

# Machine Learning - Waseda University

## Markov Chain Monte Carlo Methods

AD

July 2011

- Bayesian model: likelihood  $f(x|\theta)$  and prior distribution  $\pi(\theta)$ .
- Bayesian inference is based on the posterior distribution

$$\pi(\theta|x) = \frac{\pi(\theta) f(x|\theta)}{\pi(x)}$$

where

$$\pi(x) = \int_{\Theta} \pi(\theta) f(x|\theta) d\theta.$$

- Many point estimates require computing additional integrals, e.g.

$$\mathbb{E}[\varphi(\theta)|x] = \int_{\Theta} \varphi(\theta) \pi(\theta|x) d\theta$$

- Assume  $\theta = (\theta_1, \theta_2) \in \Theta_1 \times \Theta_2$  then, if  $\theta_2$  are some nuisance parameters, we are only interested in the marginal posterior distribution

$$\pi(\theta_1|x) = \int_{\Theta_2} \pi(\theta_1, \theta_2|x) d\theta_2.$$

- Assume  $Y|\theta \sim g(y|\theta)$  then, given the observations  $x$ , the predictive distribution is

$$\pi(y|x) = \int_{\Theta} g(y|\theta) \pi(\theta|x) d\theta.$$

- Although Bayesian inference is conceptually simple, it requires being able to compute potentially high-dimensional integrals.
- In previous lectures, we have discussed mostly examples where these calculations could be performed analytically. However, analytic tractability seriously restricts the class of models we can work with.
- Approximations such as the Laplace's approximation and BIC are rather crude and require large sample sizes. Consequently until the beginning of the 90's, Bayesian inference for all but simple models could not be implemented.
- Alternative deterministic Bayesian approximation techniques include variational and expectation-propagation but make strong assumptions.

- MCMC are a class of powerful simulation-based algorithms that allows us to sample (approximately) from any high-dimensional probability distribution.
- The availability of these algorithms has truly revolutionized many fields, allowing people to fit complex models.
- MCMC are now widely used in bioinformatics, machine learning, econometrics, applied statistics, genetics, etc.

# Introduction to Monte Carlo

- Assume you are interested in approximating an high-dimensional pdf  $\pi(\theta|x)$ .
- A Monte Carlo approximation consists of sampling a large number  $N$  of i.i.d. random variables  $\theta^{(i)} \sim \pi(\theta|x)$  and build the following approximation

$$\hat{\pi}_N(\theta|x) = \frac{1}{N} \sum_{i=1}^N \delta_{\theta^{(i)}}(\theta)$$

where  $\delta_a(\theta)$  is the delta-Dirac mass which is such that

$$\int_A \delta_a(\theta) d\theta = \begin{cases} 1 & \text{if } a \in A, \\ 0 & \text{otherwise.} \end{cases}$$

- This approach is in contrast with what is usually done in parametric statistics, *i.e.* start with samples and then introduce a distribution with an algebraic representation for the underlying population.

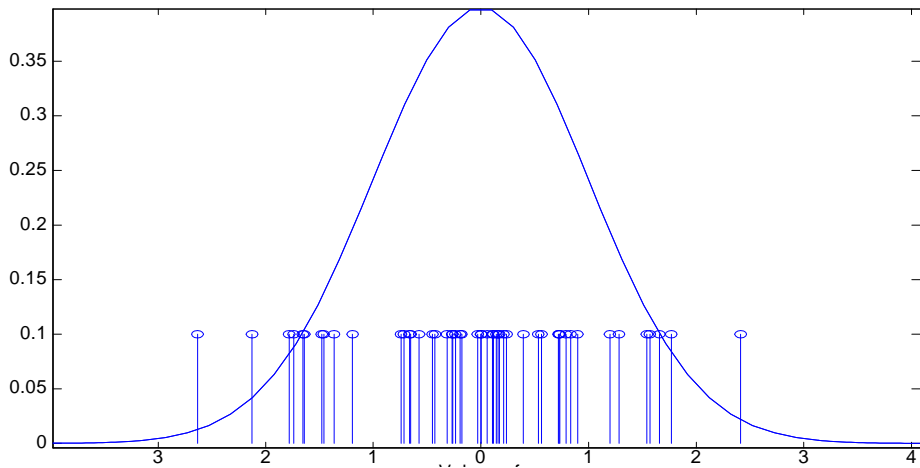


Figure: Normal  $\mathcal{N}(0, 1)$  and Monte Carlo approximation using  $N = 50$  samples

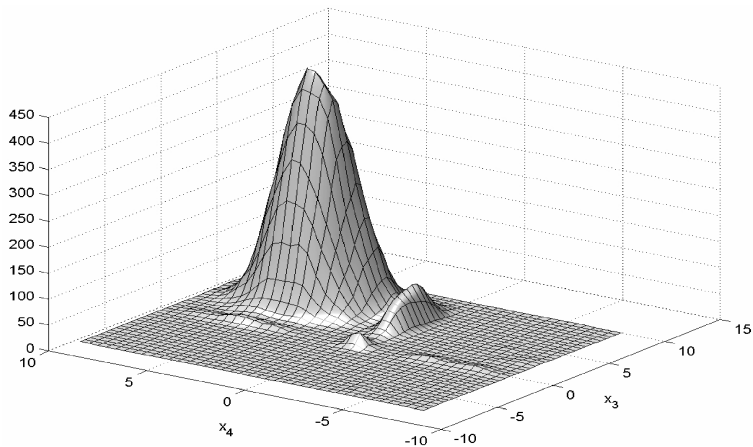


Figure: Bivariate non-standard probability density function



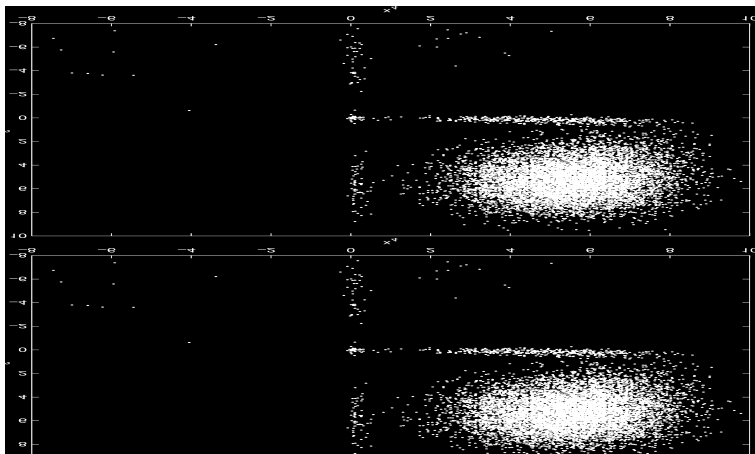


Figure: Scatterplot of  $N = 1000$  random samples

- You can think of the Monte Carlo method as an ‘clever’ discretization of the space where automatically the samples concentrate themselves in regions of high probability mass.
- An alternative deterministic approximation would consist of discretizing  $\Theta$  using a regular grid and then computing

$$\tilde{\pi}_N(\theta) = \frac{\sum_{i=1}^N \pi(\theta^{(i)} | x) \delta_{\theta^{(i)}}(\theta)}{\sum_{j=1}^N \pi(\theta^{(j)} | x)}.$$

- Such an approach will be extremely inefficient in high-dimensional spaces. If  $\Theta = \mathbb{R}^d$  and you discretize say each dimension using  $p$  values, then the total number of points in the grid will be  $p^d$ . Given we are routinely interested in problems where  $d = 100$ , even the crudest discretization  $p = 2$  would require  $2^{100} \gg 1$  points.

# Properties of Monte Carlo Estimates

- Now consider the problem of estimating

$$\mathbb{E}_{\pi}(\varphi) = \int_{\Theta} \varphi(\theta) \pi(\theta | x) d\theta$$

using Monte Carlo where  $\varphi : \Theta \rightarrow \mathbb{R}$ .

- We substitute to  $\pi(\theta | x)$  its Monte Carlo approximation  $\hat{\pi}_N$  and obtain

$$\begin{aligned} \mathbb{E}_{\hat{\pi}_N}(\varphi) &= \int_{\Theta} \varphi(\theta) \left( \sum_{i=1}^N \frac{1}{N} \delta_{\theta^{(i)}}(\theta) \right) d\theta = \sum_{i=1}^N \frac{1}{N} \int_{\Theta} \varphi(\theta) \delta_{\theta^{(i)}}(\theta) d\theta \\ &= \frac{1}{N} \sum_{i=1}^N \varphi(\theta^{(i)}). \end{aligned}$$

- This estimate is unbiased

$$\mathbb{E}_{\{\theta^{(i)}\}} [\mathbb{E}_{\hat{\pi}_N}(\varphi)] = \frac{1}{N} \sum_{i=1}^N \mathbb{E}_{\theta^{(i)}} \left( \varphi(\theta^{(i)}) \right) = \mathbb{E}_{\pi}(\varphi).$$

- The variance is given by

$$\mathbb{V}_{\{\theta^{(i)}\}} [\mathbb{E}_{\hat{\pi}_N}(\varphi)] = \frac{1}{N} \mathbb{V}_{\pi}(\varphi(\theta)).$$

- The CLT yields

$$\sqrt{N} (\mathbb{E}_{\hat{\pi}_N}(\varphi) - \mathbb{E}_{\pi_N}(\varphi)) \xrightarrow{d} \mathcal{N}(0, \mathbb{V}_{\pi}(\varphi(\theta))).$$

- The most remarkable property of the MC estimate is that the rate of convergence is independent of the dimension of  $\Theta$ .
- It is sometimes said that the MC beats the *curse of dimensionality*. This is not quite true.

- MC can easily be used to compute marginal distributions. We approximate

$$\begin{aligned}\pi(\theta_1|x) &= \int_{\Theta_2} \pi(\theta_1, \theta_2|x) d\theta_2 \approx \int_{\Theta_2} \hat{\pi}_N(\theta_1, \theta_2|x) d\theta_2 \\ &= \frac{1}{N} \int \sum_{i=1}^N \delta_{\theta_1^{(i)}, \theta_2^{(i)}}(\theta_1, \theta_2) d\theta_2 \\ &= \frac{1}{N} \sum_{i=1}^N \delta_{\theta_1^{(i)}}(\theta_1) = \hat{\pi}_N(\theta_1|x)\end{aligned}$$

- Similarly the predictive distribution can be approximated easily

$$\begin{aligned}\pi(y|x) &= \int_{\Theta} g(y|\theta) \pi(\theta|x) d\theta \approx \int_{\Theta} g(y|\theta) \hat{\pi}_N(\theta|x) d\theta \\ &= \frac{1}{N} \sum_{i=1}^N g(y|\theta^{(i)}).\end{aligned}$$

- However, the marginal likelihood  $\pi(x)$  cannot be estimated easily using samples from  $\pi(\theta|x)$ !

# Sampling from complex distributions

- There are standard methods to sample from classical distributions such as Beta, Gamma, Normal, Poisson etc. We will not detail them here.
- We are interested in problems where  $\pi(\theta|x)$  is not standard and is only known up to a normalizing constant; i.e.

$$\pi(\theta|x) = \frac{\pi(\theta) f(x|\theta)}{\pi(x)}$$

where  $\pi(\theta) f(x|\theta)$  is known pointwise whereas  $\pi(x)$  is unknown.

- There is no method available able to sample from any high dimensional probability distribution.

- The rejection method relies on a so-called proposal distribution  $q(\theta|x)$  which is selected such that it is easy to sample from it.
- We have  $q(\theta|x) \propto q^*(\theta|x)$  where  $q^*(\theta|x)$  is known pointwise.
- We need  $q^*(\theta|x)$  to 'dominate'  $\pi(\theta)f(x|\theta)$ ; i.e.

$$C = \sup_{\theta \in \Theta} \frac{\pi(\theta)f(x|\theta)}{q^*(\theta|x)} < +\infty$$

- This implies  $\pi(\theta|x) > 0 \Rightarrow q^*(\theta|x) > 0$  but also that the tails of  $q(\theta|x)$  must be thicker than the tails of  $\pi(\theta|x)$ .

Consider  $M \geq C$ . Then the accept/reject procedure proceeds as follows.

- 1 Sample  $\theta^* \sim q(\theta | x)$  and  $U \sim \mathcal{U}[0, 1]$ .
- 2 If  $U < \frac{\pi(\theta^*)f(x|\theta^*)}{Mq^*(\theta^*|x)}$  then return  $\theta^*$ ; otherwise return to step 1.



# Proof of Validity

- We have for any  $\theta \in \Theta (= \mathbb{R})$ , this is only to simplify notation)

$$\begin{aligned} & \Pr(\theta^* \leq \theta \text{ and } \theta^* \text{ accepted}) \\ &= \int_{\Theta} \int_0^1 \mathbb{I}(\theta^* \leq \theta) \mathbb{I}\left(u \leq \frac{\pi(\theta^*) f(x|\theta^*)}{M q^*(\theta^*|x)}\right) q(\theta^*|x) \times 1 du d\theta^* \\ &= \int_{-\infty}^{\theta} \frac{\pi(\theta^*) f(x|\theta^*)}{M q^*(\theta^*|x)} q(\theta^*|x) d\theta^* \\ &= \frac{\int_{-\infty}^{\theta} \pi(\theta^*) f(x|\theta^*) d\theta^*}{M \int_{\Theta} q^*(\theta^*|x) d\theta^*}. \end{aligned}$$

- The probability of being accepted is simply

$$\Pr(\theta^* \text{ accepted}) = \frac{\int_{\Theta} \pi(\theta^*) f(x|\theta^*) d\theta^*}{M \int_{\Theta} q^*(\theta^*|x) d\theta^*} = \frac{\pi(x)}{M \int_{\Theta} q^*(\theta^*|x) d\theta^*}.$$

- Finally, we have

$$\begin{aligned}\Pr(\theta^* \leq \theta \mid \theta^* \text{ accepted}) &= \frac{\Pr(\theta^* \leq \theta \text{ and } \theta^* \text{ accepted})}{\Pr(\theta^* \text{ accepted})} \\&= \frac{\int_{-\infty}^{\theta} \pi(\theta^*) f(x \mid \theta^*) d\theta^*}{\pi(x)} \\&= \int_{-\infty}^{\theta} \pi(\theta^* \mid x) d\theta^*.\end{aligned}$$

- The acceptance probability  $\Pr(\theta^* \text{ accepted})$  is a measure of efficiency.
- The number of trials before accepting a candidate follows a geometric distribution

$$\Pr(k^{\text{th}} \text{ proposal accepted}) = (1 - \rho)^{k-1} \rho$$

$$\text{where } \rho = \left( \frac{\pi(x)}{M \int_{\Theta} q^*(\theta^* | x) d\theta^*} \right)$$

thus its expected value is

$$\sum_{k=0}^{\infty} k (1 - \rho)^{k-1} \rho = \frac{1}{\rho} = \frac{1}{\Pr(\theta^* \text{ accepted})}.$$

# Toy Example

- Consider the following distribution

$$\pi(\theta|x) \propto \exp(-\theta^2/2) \left( \sin(\theta x)^2 + 3 \cos(\theta x)^2 \sin(4x)^2 + 1 \right)$$

- We use  $q^*(\theta|x) = q^*(\theta) = \exp(-\theta^2/2)$ , that is  $q(\theta) = \mathcal{N}(\theta; 0, 1)$ .
- This proposal is easy to sample and

$$\frac{\exp(-\theta^2/2) \left( \sin(\theta x)^2 + 3 \cos(\theta x)^2 \sin(4x)^2 + 1 \right)}{\exp(-\theta^2/2)} \leq 1 + 3 + 1 = 5$$

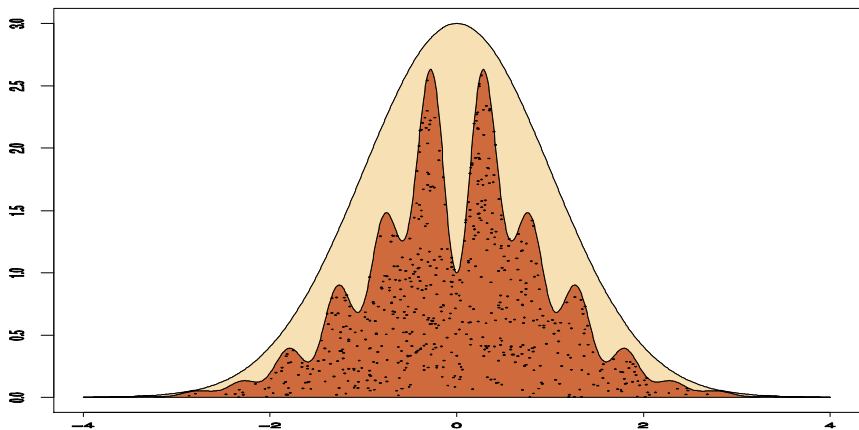


Figure:  $5q^*(\theta)$  and unnormalized target distribution  $\exp\left(-\theta^2/2\right) \left(\sin(\theta x)^2 + 3 \cos(\theta x)^2 \sin(4x)^2 + 1\right)$

# Using the Prior as Proposal

- A simple choice consists of selecting  $q(\theta|x) = q^*(\theta|x) = \pi(\theta)$ .
- This is possible if the likelihood is upper bounded as

$$\sup_{\theta \in \Theta} \frac{\pi(\theta) f(x|\theta)}{q^*(\theta|x)} = \sup_{\theta \in \Theta} f(x|\theta) \leq M$$

- In this case, expected value before acceptance is  $\rho^{-1}$  where

$$\rho = \frac{\int_{\Theta} \pi^*(\theta^*) f(x|\theta^*) d\theta^*}{M \int_{\Theta} q^*(\theta^*|x) d\theta^*} = \frac{\pi(x)}{M}$$

and provides us with an estimate of the marginal likelihood.

# Inefficiency of the Accept-Reject Strategy

- Consider  $\theta \sim \mathcal{N}(0, \tau^2)$  and  $X_i | \theta \sim \mathcal{N}(\theta, \sigma^2)$ .
- In this case, we have

$$\pi(\theta | x_{1:n}) = \mathcal{N}(\theta; m, v^2)$$

where  $v^2 = \left(\frac{1}{\tau^2} + \frac{n}{\sigma^2}\right)^{-1}$ ,  $m = v^2 \left(\frac{\sum_{i=1}^n x_i}{\sigma^2}\right)$  and

$$\pi(x_{1:n}) = \frac{v}{(2\pi\sigma^2)^{n/2} \tau} \exp\left(\frac{m^2}{2v^2} - \frac{\sum_{i=1}^n x_i^2}{2\sigma^2}\right).$$

whereas the likelihood is bounded by  $(2\pi\sigma^2)^{-n/2}$  so

$$\rho = \frac{\pi(x_{1:n})}{M} = \frac{v}{\tau} \exp\left(\frac{m^2}{2v^2} - \frac{\sum_{i=1}^n x_i^2}{2\sigma^2}\right).$$

- For  $\tau^2 \gg 1$  and  $\sigma^2 \ll 1$ , we have  $v^2 \approx \frac{\sigma^2}{n}$ ,  $m \approx \bar{x}$

$$\rho \approx \frac{\sigma^2}{n\tau} \exp\left(\frac{n(\bar{x}^2 - \overline{x^2})}{2\sigma^2}\right) \rightarrow 0$$

# Application to Genetic Linkage Model

- Consider the following genetic linkage model where observations

$$(X_1, X_2, X_3, X_4) \sim \mathcal{M} \left( n; \frac{1}{2} + \frac{\theta}{4}, \frac{1}{4} (1 - \theta), \frac{1}{4} (1 - \theta), \frac{\theta}{4} \right)$$

with  $\theta \in (0, 1)$ .

- We set an uniform prior  $\pi(\theta) = \mathcal{U}[0, 1](\theta)$  and want to estimate  $\pi(\theta | x_{1:4})$  where  $(x_1, x_2, x_3, x_4) = (125, 18, 20, 34)$ .
- We have

$$\pi(\theta | x_{1:4}) \propto (2 + \theta)^{x_1} (1 - \theta)^{x_2 + x_3} \theta^{x_4} \mathbf{1}_{(0,1)}(\theta).$$

- This univariate distribution is not standard and we sample from it using Accept-Reject.



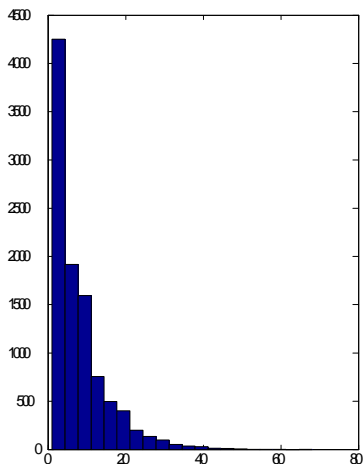
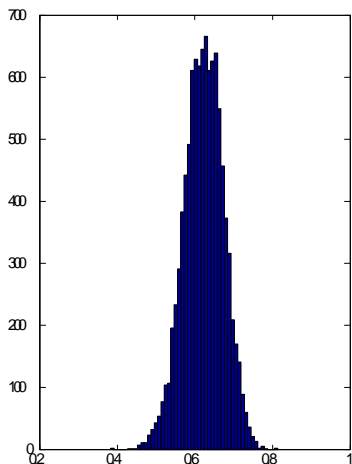
- We propose to use  $q(\theta|x) = q^*(\theta|x) = \mathcal{U}[0,1](\theta)$ .
- To apply accept-reject, we need to be able to upper bound over  $\theta \in (0,1)$  the function

$$g(\theta) = (2 + \theta)^{x_1} (1 - \theta)^{x_2 + x_3} \theta^{x_4}.$$

- Using a simple optimization algorithm (direct search or EM), we get

$$g(\theta) \leq g(\theta_{\max})$$

where  $\theta_{\max} = 0.6268$  and  $g(\theta_{\max}) = \exp(67.3841)$ .



**Figure:** Histogram approximation of  $\pi(\theta | x_{1:4})$  (left) and histogram approximation of waiting time distribution before acceptance (mean 7.8) (right)

# Limitations of Accept-Reject

- For complex problems, it will be impossible to bound the ratio between the (unnormalized) target and the proposal.
- Even if it is possible to obtain such a bound, the computational complexity typically increases exponentially fast with the dimension of the problem.
- We need to use more powerful techniques.

# The Gibbs Sampler

- The Gibbs sampler is an iterative popular method to sample from high dimensional probability distributions which has found numerous applications in Bayesian statistics.
- For sake of simplicity, consider first that we are interested in sampling from  $\pi(\theta)$  where  $\theta = (\theta_1, \theta_2) \in \Theta_1 \times \Theta_2$ .
- The Gibbs sampler relies on the fact that, although it is impossible to sample from  $\pi(\theta)$ , it is often possible to sample from the conditional distributions

$$\pi(\theta_1 | \theta_2) \text{ and } \pi(\theta_2 | \theta_1).$$

- Note that here  $\pi(\theta)$  denotes a generic pdf and could be the posterior!

# Two Component Gibbs Sampler aka Data Augmentation

- To sample from  $\pi(\theta_1, \theta_2)$ , the Gibbs sampler generates a Markov chain  $(\theta_1^{(i)}, \theta_2^{(i)})$  as follows.
- Initialization: Set  $(\theta_1^{(0)}, \theta_2^{(0)})$  deterministically or randomly.
- Iteration  $i$ ;  $i \geq 1$ .
  - Sample  $\theta_1^{(i)} \sim \pi(\theta_1 | \theta_2^{(i-1)})$
  - Sample  $\theta_2^{(i)} \sim \pi(\theta_2 | \theta_1^{(i)})$ .

- Under weak assumptions, after many such iterations (usually several hundred or thousand are required), the sampling distribution of  $(\theta_1^{(i)}, \theta_2^{(i)})$  will approximate closely  $\pi(\theta_1, \theta_2)$ .
- We should discard the first hundred/thousand simulated samples, say  $N_0$ , whose distribution might be 'far' from  $\pi(\theta_1, \theta_2)$ .
- The samples  $(\theta_1^{(i)}, \theta_2^{(i)})$  are not independent but it is still valid to consider the approximation

$$\hat{\pi}_N(\theta_1, \theta_2) = \frac{1}{N_0 - N + 1} \sum_{i=N_0}^N \delta_{(\theta_1^{(i)}, \theta_2^{(i)})}(\theta_1, \theta_2).$$

- Note that a deterministic version of this algorithm:  
 $\theta_1^{(i)} = \arg \max \pi(\theta_1 | \theta_2^{(i-1)})$  and  $\theta_2^{(i)} = \arg \max \pi(\theta_2 | \theta_1^{(i)})$  would not converge typically to the global maximum of  $\pi(\theta_1, \theta_2)$ .

# Toy Example

- Consider the distribution

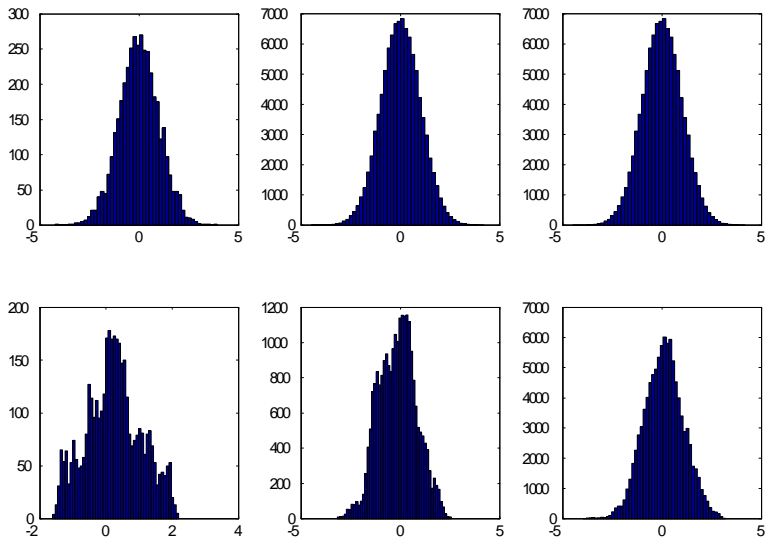
$$\begin{aligned}\pi(\theta_1, \theta_2) &= \mathcal{N}\left((\theta_1, \theta_2), \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix}\right) \\ &\propto \exp\left(-\frac{1}{2(1-\rho^2)}(\theta_1^2 - 2\rho\theta_1\theta_2 + \theta_2^2)\right)\end{aligned}$$

where  $|\rho| < 1$ .

- To sample from  $\pi(\theta_1, \theta_2)$ , the Gibbs sampler would sample iteratively and successively from

$$\begin{aligned}\pi(\theta_1 | \theta_2) &= \mathcal{N}(\theta_1; \rho\theta_2, 1 - \rho^2), \\ \pi(\theta_2 | \theta_1) &= \mathcal{N}(\theta_2; \rho\theta_1, 1 - \rho^2).\end{aligned}$$

- As  $|\rho| \rightarrow 1$  the algorithm will converge more slowly.
- Generally speaking, the more correlated  $\theta_1$  and  $\theta_2$  are the slower the algorithm converges.



**Figure:** Estimation of  $\pi(\theta_1) = \mathcal{N}(\theta_1, 0, 1)$  using the Gibbs sampler for  $N = 5000$  (left),  $25000$  (center) and  $100000$  (right) with  $N_0 = 1000$  and  $\rho = 0.2$  (top),  $\rho = 0.999$  (bottom)



# Application to Genetic Linkage Model

- In this example, we have

$$(X_1, X_2, X_3, X_4) \sim \mathcal{M} \left( n; \frac{1}{2} + \frac{\theta}{4}, \frac{1}{4} (1 - \theta), \frac{1}{4} (1 - \theta), \frac{\theta}{4} \right)$$

and the target posterior distribution is given by

$$\pi(\theta | x_{1:4}) \propto (2 + \theta)^{x_1} (1 - \theta)^{x_2 + x_3} \theta^{x_4} 1_{(0,1)}(\theta).$$

- We cannot use the Gibbs sampler in this case.
- Now assume we introduce the missing data  $(Z_1, Z_2)$  such that  $Z_1 + Z_2 = X_1$  and

$$(Z_1, Z_2, X_2, X_3, X_4) \sim \mathcal{M} \left( n; \frac{1}{2}, \frac{\theta}{4}, \frac{1}{4} (1 - \theta), \frac{1}{4} (1 - \theta), \frac{\theta}{4} \right).$$

- This complete likelihood is given by

$$g(z_{1:2}, x_{2:4} | \theta) \propto \theta^{z_2 + x_4} (1 - \theta)^{x_2 + x_3}.$$

- In the EM, we introduced these missing data to ease the maximization of the likelihood. Here we use them to ease the simulation from the posterior.
- We will now use the Gibbs sampler to generate sample from

$$\pi(\theta, z_1, z_2 | x_{1:4}).$$

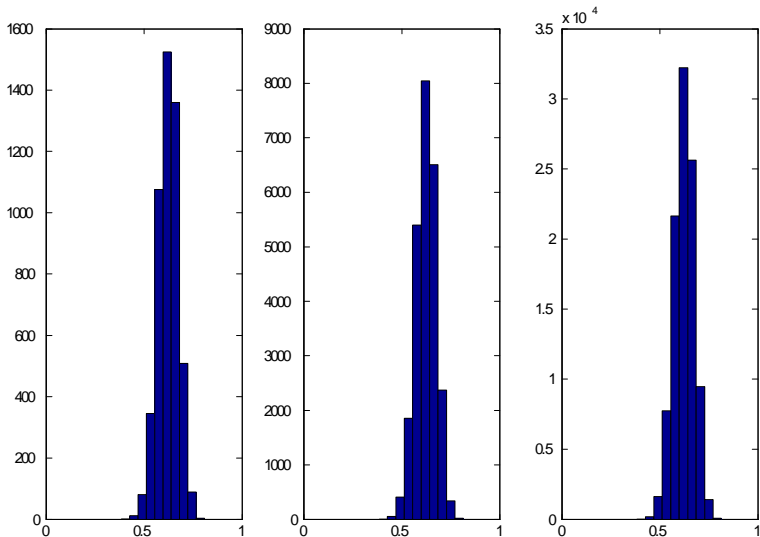
To implement the Gibbs sampler, we need to be able to sample from  $\pi(\theta | x_{1:4}, z_{1:2})$  and  $\pi(z_{1:2} | x_{1:4}, \theta)$ .

- We have

$$\begin{aligned} \pi(\theta | x_{1:4}, z_{1:2}) &\propto g(z_{1:2}, x_{2:4} | \theta) \pi(\theta) \propto \theta^{z_2 + x_4} (1 - \theta)^{x_2 + x_3} \\ &= \text{Beta}(\theta; z_2 + x_4 + 1, x_2 + x_3 + 1). \end{aligned}$$

- We have

$$\pi(z_{1:2} | x_{1:4}, \theta) = \mathcal{M}\left(z_{1:2}; x_1, \frac{\frac{1}{2}}{\frac{1}{2} + \frac{\theta}{4}}, \frac{\frac{\theta}{4}}{\frac{1}{2} + \frac{\theta}{4}}\right)$$



**Figure:** Estimation of  $\pi(\theta | x_{1:4})$  using the Gibbs sampler for  $N = 5000$  (left), 25000 (center) and 100000 (right) with  $N_0 = 1000$ .

- It is beyond the scope of this course to establish convergence results for the Gibbs sampler.
- **Proposition:** If  $\left(\theta_1^{(i-1)}, \theta_2^{(i-1)}\right) \sim \pi\left(\theta_1, \theta_2\right)$  then  $\left(\theta_1^{(i)}, \theta_2^{(i)}\right) \sim \pi\left(\theta_1, \theta_2\right)$ .
- *Proof.* The joint distribution of  $\left(\left(\theta_1^{(i-1)}, \theta_2^{(i-1)}\right), \left(\theta_1^{(i)}, \theta_2^{(i)}\right)\right)$  is given by

$$\pi\left(\theta_1^{(i-1)}, \theta_2^{(i-1)}\right) \pi\left(\theta_1^{(i)} \mid \theta_2^{(i-1)}\right) \pi\left(\theta_2^{(i)} \mid \theta_1^{(i)}\right)$$

so

$$\begin{aligned} & \int \int \pi\left(\theta_1^{(i-1)}, \theta_2^{(i-1)}\right) \pi\left(\theta_1^{(i)} \mid \theta_2^{(i-1)}\right) \pi\left(\theta_2^{(i)} \mid \theta_1^{(i)}\right) d\theta_1^{(i-1)} d\theta_2^{(i-1)} \\ &= \int \pi\left(\theta_2^{(i-1)}\right) \pi\left(\theta_1^{(i)} \mid \theta_2^{(i-1)}\right) \pi\left(\theta_2^{(i)} \mid \theta_1^{(i)}\right) d\theta_2^{(i-1)} \\ &= \int \pi\left(\theta_1^{(i)}, \theta_2^{(i-1)}\right) \pi\left(\theta_2^{(i)} \mid \theta_1^{(i)}\right) d\theta_2^{(i-1)} \\ &= \pi\left(\theta_1^{(i)}\right) \pi\left(\theta_2^{(i)} \mid \theta_1^{(i)}\right) = \pi\left(\theta_1^{(i)}, \theta_2^{(i)}\right). \end{aligned}$$

- We can extend straightforwardly the Gibbs sampler to the case where  $\theta = (\theta_1, \dots, \theta_p)$ . To sample from  $\pi(\theta)$ , the Gibbs sampler generates a Markov chain  $(\theta_1^{(i)}, \theta_2^{(i)}, \dots, \theta_p^{(i)})$  as follows.
- Initialization: Set  $(\theta_1^{(0)}, \theta_2^{(0)}, \dots, \theta_p^{(0)})$  deterministically or randomly.
- Iteration  $i$ ;  $i \geq 1$ .
  - For  $k = 1$  to  $p$   
Sample  $\theta_k^{(i)} \sim \pi(\theta_k | \theta_{-k}^{(i)})$  where  
 $\theta_{-k}^{(i)} = (\theta_1^{(i)}, \dots, \theta_{k-1}^{(i)}, \theta_{k+1}^{(i-1)}, \dots, \theta_p^{(i-1)})$ .

- Hierarchical models allow us to combine information and to combine the results of several studies addressing a set of related research hypotheses (meta-analysis).
- Consider a set of experiments/studies,  $j = 1, \dots, J$ , in which experiment  $j$  has data (vector)  $y_j$  and parameter (vector)  $\theta_j$ , with likelihood  $p(y_j | \theta_j)$ . [the method applies equally well to nonexperimental data].
- If no observation -other than the data  $\{y_j\}$ - is available to distinguish any of the  $\theta_j$ 's from any of the others, and no ordering/grouping of the parameters can be made, one must assume symmetry among the parameters in their prior; i.e. the parameters  $(\theta_1, \theta_2, \dots, \theta_J)$  should be exchangeable. This means that  $p(\theta_1, \theta_2, \dots, \theta_J)$  is invariant to permutation of the indexes  $(1, 2, \dots, J)$ .

- The simplest form of exchangeability is

$$p(\theta_1, \theta_2, \dots, \theta_J | \phi) = \prod_{j=1}^J p(\theta_j | \phi).$$

- In general  $\phi$  is unknown and we consider it random so that

$$p(\theta_1, \theta_2, \dots, \theta_J) = \int \left( \prod_{j=1}^J p(\theta_j | \phi) \right) p(\phi) d\phi.$$

- As a consequence, we have

$$p(\theta_j | y_1, y_2, \dots, y_J) \neq p(\theta_j | y_j).$$

- This model allows us to borrow information from different but related datasets.

## Example: Study of the effectiveness of cardiac treatments

- Assume we have  $J$  hospitals.
- Let  $y_j$  be data corresponding say to the number of patients having survived a cardiac treatment.
- Let  $\theta_j$  be the survival probability for patients in hospital  $j$ .
- For example, we could have a model like

$$Y_j | (n_j, \theta_j) \sim \text{Bin}(n_j, \theta_j)$$

and

$$\theta_j \sim \text{Beta}(\alpha, \beta).$$

- This model reflects the fact that the  $\theta_j$  are different but related to each other.
- Further, we put a prior distribution on the hyperparameters  $(\alpha, \beta)$ .



## Example: Baseball data

- We have the statistics of  $J = 17$  players in pre-season exhibition matches.
- The data  $y_j$  for the player  $j$  corresponds to the number of home runs in  $n_j$  times at the bat modelled through

$$Y_j | (n_j, p_j) \sim \text{Bin}(n_j, p_j).$$



$j$	McGw.	Sosa	Griffey	Castilla	Gonz.	Gala.	Palm.	Vaughn
$y_j$	7	9	4	7	3	6	2	10
$n_j$	58	59	74	84	69	63	60	54

$j$	Bond	Bag.	Piaz.	Thom.	Thom.	Mart.	Wal.	Burks	Buhner
$y_j$	2	2	4	3	2	5	3	2	6
$n_j$	53	60	66	66	72	64	42	38	58

- Following Efron & Morris (JASA, 1975), we define

$$X_j = f_{n_j} (Y_j / n_j)$$

where

$$f_a(u) = a^{1/2} \sin^{-1}(2u - 1).$$

- This is a variance stabilising transformation. It can be shown that we now have approximately

$$X_j | \theta_j \sim \mathcal{N}(\theta_j, 1)$$

where

$$\theta_j = f_{n_j}(p_j).$$

- *Remark:* We only use this transformation to compare our results to the James-Stein estimate discussed later during the lectures.

# Bayesian Model for Baseball data

- We set an exchangeable prior distribution of the form

$$\theta_j | (\mu, \tau^2) \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mu, \tau^2)$$

and

$$\pi(\mu, \tau^2) = \pi(\mu) \pi(\tau^2)$$

with

$$\pi(\tau^2) = \mathcal{IG}\left(\tau^2; \frac{a}{2}, \frac{b}{2}\right), \quad \pi(\mu) \propto 1.$$

where  $a = b = 0.001$ .

- Note that if you are a specialist of baseball, an exchangeable prior might not be appropriate.

- The full posterior distribution is given by

$$\begin{aligned}
 & \pi(\mu, \tau^2, \theta_{1:J} | x_{1:J}) \\
 \propto & \pi(\mu, \tau^2) \prod_{j=1}^J \pi(\theta_j | \mu, \tau^2) \prod_{j=1}^J f(x_j | \theta_j) \\
 \propto & \frac{1}{\tau^{J+a+2}} \exp \left( -\frac{b}{2\tau^2} - \sum_{j=1}^J \left( \frac{(\theta_j - \mu)^2}{2\tau^2} + \frac{(x_j - \theta_j)^2}{2} \right) \right)
 \end{aligned}$$

- This distribution does not admit a closed-form expression and we are going to use the Gibbs sampler by decomposing the parameter space in 3 blocks  $\mu$ ,  $\tau^2$  and  $\theta_{1:J}$ .
- The Gibbs sampler will require being able to sample from  $\pi(\mu | x_{1:J}, \tau^2, \theta_{1:J})$ ,  $\pi(\tau^2 | x_{1:J}, \mu, \theta_{1:J})$  and  $\pi(\theta_{1:J} | x_{1:J}, \mu, \tau^2)$ .

# Gibbs Sampling for Baseball Data

- Full conditional distribution for  $\mu$

$$\begin{aligned}\pi(\mu | x_{1:J}, \tau^2, \theta_{1:J}) &\propto \exp\left(-\sum_{j=1}^J \frac{(\theta_j - \mu)^2}{2\tau^2}\right) \\ &= \mathcal{N}\left(\mu; J^{-1} \sum_{j=1}^J \theta_j, J^{-1} \tau^2\right).\end{aligned}$$

- Full conditional distribution for  $\tau^2$

$$\begin{aligned}\pi(\tau^2 | x_{1:J}, \mu, \theta_{1:J}) &\propto \frac{1}{\tau^{J+a+2}} \exp\left(-\frac{b}{2\tau^2} - \sum_{j=1}^J \frac{(\theta_j - \mu)^2}{2\tau^2}\right) \\ &= \mathcal{IG}\left(\tau^2; \frac{J+a}{2}, \frac{b + \sum_{j=1}^J (\theta_j - \mu)^2}{2}\right).\end{aligned}$$

- Full conditional distribution for  $\pi(\theta_{1:J} | x_{1:J}, \mu, \tau^2)$

$$\pi(\theta_{1:J} | x_{1:J}, \mu, \tau^2) = \prod_{j=1}^J \pi(\theta_j | x_j, \mu, \tau^2)$$

where

$$\begin{aligned} \pi(\theta_j | x_j, \mu, \tau^2) &\propto \exp\left(-\left(\frac{(\theta_j - \mu)^2}{2\tau^2} + \frac{(x_j - \theta_j)^2}{2}\right)\right) \\ &= \mathcal{N}\left(\theta_j; \frac{\mu + \tau^2 x_j}{1 + \tau^2}, \frac{\tau^2}{1 + \tau^2}\right). \end{aligned}$$

- These three distributions can be sampled using standard procedures and the simulated samples  $(\mu^{(i)}, \tau^{2(i)}, \theta_{1:J}^{(i)})$  are (asymptotically) distributed from the posterior.
- We applied the Gibbs sampling using  $N = 10000$  iterations and compared it to a simpler empirical Bayesian analysis.

# Empirical Bayesian Analysis for Baseball Data

- In an empirical Bayesian analysis of the data, we just obtain a point estimate  $(\hat{\mu}, \hat{\tau}^2)$  of  $(\mu, \tau^2)$  using the data.
- For example, given that

$$\begin{aligned} f(x_j | \mu, \tau^2) &= \int f(x_j | \theta_j) \pi(\theta_j | \mu, \tau^2) d\theta_j \\ &= \mathcal{N}(x_j; \mu, 1 + \tau^2) \end{aligned}$$

when can select

$$\hat{\mu} = \frac{1}{J} \sum_{j=1}^J x_j, \quad \hat{\tau}^2 = \frac{1}{J} \sum_{j=1}^J (x_j - \hat{\mu})^2 - 1$$

- This is *not* a Bayesian calculation because it is not based on any specified full probability model. The selection of the point estimate is somewhat arbitrary (although principled) and no uncertainty about  $\mu, \tau^2$  is taken into account.

- Moreover, the distribution

$$\pi \left( \theta_j | x_j, \hat{\mu}, \hat{\tau}^2 \right) = \mathcal{N} \left( \theta_j; \frac{\hat{\mu} + \hat{\tau}^2 x_j}{1 + \hat{\tau}^2}, \frac{\hat{\tau}^2}{1 + \hat{\tau}^2} \right)$$

are necessarily Gaussian and we somehow use the data twice.

- The main advantage of the empirical Bayes approach is that it is very easy to implement compared to the 'full' Bayesian approach which relies no MCMC.



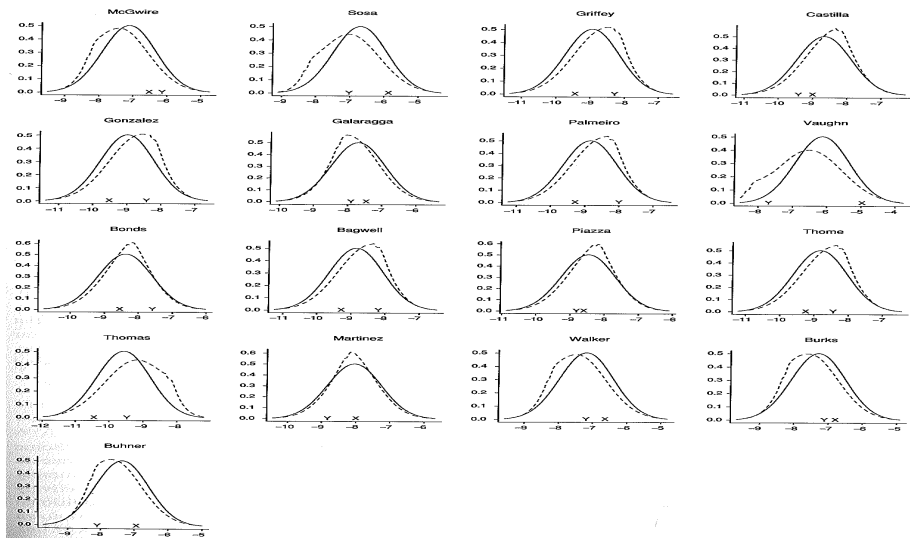


Figure: Posterior distributions  $\pi(\theta_j | x_{1:j})$  (dashed line) estimated using MCMC and  $\pi(\theta_j | x_j, \hat{\mu}, \hat{\tau}^2)$  (solid line) estimated using Empirical Bayes.

- The estimate of  $\pi(\theta_j | x_j)$  was not obtained by smoothing the histogram of the simulated values  $\theta_j^{(i)}$  but we use the fact that

$$\pi(\theta_j | x_{1:j}) = \int \pi(\theta_j | x_j, \mu, \tau^2) \pi(\mu, \tau^2 | x_{1:j}) d\mu d\tau^2$$

- So using the Monte Carlo approximation of  $\pi(\mu, \tau^2 | x_{1:j})$

$$\hat{\pi}(\mu, \tau^2 | x_{1:j}) = \frac{1}{N - N_0 + 1} \sum_{i=N_0}^N \delta_{\mu^{(i)}, \tau^{2(i)}}(\mu, \tau^2)$$

we obtain

$$\begin{aligned} \hat{\pi}(\theta_j | x_{1:j}) &= \frac{1}{N - N_0 + 1} \sum_{i=N_0}^N \pi(\theta_j | x_j, \mu^{(i)}, \tau^{2(i)}) \\ &= \frac{1}{N - N_0 + 1} \sum_{i=N_0}^N \mathcal{N}\left(\theta_j; \frac{\mu + \tau^{2(i)} x_j}{1 + \tau^{2(i)}}, \frac{\tau^{2(i)}}{1 + \tau^{2(i)}}\right). \end{aligned}$$

- This is a so-called Rao-Blackwellised estimate.

## Example: Nuclear Pump Data

- Multiple failures in a nuclear plant

Pump $j$	1	2	3	4	5
# Failures $p_j$	5	1	5	14	3
Times $t_j$	94.32	15.72	62.88	125.76	5.24
Pump $j$	6	7	8	9	10
# Failures $p_j$	19	1	1	4	22
Times $t_j$	31.44	1.05	1.05	2.10	10.48

- Model: Failures of the  $j$ -th pump follow a Poisson process with parameter  $\lambda_j$  ( $1 \leq j \leq 10$ ). For an observed time  $t_j$ , the number of failures  $p_j$  is thus a Poisson  $\mathcal{P}(\lambda_j t_j)$  random variable.
- The unknown parameters consist of  $\theta = (\lambda_1, \dots, \lambda_{10}, \beta)$ .

- Hierarchical model

$$\lambda_j | (\alpha, \beta) \stackrel{\text{iid}}{\sim} \mathcal{G}(\alpha, \beta) \text{ and } \beta \sim \mathcal{G}(\gamma, \delta)$$

with  $\alpha = 1.8$  and  $\gamma = 0.01$  and  $\delta = 1$ .

- The posterior distribution is proportional to

$$\begin{aligned} \pi(\lambda_{1:10}, \beta | p_{1:10}, t_{1:10}) &\propto \pi(\beta) \prod_{j=1}^{10} \pi(\lambda_j | \beta) \prod_{j=1}^{10} f(p_j | \lambda_j, t_j) \\ &\propto \beta^{\gamma-1} \exp(-\delta\beta) \prod_{j=1}^{10} \beta^\alpha \lambda_j^{\alpha-1} \exp(-\beta\lambda_j) \prod_{j=1}^{10} (\lambda_j t_j)^{p_j} \exp(-\lambda_j t_j) \end{aligned}$$

- This multidimensional distribution is rather complex. It is not obvious how the rejection method or importance sampling could be used in this context.

- The conditionals have a familiar form

$$\pi(\lambda_{1:10} | p_{1:10}, t_{1:10}, \beta) = \prod_{j=1}^{10} \pi(\lambda_j | p_j, t_j, \beta)$$

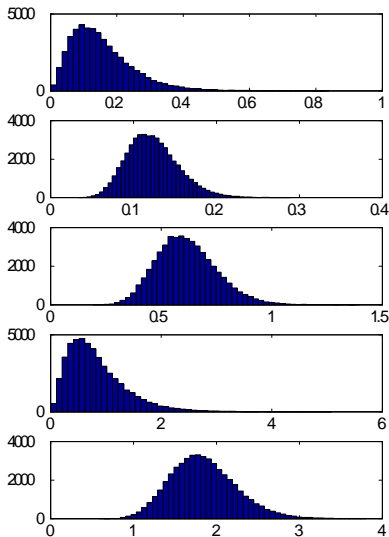
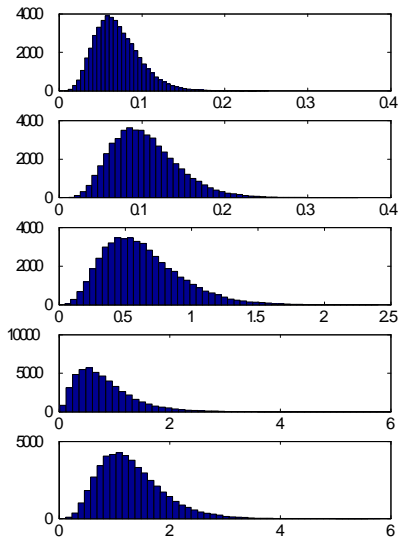
where

$$\begin{aligned} \pi(\lambda_j | p_j, t_j, \beta) &\propto \lambda_j^{p_j + \alpha - 1} \exp(-(t_j + \beta)\lambda_j) \\ &= \mathcal{G}(\lambda_j; p_j + \alpha, t_j + \beta) \end{aligned}$$

- For the hyperparameter

$$\begin{aligned} \pi(\beta | p_{1:10}, t_{1:10}, \lambda_{1:10}) &\propto \beta^{10\alpha + \gamma - 1} \exp\left(-\left(\delta + \sum_{j=1}^{10} \lambda_j\right)\beta\right) \\ &= \mathcal{G}(\beta; \gamma + 10\alpha, \delta + \sum_{j=1}^{10} \lambda_j). \end{aligned}$$

- Gibbs sampling is once more easily feasible in such cases.



**Figure:** Histogram approximations of  $\pi(\lambda_j | t_{1:10}, p_{1:10})$  obtained using  $N = 50,000$  and  $N_0 = 1000$ .

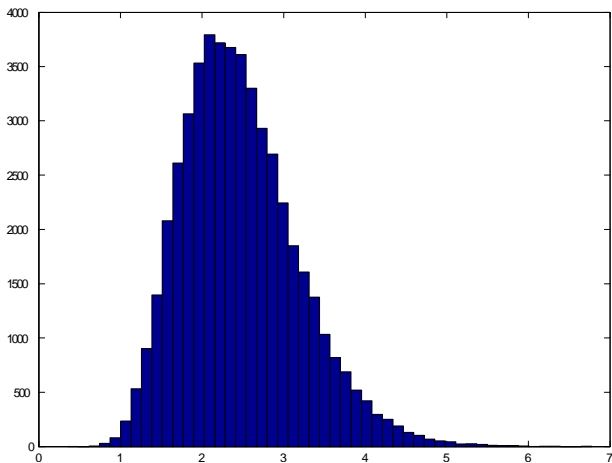


Figure: Histogram approximation of  $\pi(\beta | p_{1:10}, t_{1:10})$ .

- The posterior distribution of the hyperparameter is quite diffuse. Hence the results we obtained using a full Bayesian approach are significantly different from an empirical Bayes approach.

# Limitation of Gibbs Sampling

- The Gibbs sampler requires sampling from the full conditional distributions

$$\pi(\theta_k | \theta_{-k}).$$

- For many complex models, it is impossible to sample from several of these “full” conditional distributions.
- Even if it is possible to implement the Gibbs sampler, the algorithm might be very inefficient because the variables are very correlated or sampling from the full conditionals is extremely expensive/inefficient.



# Metropolis-Hastings Algorithm

- The Metropolis-Hastings algorithm is an alternative algorithm to sample from any probability distribution  $\pi(\theta)$  known up to a normalizing constant.
- This can be interpreted as the basis of all MCMC algorithms: It provides a generic way to build a Markov kernel admitting  $\pi(\theta)$  as an invariant distribution.
- The Metropolis algorithm was named the “Top algorithm of the 20th century” by computer scientists, mathematicians, physicists.

- Introduce a proposal distribution/kernel  $q(\theta' | \theta)$ , i.e.

$$\int q(\theta' | \theta) d\theta' = 1 \text{ for any } \theta.$$

- The basic idea of the MH algorithm is to propose a new candidate  $\theta'$  based on the current state of the Markov chain  $\theta$ .
- We only accept this algorithm with respect to a probability  $\alpha(\theta, \theta')$  which ensures that the invariant distribution of the transition kernel is the target distribution  $\pi(\theta)$ .

- Initialization: Select deterministically or randomly  $\theta^{(0)}$ .
- Iteration  $i$ ;  $i \geq 1$ :
  - Sample  $\theta^* \sim q\left(\theta \mid \theta^{(i-1)}\right)$  and compute

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

- With probability  $\alpha\left(\theta^{(i-1)}, \theta^*\right)$ , set  $\theta^{(i)}=\theta^*$ ; otherwise set  $\theta^{(i)}=\theta^{(i-1)}$ .
- Simulated annealing is an extremely popular optimization algorithm: it is a simple nonhomogeneous version of MH where at iteration  $i$  the target distribution is  $\pi_i(\theta) \propto[\pi(\theta)]^{\gamma_i}$  where  $\gamma_i$  is an increasing sequence going to  $\infty$ .

- It is not necessary to know the normalizing constant of  $\pi(\theta)$  to implement the algorithm.
- This algorithm is extremely general:  $q(\theta' | \theta)$  can be any proposal distribution. So in practice, we can select it so that it is easy to sample from it.
- There is much more freedom than in the Gibbs sampler where the proposal distributions are fixed once we have partitioned the vector parameter.

# Random Walk Metropolis

- The original Metropolis algorithm corresponds to the following choice for  $q(\theta' | \theta)$

$$\theta' = \theta + Z \text{ where } Z \sim g;$$

i.e. this is a so-called *random walk proposal*.

- The distribution  $g(z)$  is the distribution of the random walk increments  $Z$  and

$$q(\theta' | \theta) = g(\theta' - \theta) \Rightarrow \alpha(\theta, \theta') = \min \left( 1, \frac{\pi(\theta') g(\theta - \theta')}{\pi(\theta) g(\theta' - \theta)} \right).$$

- If  $g(z) = g(-z)$  - e.g.  $Z \sim \mathcal{N}(0, \Sigma)$  - then

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

- There is no clear guideline how to select the proposal distribution.
- When the variance of the random walk increments (if it exists) is very small then the acceptance rate can be expected to be around 0.5-0.7.
- You would like to scale the random walk moves such that it is possible to move reasonably fast in regions of positive probability masses under  $\pi$ .

- In the Bayesian framework where  $\pi(\theta) = p(\theta|x) \propto f(x|\theta)p(\theta)$  and  $q(\theta'|\theta) = q(\theta|\theta') = g(\theta' - \theta)$  then

$$\alpha(\theta, \theta') = \min \left( 1, \frac{f(x|\theta')p(\theta')}{f(x|\theta)p(\theta)} \right).$$

- Assuming  $g(z) = \mathcal{N}(z; 0, \Sigma)$  then the selection of  $\Sigma$  will be difficult.
- When the Fisher/observed information matrix at  $\theta$  is available, then we typically select  $\Sigma$  as the inverse of it.

# Toy Example

- Consider the case where

$$\pi(\theta) \propto \exp\left(-\frac{\theta^2}{2}\right).$$

- We implement the MH algorithm for

$$q_1(\theta' | \theta) \propto \exp\left(-\frac{(\theta' - \theta)^2}{2(0.2)^2}\right),$$

$$q_2(\theta' | \theta) \propto \exp\left(-\frac{(\theta' - \theta)^2}{2(5)^2}\right),$$

$$q_3(\theta' | \theta) \propto \exp\left(-\frac{(\theta' - \theta)^2}{2(0.02)^2}\right).$$



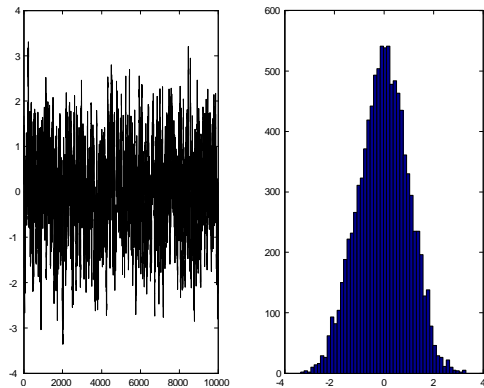


Figure: MCMC output for  $q_1$ , we estimate  $\mathbb{E}(\theta) = -0.02$  and  $\mathbb{V}(\theta) = 0.99$

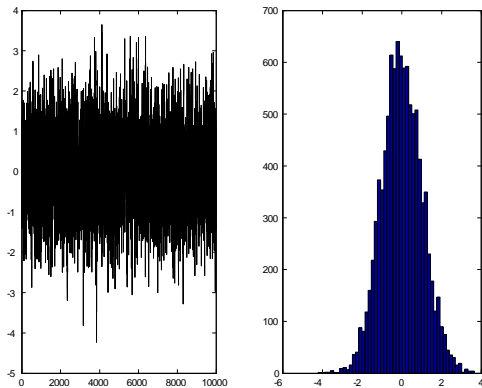


Figure: MCMC output for  $q_2$ , we estimate  $\mathbb{E}(\theta) = 0.00$  and  $\mathbb{V}(\theta) = 1.02$ .

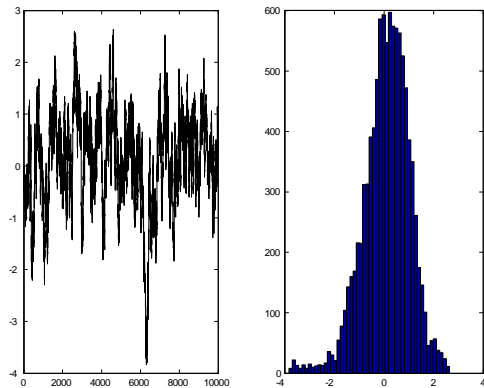


Figure: MCMC output for  $q_3$ , we estimate  $\mathbb{E}(\theta) = 0.10$  and  $\mathbb{V}(\theta) = 0.92$ .

# Independent Metropolis-Hastings

- The Hastings' generalization corresponds to the following choice for  $q(\theta' | \theta)$

$$q(\theta' | \theta) = q(\theta') ;$$

i.e. this is a so-called *independent proposal*.

- In this case, the acceptance probability is given by

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

- The ratio  $\pi(\theta) / q(\theta)$  also appears the Accept/Reject method.
- The optimal independent proposal is clearly  $q(\theta) = \pi(\theta)$ !

- In the Bayesian framework where  $\pi(\theta) = p(\theta|x) \propto f(x|\theta)p(\theta)$  and  $q(\theta) = p(\theta)$  then

$$\alpha(\theta, \theta') = \min \left( 1, \frac{f(x|\theta')}{f(x|\theta)} \right).$$

- Like the accept/reject method, this approach will be inefficient if the prior and the posterior are very different from each other.
- The MH is very flexible and we could use for example

$$q(\theta) = \mathcal{N}(\theta; \hat{\theta}_{MLE}, \sigma^2)$$

or a distribution with thicker tails.

- Any heuristic can be made rigorous using the MH algorithm.

# Toy example

- **Example:** Consider the case where

$$\pi(\theta) \propto \exp\left(-\frac{\theta^2}{2}\right).$$

- We implement the MH algorithm for

$$q_1(\theta) \propto \exp\left(-\frac{\theta^2}{2(0.2)^2}\right)$$

so  $\pi(\theta) / q_1(\theta) \rightarrow \infty$  as  $\theta \rightarrow \infty$  and for

$$q_2(\theta) \propto \exp\left(-\frac{\theta^2}{2(5)^2}\right)$$

so  $\pi(\theta) / q_2(\theta) \leq C < \infty$  for all  $\theta$ .

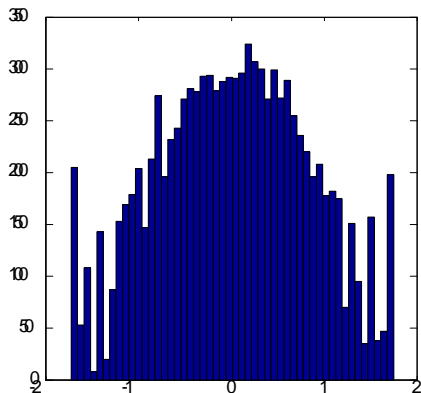
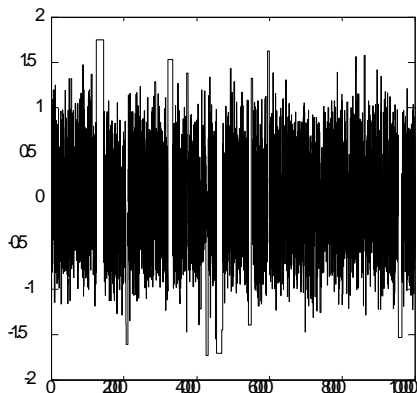


Figure: MCMC output for  $q_1$ , we estimate  $\mathbb{E}(\theta) = 0.0206$  and  $\mathbb{V}(\theta) = 0.83$ .

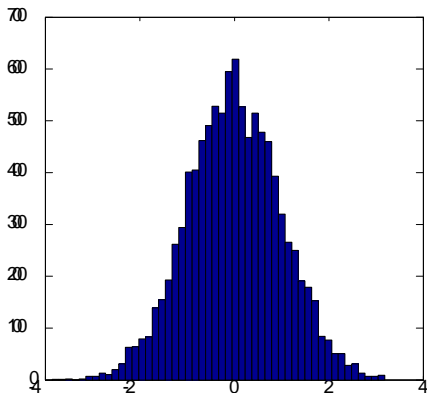
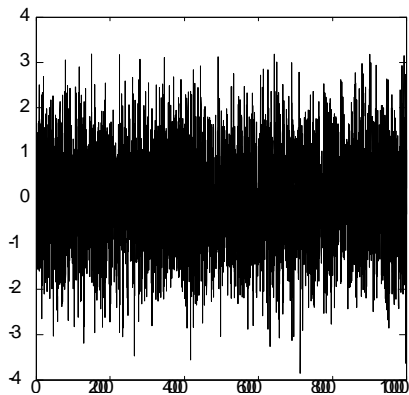


Figure: MCMC output for  $q_2$ , we estimate  $\mathbb{E}(\theta) = -0.004$  and  $\mathbb{V}(\theta) = 1.00$ .



- To establish that the MH chain converges towards the required target, we need to establish that
  - If  $\theta^{(i-1)} \sim \pi(\theta)$  then  $\theta^{(i)} \sim \pi(\theta)$ , i.e.  $\pi$  is an invariant distribution of the Markov kernel associated to the MH algorithm.
  - The Markov chain is irreducible; i.e. we can reach any set  $A$  such that  $\pi(A) > 0$ .
  - The Markov chain is aperiodic; i.e. we do not visit the state-space in a periodic way.

# Metropolis Hastings kernel

- The transition kernel associated to the MH algorithm can be rewritten as

$$K(\theta' | \theta) = \alpha(\theta, \theta') q(\theta' | \theta) + \underbrace{\left(1 - \int \alpha(\theta, u) q(u | \theta) du\right)}_{\text{rejection probability}} \delta_{\theta}(\theta')$$

- Clearly we have

$$\begin{aligned} \int K(\theta' | \theta) d\theta' &= \int \alpha(\theta, \theta') q(\theta' | \theta) d\theta' \\ &\quad + \left(1 - \int \alpha(\theta, u) q(u | \theta) du\right) \int \delta_{\theta}(\theta') d\theta' \\ &= 1. \end{aligned}$$

- **Proposition:** If  $\theta^{(i-1)} \sim \pi(\theta)$  then  $\theta^{(i)} \sim \pi(\theta)$ .
- To prove the result, we are going to establish the stronger *reversibility property*: for any  $\theta, \theta'$

$$\pi(\theta) K(\theta' | \theta) = \pi(\theta') K(\theta | \theta') ;$$

i.e. the probability of being in  $A$  and moving to  $B$  is equal to the probability of being in  $B$  and moving to  $A$ .

- Indeed the reversibility condition implies that

$$\begin{aligned} \int \pi(\theta) K(\theta' | \theta) d\theta &= \int \pi(\theta') K(\theta | \theta') d\theta \\ &= \pi(\theta') \int K(\theta | \theta') d\theta \\ &= \pi(\theta') \end{aligned}$$

- By definition of the kernel, we have

$$\begin{aligned} \pi(\theta) K(\theta' | \theta) &= \pi(\theta) \alpha(\theta, \theta') q(\theta' | \theta) \\ &\quad + \pi(\theta) \left( 1 - \int \alpha(\theta, u) q(u | \theta) du \right) \delta_{\theta}(\theta'). \end{aligned}$$

- Then

$$\begin{aligned} \pi(\theta) \alpha(\theta, \theta') q(\theta' | \theta) &= \pi(\theta) q(\theta' | \theta) \min \left( 1, \frac{\pi(\theta') q(\theta | \theta')}{\pi(\theta) q(\theta' | \theta)} \right) \\ &= \min(\pi(\theta) q(\theta' | \theta), \pi(\theta') q(\theta | \theta')) \\ &= \pi(\theta') q(\theta | \theta') \min \left( 1, \frac{\pi(\theta) q(\theta' | \theta)}{\pi(\theta') q(\theta | \theta')} \right) \\ &= \pi(\theta') \alpha(\theta', \theta) q(\theta | \theta'). \end{aligned}$$

- We have obviously

$$\begin{aligned} & \left( 1 - \int \alpha(\theta, u) q(u|\theta) du \right) \delta_{\theta}(\theta') \pi(\theta) \\ &= \left( 1 - \int \alpha(\theta', u) q(u|\theta') du \right) \delta_{\theta'}(\theta) \pi(\theta'). \end{aligned}$$

- It follows that

$$\pi(\theta) K(\theta'|\theta) = \pi(\theta') K(\theta|\theta').$$

- Hence,  $\pi$  is the invariant distribution of the transition kernel  $K$ .

- To ensure irreducibility, a sufficient but not necessary condition is that

$$\pi(\theta') > 0 \Rightarrow q(\theta' | \theta) > 0 \text{ for any } \theta$$

- Aperiodicity is automatically ensured as there is always a strictly positive probability to reject the candidate.
- Theoretically, the MH algorithm converges under very weak assumptions to the target distribution  $\pi$ . In practice, this convergence can be so slow that the algorithm is useless.

# Application to Genetic Linkage Model

- In this example, we have

$$(X_1, X_2, X_3, X_4) \sim \mathcal{M} \left( n; \frac{1}{2} + \frac{\theta}{4}, \frac{1}{4} (1 - \theta), \frac{1}{4} (1 - \theta), \frac{\theta}{4} \right)$$

and the target posterior distribution is given by

$$\pi(\theta | x_{1:4}) \propto (2 + \theta)^{x_1} (1 - \theta)^{x_2 + x_3} \theta^{x_4} 1_{(0,1)}(\theta).$$

- Accept/Reject requires maximizing the likelihood whereas Gibbs sampling requires introducing missing data.
- Alternatively, we can use a simple MH algorithm with proposal distribution  $q(\theta' | \theta) = \mathcal{N}(\theta'; \theta, \sigma^2)$ . See computer simulations.

# Application to Probit Regression

- We consider the following example: we take 4 measurements from 100 genuine Swiss banknotes and 100 counterfeit ones (Marin & Robert, *Bayesian Core*, Springer-Verlag, 2007).
- The response variable  $y$  is 0 for *genuine* and 1 for *counterfeit* and the explanatory variables are
  - $x^1$  the length,
  - $x^2$ : the width of the left edge
  - $x^3$ : the width of the right edge
  - $x^4$ : the bottom margin width



- As  $Y \in \{0, 1\}$ , we cannot have a linear model

$$\begin{aligned} Y &= x^1 \beta_1 + \dots + x^4 \beta_4 + \varepsilon \\ &= \mathbf{x}^\top \boldsymbol{\beta} + \varepsilon \end{aligned}$$

where  $\mathbf{x} = (x^1 \ x^2 \ x^3 \ x^4)^\top$  and  $\boldsymbol{\beta} = (\beta_1 \ \beta_2 \ \beta_3 \ \beta_4)^\top$ .

- We select a probit link here

$$\Pr(Y = 1 | \boldsymbol{\beta}, \mathbf{x}) = 1 - \Pr(Y = 0 | \boldsymbol{\beta}, \mathbf{x}) = \Phi(\mathbf{x}^\top \boldsymbol{\beta})$$

where

$$\Phi(u) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^u \exp\left(-\frac{v^2}{2}\right) dv$$

- For  $n = 200$  data, the likelihood is then given by

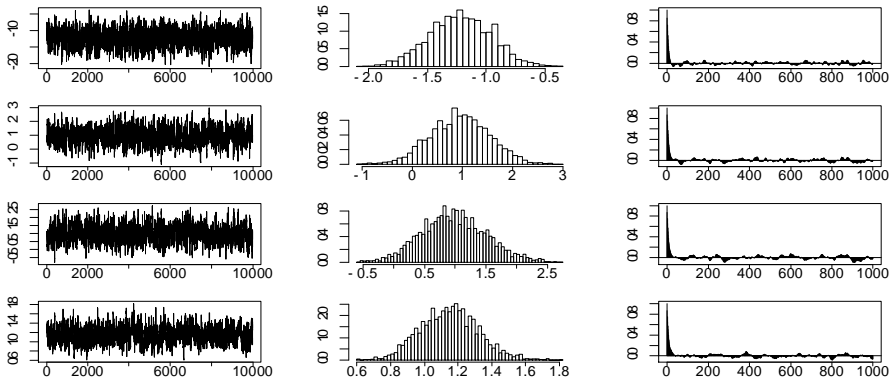
$$f(y_{1:n} | \boldsymbol{\beta}, \mathbf{x}_{1:n}) = \prod_{i=1}^n \Phi(\mathbf{x}_i^\top \boldsymbol{\beta})^{y_i} (1 - \Phi(\mathbf{x}_i^\top \boldsymbol{\beta}))^{1-y_i}.$$

- We assume a vague prior  $p(\beta) = \mathcal{N}(\beta; 0, 100I_4)$ .
- We use a simple random walk sampler where  $q(\beta' | \beta) = \mathcal{N}(\beta'; \beta, \tau^2 \Sigma)$  with  $\Sigma = \Omega^{-1}$ ,  

$$\Omega_{ij} = \left[ -\frac{\partial^2 \log f(y_{1:n} | \beta, \mathbf{x}_{1:n})}{\partial \beta_i \partial \beta_j} \right] \Big|_{\beta_{\text{MLE}}}.$$
- The algorithm is thus simply given at iteration  $i$  by
  - Sample  $\beta^* \sim \mathcal{N}(\beta^{(i-1)}, \tau^2 \Sigma)$  and compute

$$\begin{aligned} \alpha(\beta^{(i-1)}, \beta^*) &= \min \left( 1, \frac{p(\beta^* | y_{1:n}, \mathbf{x}_{1:n})}{p(\beta^{(i-1)} | y_{1:n}, \mathbf{x}_{1:n})} \right) \\ &= \min \left( 1, \frac{f(y_{1:n} | \beta^*, \mathbf{x}_{1:n}) p(\beta^*)}{f(y_{1:n} | \beta^{(i-1)}, \mathbf{x}_{1:n}) p(\beta^{(i-1)})} \right). \end{aligned}$$

- Set  $\beta^{(i)} = \beta^*$  with probability  $\alpha(\beta^{(i-1)}, \beta^*)$  and  $\beta^{(i)} = \beta^{(i-1)}$  otherwise.
- Algorithm tested with  $\tau^2 = 1, 10$  and  $0.1$ . Best results obtained with  $\tau^2 = 1$ .



**Figure:** Traces (left), Histograms (middle) and Autocorrelations (right) for  $(\beta_1^{(i)}, \dots, \beta_4^{(i)})$  for  $N = 10000$  and  $N_0 = 1000$ .

- We found for  $\hat{\beta} = \mathbb{E}(\beta | y_{1:n}, \mathbf{x}_{1:n}) = (-1.22, 0.95, 0.96, 1.15)$  so a simple plug-in estimate of the predictive probability of a counterfeit bill is

$$\Pr(Y = 1 | \mathbf{x}, \hat{\beta}) = \Phi(\mathbf{x}^\top \hat{\beta})$$

- The predictive distribution is given by

$$\Pr(Y = 1 | \mathbf{x}, \mathbf{x}_{1:n}) = \int \Phi(\mathbf{x}^\top \beta) \pi(\beta | y_{1:n}, \mathbf{x}_{1:n}) d\beta.$$

- We rerun the algorithm on only  $n = 100$  randomly selected data (50 genuine and 50 counterfeit) and use the results to classify the remaining 100 banknotes. The missclassification rate was 13% for the plug-in classifier and 7% for the full Bayesian approach.

# Limitations of the MH algorithm

- The MH algorithm is a simple and very general algorithm to sample from a target distribution  $\pi(\theta)$ .
- In practice, the choice of the proposal distribution has a crucial impact on the performance of the algorithm.
- In high dimensional problems, a simple MH algorithm will be useless. It will be necessary to use a combination of MH kernels.

# Towards more flexible algorithms

- It is possible to combine several MH kernels, say using

$$K(\theta' | \theta) = \lambda K_1(\theta' | \theta) + (1 - \lambda) K_2(\theta' | \theta),$$

$$K(\theta' | \theta) = \int K_1(\theta' | \theta'') K_2(\theta'' | \theta) d\theta''$$

where  $K_1$  (resp.  $K_2$ ) is an MH algorithm of proposal  $q_1$  (resp.  $q_2$ ).

- Each proposal can modify only a subset of the components of  $\theta$ . That is, if  $\theta = (\theta_1, \theta_2)$ , then we can have  $q_1(\theta' | \theta) = q_1(\theta'_1 | \theta) \delta_{\theta_2}(\theta'_2)$  and  $q_2(\theta' | \theta) = \delta_{\theta_1}(\theta'_1) q_2(\theta'_2 | \theta)$ .
- This algorithm satisfies

$$\begin{aligned} & \int \pi(\theta) K(\theta' | \theta) d\theta \\ &= \lambda \int \pi(\theta) K_1(\theta' | \theta) d\theta + (1 - \lambda) \int \pi(\theta) K_2(\theta' | \theta) d\theta \\ &= \lambda \pi(\theta') + (1 - \lambda) \pi(\theta') = \pi(\theta'). \end{aligned}$$

- To sample from  $\pi(\theta)$  where  $\theta = (\theta_1, \dots, \theta_p)$ , we can use the following algorithm at iteration  $i$ .
- Iteration  $i$ ;  $i \geq 1$ :
  - For  $k = 1 : p$ 
    - Sample  $\theta_k^{(i)}$  using an MH step of proposal distribution  $q_k(\theta_k | (\theta_{-k}^{(i)}, \theta_k^{(i-1)}))$  and target  $\pi(\theta_k | \theta_{-k}^{(i)})$  where  $\theta_{-k}^{(i)} = (\theta_1^{(i)}, \dots, \theta_{k-1}^{(i)}, \theta_{k+1}^{(i-1)}, \dots, \theta_p^{(i-1)})$ .
- *Remark:* It is possible to rewrite each MH step as a proposal  $q_k(\theta_k | (\theta_{-k}^{(i)}, \theta_k^{(i-1)})) \delta_{\theta_{-k}^{(i)}}(\theta_{-k})$  and target  $\pi(\theta_k, \theta_{-k})$ .

- The acceptance ratio of the MH step updating  $\theta_k$  is given by

$$\begin{aligned} & \alpha \left( \theta_{-k}^{(i)}, \theta_k^{(i-1)}, \theta_k^* \right) \\ &= \min \left( 1, \frac{\pi \left( \theta_k^* | \theta_{-k}^{(i)} \right) q_k \left( \theta_k^{(i-1)} | \left( \theta_{-k}^{(i)}, \theta_k^* \right) \right)}{\pi \left( \theta_k^{(i-1)} | \theta_{-k}^{(i)} \right) q_k \left( \theta_k^* | \left( \theta_{-k}^{(i)}, \theta_k^{(i-1)} \right) \right)} \right) \\ &= \min \left( 1, \frac{\pi \left( \theta_k^*, \theta_{-k}^{(i)} \right) q_k \left( \theta_k^{(i-1)} | \left( \theta_{-k}^{(i)}, \theta_k^* \right) \right)}{\pi \left( \theta_k^{(i-1)}, \theta_{-k}^{(i)} \right) q_k \left( \theta_k^* | \left( \theta_{-k}^{(i)}, \theta_k^{(i-1)} \right) \right)} \right) \end{aligned}$$

- If we select  $q_k \left( \theta_k | \left( \theta_{-k}^{(i)}, \theta_k^{(i-1)} \right) \right) = \pi \left( \theta_k | \theta_{-k}^{(i)} \right)$  then  $\alpha \left( \theta_{-k}^{(i)}, \theta_k^{(i-1)}, \theta_k^* \right) = 1$  and we are back to the Gibbs sampler.
- Example:** Assume we have  $\pi(\theta_1, \theta_2)$  where it is easy to sample from  $\pi(\theta_1 | \theta_2)$  and  $\pi(\theta_2 | \theta_1)$  is not standard. We then use a Gibbs step to update  $\theta_1$  and an MH step to update  $\theta_2$ .



## Example: Baseball data revisited

- We have the statistics of  $J = 17$  players in pre-season exhibition matches where the data  $y_j$  for the player  $j$  corresponds to the number of home runs in  $n_j$  times at the bat modelled through

$$Y_j | (n_j, p_j) \sim \text{Bin}(n_j, p_j).$$

- We use the logit transformation and write

$$\theta_j = \log \left( \frac{p_j}{1 - p_j} \right) \Leftrightarrow p_j = \frac{\exp(\theta_j)}{1 + \exp(\theta_j)}$$

which translates the parameter range  $(0, 1)$  for  $p_j$  to  $(-\infty, \infty)$  for  $\theta_j$ .

- We set an exchangeable prior where

$$\theta_j | (\mu, \tau) \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mu, \tau^2).$$

and the hyperprior is selected as before with  $a = b = 0.001$ .

$$\pi(\tau^2) = \mathcal{IG}\left(\tau^2; \frac{a}{2}, \frac{b}{2}\right), \quad \pi(\mu) \propto 1.$$

- The full posterior distribution is given by

$$\begin{aligned}
 & \pi(\mu, \tau^2, \theta_{1:J} | y_{1:J}) \\
 \propto & \pi(\mu, \tau^2) \prod_{j=1}^J \pi(\theta_j | \mu, \tau^2) \prod_{j=1}^J f(y_j | \theta_j, n_j) \\
 \propto & \frac{1}{\tau^{J+a+2}} \exp\left(-\frac{b}{2\tau^2} - \sum_{j=1}^J \frac{(\theta_j - \mu)^2}{2\tau^2}\right) \\
 & \times \prod_{j=1}^J \left(\frac{\exp(\theta_j)}{1 + \exp(\theta_j)}\right)^{y_j} \left(\frac{1}{1 + \exp(\theta_j)}\right)^{n_j - y_j}
 \end{aligned}$$

- This distribution does not admit a closed-form expression and we are going to use the Gibbs sampler by decomposing the parameter space in  $J + 2$  blocks  $\mu, \tau^2, \theta_1, \theta_2, \dots, \theta_J$ .
- The Gibbs sampler will require being able to be able to sample from  $\pi(\mu | y_{1:J}, \tau^2, \theta_{1:J})$ ,  $\pi(\tau^2 | y_{1:J}, \mu, \theta_{1:J})$  and  $\pi(\theta_j | y_{1:J}, \mu, \tau^2)$ .

- Full conditional distribution for  $\mu$

$$\begin{aligned}\pi(\mu | y_{1:J}, \tau^2, \theta_{1:J}) &= \pi(\mu | \tau^2, \theta_{1:J}) \\ &= \mathcal{N}\left(\mu; J^{-1} \sum_{j=1}^J \theta_j, J^{-1} \tau^2\right).\end{aligned}$$

- Full conditional distribution for  $\tau^2$

$$\begin{aligned}\pi(\tau^2 | y_{1:J}, \mu, \theta_{1:J}) &= \pi(\tau^2 | \mu, \theta_{1:J}) \\ &= \mathcal{IG}\left(\tau^2; \frac{a+J}{2}, \frac{b + \sum_{j=1}^J (\theta_j - \mu)^2}{2}\right).\end{aligned}$$

- These two distributions are standard.

- Full conditional distribution for  $\theta_j$

$$\pi(\theta_j | y_{1:j}, \tau^2, \mu) \propto \frac{\exp(\theta_j)^{y_j}}{(1 + \exp(\theta_j))^{n_j}} \exp\left(-\frac{(\theta_j - \mu)^2}{2\tau^2}\right).$$

- This distribution does not admit a closed-form expression so we use a Metropolis algorithm where

$$q(\theta'_j | \theta_j) = \mathcal{U}(\theta'_j; [\theta_j - \delta, \theta_j + \delta])$$

where  $\delta$  is a parameter of the algorithm.

- We run  $N = 100,000$  iterations of the Metropolis-Hastings one-at-a time and discard the first  $N_0 = 10,000$  samples. The acceptance rates for the Metropolis algorithm was given by

$\delta$	0.1	0.5	1	2	10
Average acceptance proba.	0.87	0.55	0.35	0.17	0.04

- Based on these results we computed the predictive distribution of  $Z_i$  the number of homes runs in the full season where  $m_j$  is the number of bats in the full seasons. We assume

$$Z_i \sim \text{Bin}(m_j, p_j) .$$

- The predictive distribution is given by

$$\Pr(Z_i = k | y_{1:J}, n_{1:J}, m_{1:J}) = \int \Pr(Z_i = k | m_j, p_j) \pi(p_j | y_{1:J}, n_{1:J}) dp_j$$

and can be approximated through the output of the MCMC algorithm

$$\frac{1}{N - N_0 + 1} \sum_{i=N_0}^N \Pr(Z_i = k | m_j^{(i)}, p_j^{(i)}) .$$

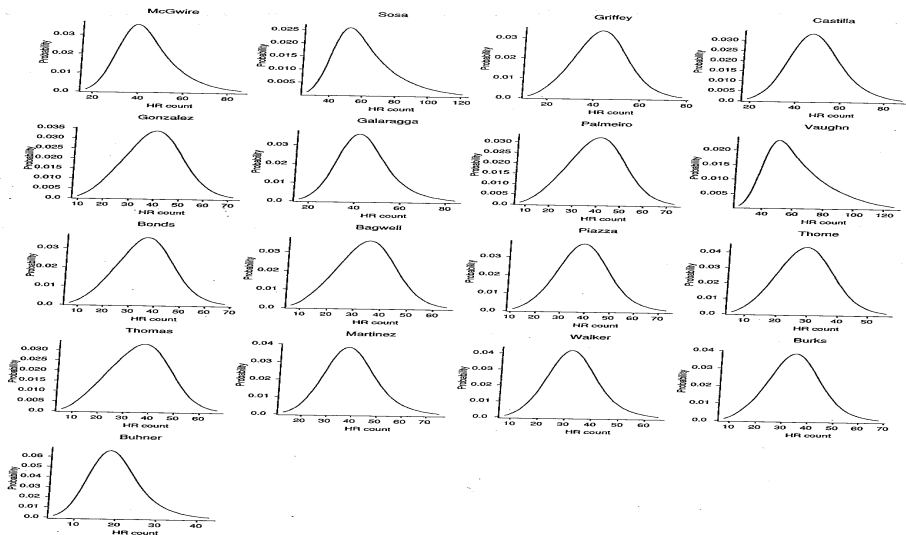


Figure: Predictive distribution for the number of home runs scored by each batter.