmatrix} \quad \mu_2 = \begin{pmatrix} 1 \\ 4 \end{pmatrix}$$

and

$$\Sigma_1 = \begin{pmatrix} 1 & 0.7 \\ 0.7 & 1 \end{pmatrix} \quad \Sigma_2 = \begin{pmatrix} 1 & -0.7 \\ -0.7 & 1 \end{pmatrix}$$

Assume one wants to use

$$q(\mathbf{X}^{(i)} | \mathbf{X}^{(i-1)}) \sim N(\mathbf{X}^{(i-1)}, \mathbf{V})$$

as the proposal of the Metropolis-Hastings algorithm.

JAGS

Just Another Gibbs Sampler

<http://mcmc-jags.sourceforge.net/>

It is a program for analysis of Bayesian hierarchical models using Markov Chain Monte Carlo (MCMC) simulation not wholly unlike BUGS. JAGS was written with three aims in mind:

- To have a cross-platform engine for the BUGS language
- To be extensible, allowing users to write their own functions, distributions and samplers
- To be a platform for experimentation with ideas in Bayesian modelling

NIMBLE

Numerical Inference for statistical Models for Bayesian and Likelihood Estimation

<https://r-nimble.org/>

NIMBLE is a system for building and sharing analysis methods for statistical models, especially for hierarchical models and computationally-intensive methods. NIMBLE is built in R but compiles your models and algorithms using C++ for speed. It includes three components:

- A system for using models written in the BUGS model language as programmable objects in R.
- An initial library of algorithms for models written in BUGS, including basic MCMC, which can be used directly or can be customized from R before being compiled and run.
- A language embedded in R for programming algorithms for models, both of which are compiled through C++ code and loaded into R.

NIMBLE can also be used without BUGS models as a way to compile simple R-like code into C++, which is then compiled and loaded into R with an interface function or object.

<http://mc-stan.org/>

Stan is a state-of-the-art platform for statistical modeling and high-performance statistical computation.

Users specify log density functions in Stan's probabilistic programming language and get:

- full Bayesian statistical inference with MCMC sampling (NUTS, HMC)
- approximate Bayesian inference with variational inference (ADVI)
- penalized maximum likelihood estimation with optimization (L-BFGS)
- Stan's math library provides differentiable probability functions & linear algebra (C++ autodiff)
- Additional R packages provide expression-based linear modeling, posterior visualization, and leave-one-out cross-validation
- It does not allow inference for discrete parameters

¹Named after Stanislaw Ulam, one of the developers of Monte Carlo methods in the

<http://www.r-inla.org/>

Based on the Integrated Nested Laplace Approach (Rue et al., 2009)

- Makes use of an integrated nested Laplace approximation and its simplified version
- Provides accurate approximations to the posterior *marginals*
- The main benefit of these approximations is computational: where Markov chain Monte Carlo algorithms need hours or days to run, INLA provides more precise estimates in seconds or minutes

- Consider a model with parameters θ_1 that are assigned normal priors, with the remaining parameters being denoted θ_2 with $G = \dim(\theta_1)$ and $V = \dim(\theta_2)$
- For ease of explanation, assume $\theta_1 \sim N_G(\mathbf{0}, \Sigma)$ where Σ depends on elements in θ_2
- The posterior is proportional to

$$\begin{aligned}
 \pi(\theta_1, \theta_2 \mid \mathbf{y}) &\propto \pi(\theta_1 \mid \theta_2) \pi(\theta_2) \prod_{i=1}^n p(\mathbf{y}_i \mid \theta_1, \theta_2) \\
 &\propto \pi(\theta_2) \mid \Sigma(\theta_2) \mid^{-1/2} \exp \left\{ -\frac{1}{2} \theta_1^T \Sigma(\theta_2)^{-1} \theta_1 \right. \\
 &\quad \left. + \sum_{i=1}^n \log p(\mathbf{y}_i \mid \theta_1, \theta_2) \right\} \quad (3)
 \end{aligned}$$

- Of particular interest are the posterior univariate marginal distributions $\pi(\theta_{1g} \mid \mathbf{y})$, $g = 1, \dots, G$, and $\pi(\theta_{2v} \mid \mathbf{y})$, $v = 1, 2, \dots, V$

INLA

- The normal parameters θ_1 are dealt with by analytical approximations (as applied to the term in the exponent of (3), conditional on specific values of θ_2)
- Numerical integration techniques are applied to θ_2 , so that V should not be too large for accurate inference
- For elements of θ_1 we write

$$\pi(\theta_{1g} | \mathbf{y}) = \int \pi(\theta_1 | \theta_2, \mathbf{y}) \times \pi(\theta_2 | \mathbf{y}) d\theta_2$$

which may be evaluated via the approximation

$$\begin{aligned} \tilde{\pi}(\theta_{1g} | \mathbf{y}) &= \int \tilde{\pi}(\theta_{1g} | \theta_2, \mathbf{y}) \times \tilde{\pi}(\theta_2 | \mathbf{y}) d\theta_2 \\ &\approx \sum_{k=1}^K \tilde{\pi}(\theta_{1g} | \theta_2, \mathbf{y}) \times \tilde{\pi}(\theta_2 | \mathbf{y}) \Delta_k, \end{aligned}$$

for a set of weights Δ_k , $k = 1, 2, \dots, K$. Laplace or related analytical approximations are applied to carry out the integration (over $\theta_{1g'}$, $g' \neq g$) required for evaluation of $\tilde{\pi}(\theta_{1g} | \theta_2, \mathbf{y})$.

- To produce the grid of points $\{\theta_2^{(k)}, k = 1, 2, \dots, K\}$ which numerical integration is performed, the mode of $\tilde{\pi}(\theta_2 | \mathbf{y})$ is located and the Hessian is approximated, from which the grid of points $\{\theta_2^{(k)}, k = 1, 2, \dots, K\}$, with associated weights Δ_k , is created and used in the approximation above
- The output of INLA consists of posterior *marginal* distributions, which can be summarized via means, variances, and quantiles

Budworms example

- The larvae of the tobacco budworm causes much damages to crops in the US and Central and South America. (Bad for the farmers, good for public health?!)
- A study is conducted to investigate the dose of drug needed to kill the adult moths. Six different doses were applied to 20 male and 20 female moths. The number knocked down (uncoordinated) or dead 72 hours after exposure was recorded.
- We fit a Bayesian logistic regression on dose.

Budworms example

$$Y_i \sim \text{Binom}(n_i, p_i)$$
$$\log\left(\frac{p_i}{1-p_i}\right) = \beta_0 + \beta_1 \text{dose}_i + \beta_2 \text{male}_i$$

- Bayesian model is complete after assigning a prior distribution for $\beta = (\beta_0, \beta_1, \beta_2)$
- Regardless of the prior distribution you assign, the posterior does not have a closed analytical form
- if one assumes that components of β are independent, each normally distributed with some known mean c and variance C , we have that

$$p(\beta_0, \beta_1, \beta_2 \mid \mathbf{y}) \propto l(\mathbf{y} \mid \beta) p(\beta)$$
$$\propto \prod_{j=1}^k p_j^{y_j} (1-p_j)^{(n_j-y_j)} \prod_{i=0}^2 \exp\left\{-\frac{1}{2C}(\beta_i - c)^2\right\}$$