

CPSC 535

Metropolis-Hastings: Applications

AD

March 2007

- Initialization: Select deterministically or randomly
 $\theta = (\theta_1^{(0)}, \dots, \theta_p^{(0)})$.

Metropolis-Hastings one-at-a time

- Initialization: Select deterministically or randomly
 $\theta = (\theta_1^{(0)}, \dots, \theta_p^{(0)})$.
- Iteration i ; $i \geq 1$:

Metropolis-Hastings one-at-a time

- Initialization: Select deterministically or randomly
 $\theta = \left(\theta_1^{(0)}, \dots, \theta_p^{(0)} \right)$.
- Iteration i ; $i \geq 1$:
 - For $k = 1 : p$

- Initialization: Select deterministically or randomly

$$\theta = \left(\theta_1^{(0)}, \dots, \theta_p^{(0)} \right).$$

- Iteration i ; $i \geq 1$:

- For $k = 1 : p$

- Sample $\theta_k^{(i)}$ using an MH step of proposal distribution $q_k \left(\left(\theta_{-k}^{(i)}, \theta_k^{(i-1)} \right), \theta_k' \right)$ and target $\pi \left(\theta_k | \theta_{-k}^{(i)} \right)$ where $\theta_{-k}^{(i)} = \left(\theta_1^{(i)}, \dots, \theta_{k-1}^{(i)}, \theta_{k+1}^{(i-1)}, \dots, \theta_p^{(i-1)} \right)$.

- In 1986, Challenger exploded; the explosion being the result of an O-ring failure. It was believed to be a result of a cold weather at the departure time: 31°F.

- In 1986, Challenger exploded; the explosion being the result of an O-ring failure. It was believed to be a result of a cold weather at the departure time: 31°F .
- We have access to the data of 23 previous flights which give for flight i : Temperature at flight time x_i and $y_i = 1$ failure and zero otherwise (Robert & Casella, p. 15).

- In 1986, Challenger exploded; the explosion being the result of an O-ring failure. It was believed to be a result of a cold weather at the departure time: 31°F.
- We have access to the data of 23 previous flights which give for flight i : Temperature at flight time x_i and $y_i = 1$ failure and zero otherwise (Robert & Casella, p. 15).
- We want to have a model relating Y to x . Obviously this cannot be a linear model $Y = \alpha + x\beta$ as we want $Y \in \{0, 1\}$.

- We select a simple logistic regression model

$$\Pr(Y = 1|x) = 1 - \Pr(Y = 0|x) = \frac{\exp(\alpha + x\beta)}{1 + \exp(\alpha + x\beta)}.$$

- We select a simple logistic regression model

$$\Pr(Y = 1|x) = 1 - \Pr(Y = 0|x) = \frac{\exp(\alpha + x\beta)}{1 + \exp(\alpha + x\beta)}.$$

- Equivalently we have

$$\text{logit} = \log\left(\frac{\Pr(Y = 1|x)}{\Pr(Y = 0|x)}\right) = \alpha + x\beta.$$

- We select a simple logistic regression model

$$\Pr(Y = 1|x) = 1 - \Pr(Y = 0|x) = \frac{\exp(\alpha + x\beta)}{1 + \exp(\alpha + x\beta)}.$$

- Equivalently we have

$$\text{logit} = \log\left(\frac{\Pr(Y = 1|x)}{\Pr(Y = 0|x)}\right) = \alpha + x\beta.$$

- This ensures that the response is binary.

- We follow a Bayesian approach and select

$\pi(\alpha, \beta) = \pi(\alpha|b) \pi(\beta) = b^{-1} \exp(\alpha) \exp(-b^{-1} \exp(\alpha))$; i.e. exponential prior on $\exp(\alpha)$ and flat prior on β .

- We follow a Bayesian approach and select $\pi(\alpha, \beta) = \pi(\alpha|b) \pi(\beta) = b^{-1} \exp(\alpha) \exp(-b^{-1} \exp(\alpha))$; i.e. exponential prior on $\exp(\alpha)$ and flat prior on β .
- b is selected as the data-dependent prior such that $\mathbb{E}(\alpha) = \hat{\alpha}$ where $\hat{\alpha}$ is the MLE of α (Robert & Casella).

- We follow a Bayesian approach and select $\pi(\alpha, \beta) = \pi(\alpha|b)\pi(\beta) = b^{-1} \exp(\alpha) \exp(-b^{-1} \exp(\alpha))$; i.e. exponential prior on $\exp(\alpha)$ and flat prior on β .
- b is selected as the data-dependent prior such that $\mathbb{E}(\alpha) = \hat{\alpha}$ where $\hat{\alpha}$ is the MLE of α (Robert & Casella).
- As a simple proposal distribution, we use

$$q((\alpha, \beta), (\alpha', \beta')) = \pi(\alpha'|b) \mathcal{N}(\beta'; \beta, \hat{\sigma}_\beta^2)$$

where $\hat{\sigma}_\beta^2$ is the variance associated to the MLE $\hat{\beta}$.

The algorithm proceeds as follows at iteration i

- Sample $(\alpha^*, \beta^*) \sim \pi(\alpha | b) \mathcal{N}(\beta; \beta^{(i-1)}, \hat{\sigma}_\beta^2)$ and compute

$$\begin{aligned} & \zeta\left(\left(\alpha^{(i-1)}, \beta^{(i-1)}\right), \left(\alpha^*, \beta^*\right)\right) \\ = & \min\left(1, \frac{\pi(\alpha^*, \beta^* | \text{data}) \pi(\alpha^{(i-1)} | b)}{\pi(\alpha^{(i-1)}, \beta^{(i-1)} | \text{data}) \pi(\alpha^* | b)}\right) \end{aligned}$$

The algorithm proceeds as follows at iteration i

- Sample $(\alpha^*, \beta^*) \sim \pi(\alpha | b) \mathcal{N}(\beta; \beta^{(i-1)}, \hat{\sigma}_\beta^2)$ and compute

$$\begin{aligned} & \zeta\left(\left(\alpha^{(i-1)}, \beta^{(i-1)}\right), \left(\alpha^*, \beta^*\right)\right) \\ &= \min\left(1, \frac{\pi(\alpha^*, \beta^* | \text{data}) \pi(\alpha^{(i-1)} | b)}{\pi(\alpha^{(i-1)}, \beta^{(i-1)} | \text{data}) \pi(\alpha^* | b)}\right) \end{aligned}$$

- Set $(\alpha^{(i)}, \beta^{(i)}) = (\alpha^*, \beta^*)$ with probability $\zeta\left(\left(\alpha^{(i-1)}, \beta^{(i-1)}\right), \left(\alpha^*, \beta^*\right)\right)$, otherwise set $(\alpha^{(i)}, \beta^{(i)}) = (\alpha^{(i-1)}, \beta^{(i-1)})$.

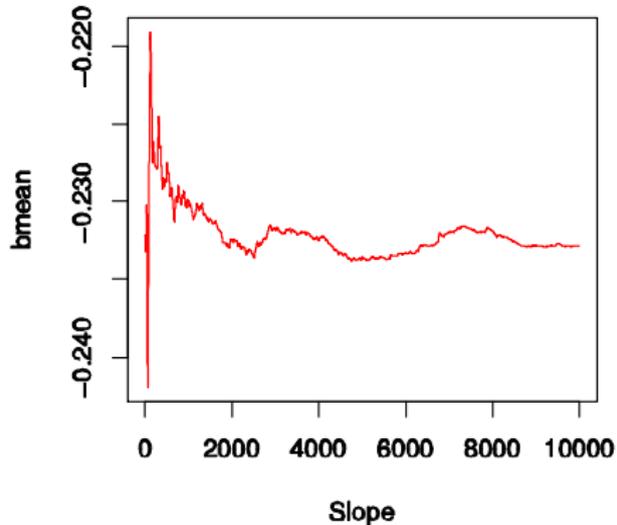
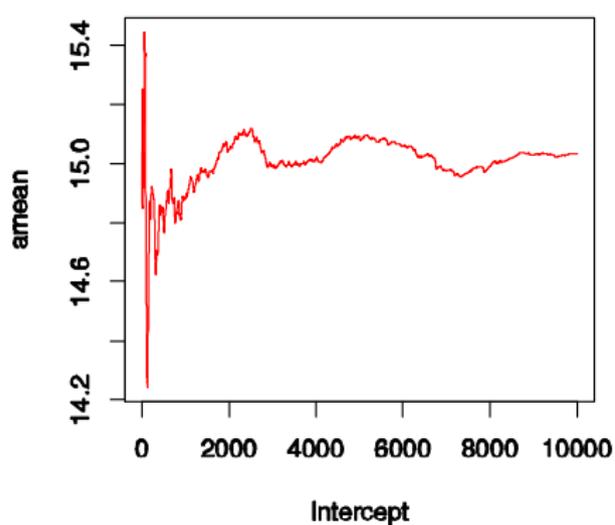


Figure: Plots of $\frac{1}{k} \sum_{i=1}^k \alpha^{(i)}$ (left) and $\frac{1}{k} \sum_{i=1}^k \beta^{(i)}$ (right).

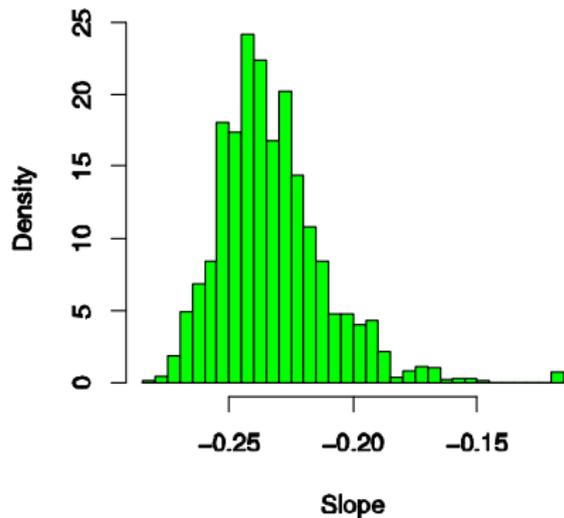
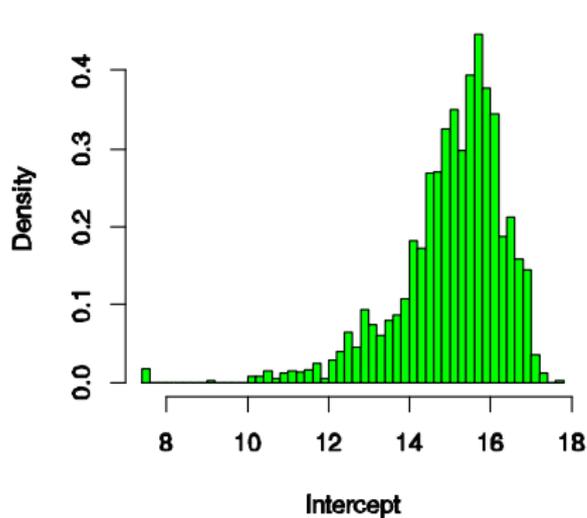


Figure: Histogram estimates of $p(\alpha | data)$ (left) and $p(\beta | data)$ (right).

- We consider the following example: we take 4 measurements from 100 genuine Swiss banknotes and 100 counterfeit ones.

- We consider the following example: we take 4 measurements from 100 genuine Swiss banknotes and 100 counterfeit ones.
- The response variable y is 0 for genuine and 1 for counterfeit and the explanatory variables are

- We consider the following example: we take 4 measurements from 100 genuine Swiss banknotes and 100 counterfeit ones.
- The response variable y is 0 for genuine and 1 for counterfeit and the explanatory variables are
 - x^1 the length,

- We consider the following example: we take 4 measurements from 100 genuine Swiss banknotes and 100 counterfeit ones.
- 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

- We consider the following example: we take 4 measurements from 100 genuine Swiss banknotes and 100 counterfeit ones.
- 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

- We consider the following example: we take 4 measurements from 100 genuine Swiss banknotes and 100 counterfeit ones.
- 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

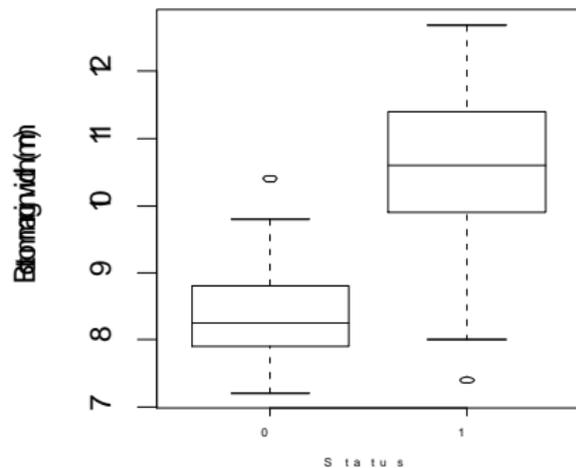
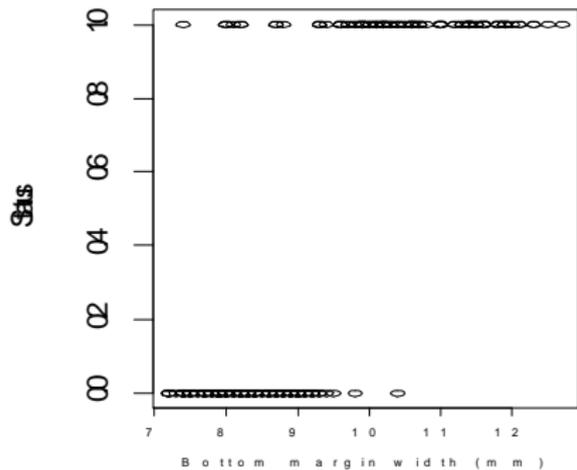


Figure: Left: Plot of the status indicator versus the bottom margin width. Right: Boxplots of the bottom margin width for both counterfeit status.

- Instead of selecting a logistic link, we select a probit one here

$$\Pr(Y = 1|x) = \Phi(x^1\beta_1 + \dots + x^4\beta_4)$$

where

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

- Instead of selecting a logistic link, we select a probit one here

$$\Pr(Y = 1 | x) = \Phi(x^1 \beta_1 + \dots + x^4 \beta_4)$$

where

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

- For n data, the likelihood is then given by

$$f(y_{1:n} | \beta, x_{1:n}) = \prod_{i=1}^n \Phi(x_i^T \beta)^{y_i} (1 - \Phi(x_i^T \beta))^{1-y_i}.$$

- We assume a vague prior where $\beta \sim \mathcal{N}(0, 100I_4)$ and we use a simple random walk sampler with $\hat{\Sigma}$ the covariance matrix associated to the MLE (estimated using simple deterministic method).

- We assume a vague prior where $\beta \sim \mathcal{N}(0, 100I_4)$ and we use a simple random walk sampler with $\hat{\Sigma}$ the covariance matrix associated to the MLE (estimated using simple deterministic method).
- The algorithm is thus simply given at iteration i by

- We assume a vague prior where $\beta \sim \mathcal{N}(0, 100I_4)$ and we use a simple random walk sampler with $\hat{\Sigma}$ the covariance matrix associated to the MLE (estimated using simple deterministic method).
- The algorithm is thus simply given at iteration i by
 - Sample $\beta^* \sim \mathcal{N}(\beta^{(i-1)}, \tau^2 \hat{\Sigma})$ and compute

$$\alpha(\beta^{(i-1)}, \beta^*) = \min \left(1, \frac{\pi(\beta^* | y_{1:n}, x_{1:n})}{\pi(\beta^{(i-1)} | y_{1:n}, x_{1:n})} \right).$$

- We assume a vague prior where $\beta \sim \mathcal{N}(0, 100I_4)$ and we use a simple random walk sampler with $\hat{\Sigma}$ the covariance matrix associated to the MLE (estimated using simple deterministic method).
- The algorithm is thus simply given at iteration i by
 - Sample $\beta^* \sim \mathcal{N}(\beta^{(i-1)}, \tau^2 \hat{\Sigma})$ and compute

$$\alpha(\beta^{(i-1)}, \beta^*) = \min \left(1, \frac{\pi(\beta^* | y_{1:n}, x_{1:n})}{\pi(\beta^{(i-1)} | y_{1:n}, x_{1:n})} \right).$$

- Set $\beta^{(i)} = \beta^*$ with probability $\alpha(\beta^{(i-1)}, \beta^*)$ and $\beta^{(i)} = \beta^{(i-1)}$ otherwise.

- We assume a vague prior where $\beta \sim \mathcal{N}(0, 100I_4)$ and we use a simple random walk sampler with $\hat{\Sigma}$ the covariance matrix associated to the MLE (estimated using simple deterministic method).
- The algorithm is thus simply given at iteration i by
 - Sample $\beta^* \sim \mathcal{N}(\beta^{(i-1)}, \tau^2 \hat{\Sigma})$ and compute

$$\alpha(\beta^{(i-1)}, \beta^*) = \min \left(1, \frac{\pi(\beta^* | y_{1:n}, x_{1:n})}{\pi(\beta^{(i-1)} | y_{1:n}, x_{1:n})} \right).$$

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

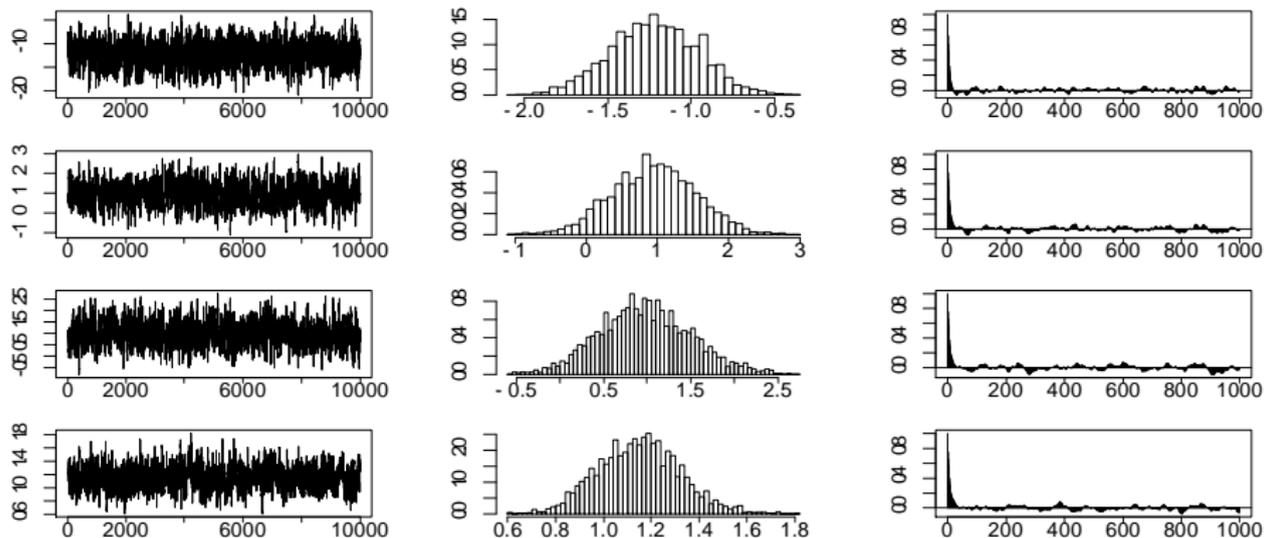


Figure: Traces (left), Histograms (middle) and Autocorrelations (right) for $(\beta_1^{(i)}, \dots, \beta_4^{(i)})$.

- One way to monitor the performance of the algorithm of the chain $\{X^{(i)}\}$ consists of displaying $\rho_k = \text{cov} [X^{(i)}, X^{(i+k)}] / \text{var} (X^{(i)})$ which can be estimated from the chain, at least for small values of k .

- One way to monitor the performance of the algorithm of the chain $\{X^{(i)}\}$ consists of displaying $\rho_k = \text{cov} [X^{(i)}, X^{(i+k)}] / \text{var} (X^{(i)})$ which can be estimated from the chain, at least for small values of k .
- Sometimes one uses an effective sample size measure

$$N^{\text{ess}} = N \left(1 + 2 \sum_{k=1}^{N_0} \hat{\rho}_k \right)^{-1/2} .$$

This represents approximately the sample size of an equivalent i.i.d. samples.

- One way to monitor the performance of the algorithm of the chain $\{X^{(i)}\}$ consists of displaying $\rho_k = \text{cov} [X^{(i)}, X^{(i+k)}] / \text{var} (X^{(i)})$ which can be estimated from the chain, at least for small values of k .
- Sometimes one uses an effective sample size measure

$$N^{\text{ess}} = N \left(1 + 2 \sum_{k=1}^{N_0} \hat{\rho}_k \right)^{-1/2} .$$

This represents approximately the sample size of an equivalent i.i.d. samples.

- One should be very careful with such measures which can be very misleading.

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

$$\hat{p} = \Phi(-1.22x^1 + 0.95x^2 + 0.96x^3 + 1.15x^4)$$

For $x = (214.9, 130.1, 129.9, 9.5)$, we obtain $\hat{p} = 0.59$.

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

$$\hat{p} = \Phi(-1.22x^1 + 0.95x^2 + 0.96x^3 + 1.15x^4)$$

For $x = (214.9, 130.1, 129.9, 9.5)$, we obtain $\hat{p} = 0.59$.

- A better estimate is obtained by

$$\int \Phi(\beta_1 x^1 + \beta_2 x^2 + \beta_3 x^3 + \beta_4 x^4) \pi(\beta | y_{1:n}, x_{1:n}) d\beta$$

Auxiliary Variables for Probit Regression

- It is impossible to use Gibbs to sample directly from $\pi(\beta | y_{1:n}, x_{1:n})$.

Auxiliary Variables for Probit Regression

- It is impossible to use Gibbs to sample directly from $\pi(\beta | y_{1:n}, x_{1:n})$.
- Introduce the following unobserved latent variables

$$Z_i \sim \mathcal{N}(x_i^T \beta, 1),$$
$$Y_i = \begin{cases} 1 & \text{if } Z_i > 0 \\ 0 & \text{otherwise.} \end{cases}$$

Auxiliary Variables for Probit Regression

- It is impossible to use Gibbs to sample directly from $\pi(\beta | y_{1:n}, x_{1:n})$.
- Introduce the following unobserved latent variables

$$Z_i \sim \mathcal{N}(x_i^\top \beta, 1),$$
$$Y_i = \begin{cases} 1 & \text{if } Z_i > 0 \\ 0 & \text{otherwise.} \end{cases}$$

- We have now define a joint distribution

$$f(y_i, z_i | \beta, x_i) = f(y_i | z_i) f(z_i | \beta, x_i).$$

- Now we can check that

$$\begin{aligned} f(y_i = 1 | x_i, \beta) &= \int f(y_i, z_i | \beta, x_i) dz_i \\ &= \int_0^{\infty} f(z_i | \beta, x_i) dz_i = \Phi(x_i^T \beta). \end{aligned}$$

⇒ We haven't changed the model!

- Now we can check that

$$\begin{aligned} f(y_i = 1 | x_i, \beta) &= \int f(y_i, z_i | \beta, x_i) dz_i \\ &= \int_0^\infty f(z_i | \beta, x_i) dz_i = \Phi(x_i^\top \beta). \end{aligned}$$

⇒ We haven't changed the model!

- We are now going to sample from $\pi(\beta, z_{1:n} | x_{1:n}, y_{1:n})$ instead of $\pi(\beta | x_{1:n}, y_{1:n})$ because the full conditional distributions are simple

$$\pi(\beta | y_{1:n}, x_{1:n}, z_{1:n}) = \pi(\beta | x_{1:n}, z_{1:n}) \text{ (standard Gaussian!),}$$

$$\pi(z_{1:n} | y_{1:n}, x_{1:n}, \beta) = \prod_{i=1}^n \pi(z_i | y_i, x_i, \beta)$$

where

$$z_k | y_k, x_k, \beta \sim \begin{cases} \mathcal{N}_+ (x_k^\top \beta, 1) & \text{if } y_k = 1 \\ \mathcal{N}_- (x_k^\top \beta, 1) & \text{if } y_k = 0. \end{cases}$$

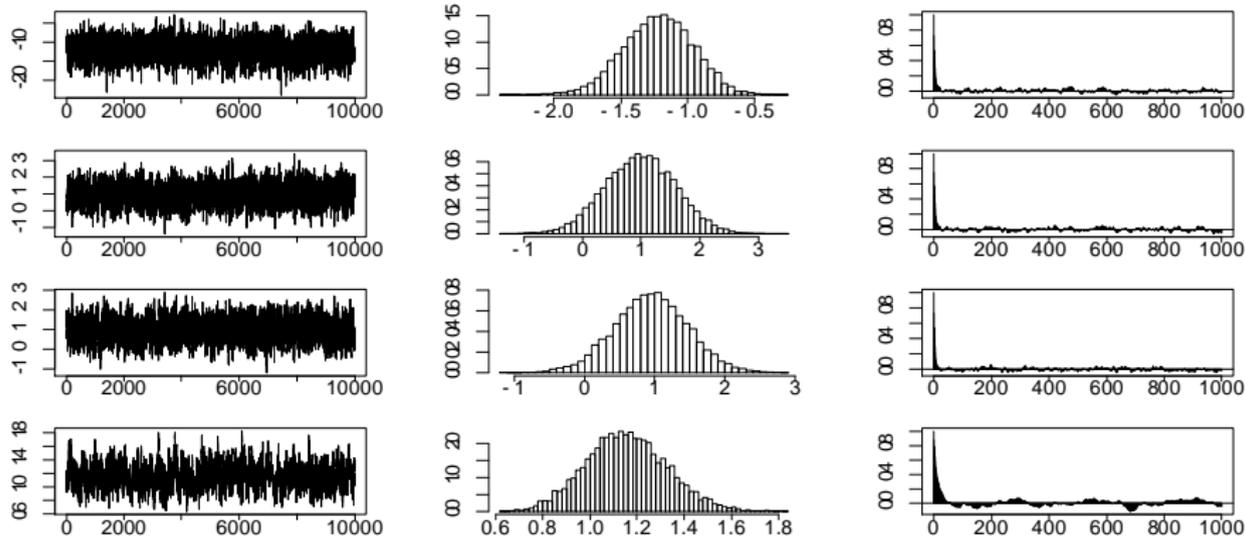


Figure: Traces (left), Histograms (middle) and Autocorrelations (right) for $(\beta_1^{(i)}, \dots, \beta_4^{(i)})$.

- The results obtained through Gibbs are very similar to MH.

- The results obtained through Gibbs are very similar to MH.
- We can also adopt an Zellner's type prior and obtain very similar results.

- The results obtained through Gibbs are very similar to MH.
- We can also adopt an Zellner's type prior and obtain very similar results.
- Very similar were also obtained using a logistic function using the MH (Gibbs is feasible but more difficult).

Warning

- Although the introduction of latent variables can be attractive, it can be also very inefficient.

Warning

- Although the introduction of latent variables can be attractive, it can be also very inefficient.
- It is not because you can use the Gibbs sampler that everything works well!

Warning

- Although the introduction of latent variables can be attractive, it can be also very inefficient.
- It is not because you can use the Gibbs sampler that everything works well!
- Consider the following simple generalization of the previous model

$$Z_i \sim \mathcal{N}(x_i\beta, \sigma^2), \quad Y_i = \begin{cases} 1 & \text{if } Z_i > 0 \\ 0 & \text{otherwise.} \end{cases}$$

Warning

- Although the introduction of latent variables can be attractive, it can be also very inefficient.
- It is not because you can use the Gibbs sampler that everything works well!
- Consider the following simple generalization of the previous model

$$Z_i \sim \mathcal{N}(x_i \beta, \sigma^2), \quad Y_i = \begin{cases} 1 & \text{if } Z_i > 0 \\ 0 & \text{otherwise.} \end{cases}$$

- We complete the model by $\sigma^2 \sim \mathcal{IG}(1.5, 1.5)$ and $\beta | \sigma^2 \sim \mathcal{N}(0, 100)$.

- Not only the data Z_i and (β, σ^2) are very correlated but we have

$$\Pr(Y_i = 1 | x_i, \beta, \sigma^2) = \Phi\left(\frac{x_i \beta}{\sigma}\right)$$

- Not only the data Z_i and (β, σ^2) are very correlated but we have

$$\Pr(Y_i = 1 | x_i, \beta, \sigma^2) = \Phi\left(\frac{x_i \beta}{\sigma}\right)$$

- The likelihood only depends on β/σ and the parameters β and σ are not identifiable.

- Not only the data Z_i and (β, σ^2) are very correlated but we have

$$\Pr(Y_i = 1 | x_i, \beta, \sigma^2) = \Phi\left(\frac{x_i \beta}{\sigma}\right)$$

- The likelihood only depends on β/σ and the parameters β and σ are not identifiable.
- One way to improve the mixing consists of using an additional MH step that proposes to randomly rescale the current value.

- Not only the data Z_i and (β, σ^2) are very correlated but we have

$$\Pr(Y_i = 1 | x_i, \beta, \sigma^2) = \Phi\left(\frac{x_i \beta}{\sigma}\right)$$

- The likelihood only depends on β/σ and the parameters β and σ are not identifiable.
- One way to improve the mixing consists of using an additional MH step that proposes to randomly rescale the current value.
- We use a proposal distribution such that

$$(\beta', \sigma') = \lambda (\beta, \sigma) \text{ with } \lambda \sim \mathcal{Exp}(1)$$

that proposes to randomly rescale the current value.

Hidden Markov Model

- Consider the following hidden Markov model

$$\begin{aligned} X_k | (X_{k-1} = x_{k-1}) &\sim f_{\theta}(\cdot | x_{k-1}), \quad X_1 \sim \mu \\ Y_n | (X_k = x_k) &\sim g_{\theta}(\cdot | x_k), \end{aligned}$$

and we set a prior $\pi(\theta)$ on the unknown hyperparameters θ .

Hidden Markov Model

- Consider the following hidden Markov model

$$\begin{aligned}X_k | (X_{k-1} = x_{k-1}) &\sim f_\theta(\cdot | x_{k-1}), \quad X_1 \sim \mu \\ Y_n | (X_k = x_k) &\sim g_\theta(\cdot | x_k),\end{aligned}$$

and we set a prior $\pi(\theta)$ on the unknown hyperparameters θ .

- Given n data, we are interested in the joint posterior

$$\pi(\theta, x_{1:n} | y_{1:n})$$

Hidden Markov Model

- Consider the following hidden Markov model

$$\begin{aligned}X_k | (X_{k-1} = x_{k-1}) &\sim f_\theta(\cdot | x_{k-1}), \quad X_1 \sim \mu \\ Y_n | (X_k = x_k) &\sim g_\theta(\cdot | x_k),\end{aligned}$$

and we set a prior $\pi(\theta)$ on the unknown hyperparameters θ .

- Given n data, we are interested in the joint posterior

$$\pi(\theta, x_{1:n} | y_{1:n})$$

- There is no closed-form expression for this joint distribution even if the model is linear Gaussian or for finite state-space model.

- In previous lectures, we propose sampling from $\pi(\theta, x_{1:n} | y_{1:n})$ using the Gibbs sampler where the variables are updated according to

$$X_k \sim \pi(x_k | y_{1:n}, x_{-k}, \theta)$$

- In previous lectures, we propose sampling from $\pi(\theta, x_{1:n} | y_{1:n})$ using the Gibbs sampler where the variables are updated according to

$$X_k \sim \pi(x_k | y_{1:n}, x_{-k}, \theta)$$

- For $2 < k < n$, we have

$$\begin{aligned} \pi(x_k | y_{1:n}, x_{-k}, \theta) &\propto \pi(x_{1:n}, y_{1:n}, \theta) \\ &\propto \underbrace{\pi(\theta) \mu(x_1) \prod_{i=2}^n f_\theta(x_i | x_{i-1})}_{\text{prior}} \underbrace{\prod_{i=1}^n g_\theta(y_i | x_i)}_{\text{likelihood}} \\ &\propto f_\theta(x_k | x_{k-1}) f_\theta(x_{k+1} | x_k) g_\theta(y_k | x_k) \end{aligned}$$

and $\theta \sim \pi(\theta | y_{1:n}, x_{1:n})$ (or by subblocks).

- It is often possible to implement the Gibbs sampler even if this can be expensive; e.g. if you use Accept/Reject to sample from $\pi(x_k | y_{1:n}, x_{-k}, \theta)$ using the proposal $\pi(x_k | x_{-k}, \theta) \propto f_\theta(x_k | x_{k-1}) f_\theta(x_{k+1} | x_k)$.

- It is often possible to implement the Gibbs sampler even if this can be expensive; e.g. if you use Accept/Reject to sample from $\pi(x_k | y_{1:n}, x_{-k}, \theta)$ using the proposal $\pi(x_k | x_{-k}, \theta) \propto f_\theta(x_k | x_{k-1}) f_\theta(x_{k+1} | x_k)$.
- Even if it is possible to implement the Gibbs sampler, one can expect a very slow convergence of the algorithm if successive variables are highly correlated.

- It is often possible to implement the Gibbs sampler even if this can be expensive; e.g. if you use Accept/Reject to sample from $\pi(x_k | y_{1:n}, x_{-k}, \theta)$ using the proposal $\pi(x_k | x_{-k}, \theta) \propto f_\theta(x_k | x_{k-1}) f_\theta(x_{k+1} | x_k)$.
- Even if it is possible to implement the Gibbs sampler, one can expect a very slow convergence of the algorithm if successive variables are highly correlated.
- Indeed, as you update x_k with x_{k-1} and x_{k+1} being fixed, then you cannot move much into the space.

- Consider the very simple case where $\theta = (\sigma_v^2, \sigma_w^2)$

$$X_k = X_{k-1} + V_k \text{ where } V_k \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma_v^2),$$

$$Y_k = X_k + W_k \text{ where } W_k \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma_w^2)$$

then we have

$$\begin{aligned} \pi(x_k | x_{-k}, \theta) &\propto f_\theta(x_k | x_{k-1}) f_\theta(x_{k+1} | x_k) \\ &= \mathcal{N}\left(x_k; \frac{x_{k-1} + x_{k+1}}{2}, \frac{\sigma_v^2}{2}\right) \end{aligned}$$

and

$$\begin{aligned} &\pi(x_k | y_{1:n}, x_{-k}, \theta) \\ &\propto \pi(x_k | x_{-k}, \theta) g_\theta(y_k | x_k) \\ &= \mathcal{N}\left(x_k; \frac{\sigma_v^2 \sigma_w^2}{\sigma_v^2 + 2\sigma_w^2} \left(\frac{x_{k-1} + x_{k+1}}{\sigma_v^2} + \frac{y_k}{\sigma_w^2}\right), \frac{\sigma_v^2 \sigma_w^2}{\sigma_v^2 + 2\sigma_w^2}\right) \end{aligned}$$

- Assume for the time being that instead of sampling from $\pi(x_k | y_{1:n}, x_{-k}, \theta)$ directly, we use rejection sampling with $\pi(x_k | x_{-k}, \theta)$ as a proposal distribution.

- Assume for the time being that instead of sampling from $\pi(x_k | y_{1:n}, x_{-k}, \theta)$ directly, we use rejection sampling with $\pi(x_k | x_{-k}, \theta)$ as a proposal distribution.
- In this case we have to bound

$$g_{\theta}(y_k | x_k) = \frac{1}{\sqrt{2\pi}\sigma_w} \exp\left(-\frac{(y_k - x_k)^2}{2\sigma_w^2}\right) \leq \frac{1}{\sqrt{2\pi}\sigma_w}.$$

- Assume for the time being that instead of sampling from $\pi(x_k | y_{1:n}, x_{-k}, \theta)$ directly, we use rejection sampling with $\pi(x_k | x_{-k}, \theta)$ as a proposal distribution.
- In this case we have to bound

$$g_{\theta}(y_k | x_k) = \frac{1}{\sqrt{2\pi}\sigma_w} \exp\left(-\frac{(y_k - x_k)^2}{2\sigma_w^2}\right) \leq \frac{1}{\sqrt{2\pi}\sigma_w}.$$

- We accept each proposal $X^* \sim \pi(x_k | x_{-k}, \theta)$ with probability $\exp\left(-\frac{(y_k - X^*)^2}{2\sigma_w^2}\right)$, so the (unconditional) acceptance probability is given by

$$\begin{aligned} & \int \pi(x_k | x_{-k}, \theta) \exp\left(-\frac{(y_k - x_k)^2}{2\sigma_w^2}\right) dx_k \\ = & \frac{\sigma_w \exp\left(-\frac{1}{2}\left(y_k^2/\sigma_w^2 - (x_{k-1} + x_{k+1})^2/\sigma_v^2\right)\right)}{\sqrt{\sigma_v^2 + 2\sigma_w^2}}. \end{aligned}$$

- To improve the algorithm, we would like to be able to sample a whole block of variables simultaneously; i.e. being able to sample for $1 < k < k + L < n$ from

$$\begin{aligned} \pi \left(x_{k:k+L} \mid y_{1:n}, x_{-(k:k+L)}, \theta \right) &= \pi \left(x_{k:k+L} \mid y_{k:k+L}, x_{k-1}, x_{k+L+1}, \theta \right) \\ &\propto \prod_{i=k}^{k+L+1} f_{\theta} \left(x_i \mid x_{i-1} \right) \prod_{i=k}^{k+L} g_{\theta} \left(y_i \mid x_i \right). \end{aligned}$$

- To improve the algorithm, we would like to be able to sample a whole block of variables simultaneously; i.e. being able to sample for $1 < k < k + L < n$ from

$$\begin{aligned} \pi \left(x_{k:k+L} \mid y_{1:n}, x_{-(k:k+L)}, \theta \right) &= \pi \left(x_{k:k+L} \mid y_{k:k+L}, x_{k-1}, x_{k+L+1}, \theta \right) \\ &\propto \prod_{i=k}^{k+L+1} f_{\theta} \left(x_i \mid x_{i-1} \right) \prod_{i=k}^{k+L} g_{\theta} \left(y_i \mid x_i \right). \end{aligned}$$

- In this case, it is typically impossible to sample from $\pi \left(x_{k:k+L} \mid y_{1:n}, x_{-(k:k+L)}, \theta \right)$ exactly as L is large, say 5 or 10.

- To improve the algorithm, we would like to be able to sample a whole block of variables simultaneously; i.e. being able to sample for $1 < k < k + L < n$ from

$$\begin{aligned} \pi \left(x_{k:k+L} \mid y_{1:n}, x_{-(k:k+L)}, \theta \right) &= \pi \left(x_{k:k+L} \mid y_{k:k+L}, x_{k-1}, x_{k+L+1}, \theta \right) \\ &\propto \prod_{i=k}^{k+L+1} f_{\theta} \left(x_i \mid x_{i-1} \right) \prod_{i=k}^{k+L} g_{\theta} \left(y_i \mid x_i \right). \end{aligned}$$

- In this case, it is typically impossible to sample from $\pi \left(x_{k:k+L} \mid y_{1:n}, x_{-(k:k+L)}, \theta \right)$ exactly as L is large, say 5 or 10.
- We propose to use a MH step of invariant distribution $\pi \left(x_{k:k+L} \mid y_{1:n}, x_{-(k:k+L)}, \theta \right)$ instead, hence we need to build a proposal distribution $q \left((x_{1:n}, \theta), x'_{k:k+L} \right)$.

- We first propose to use the conditional prior

$$\begin{aligned} q((x_{1:n}, \theta), x'_{k:k+L}) &= \pi(x_{k:k+L} | x_{-(k:k+L)}, \theta) \\ &= \pi(x_{k:k+L} | x_{k-1}, x_{k+L+1}, \theta) \\ &\propto \prod_{i=k}^{k+L+1} f_{\theta}(x_i | x_{i-1}). \end{aligned}$$

- We first propose to use the conditional prior

$$\begin{aligned}
 q((x_{1:n}, \theta), x'_{k:k+L}) &= \pi(x_{k:k+L} | x_{-(k:k+L)}, \theta) \\
 &= \pi(x_{k:k+L} | x_{k-1}, x_{k+L+1}, \theta) \\
 &\propto \prod_{i=k}^{k+L+1} f_{\theta}(x_i | x_{i-1}).
 \end{aligned}$$

- In this case, the candidate $X'_{k:k+L} \sim \pi(x_{k:k+L} | x_{k-1}, x_{k+L+1}, \theta)$ is accepted with probability

$$\begin{aligned}
 &\min \left(1, \frac{\pi(x'_{k:k+L} | y_{k:k+L}, x_{k-1}, x_{k+L+1}, \theta) \pi(x_{k:k+L} | x_{k-1}, x_{k+L+1}, \theta)}{\pi(x_{k:k+L} | y_{k:k+L}, x_{k-1}, x_{k+L+1}, \theta) \pi(x'_{k:k+L} | x_{k-1}, x_{k+L+1}, \theta)} \right) \\
 &= \min \left(1, \frac{\prod_{i=k}^{k+L} g_{\theta}(y_i | x'_i)}{\prod_{i=k}^{k+L} g_{\theta}(y_i | x_i)} \right)
 \end{aligned}$$

- We first propose to use the conditional prior

$$\begin{aligned}
 q((x_{1:n}, \theta), x'_{k:k+L}) &= \pi(x_{k:k+L} | x_{-(k:k+L)}, \theta) \\
 &= \pi(x_{k:k+L} | x_{k-1}, x_{k+L+1}, \theta) \\
 &\propto \prod_{i=k}^{k+L+1} f_{\theta}(x_i | x_{i-1}).
 \end{aligned}$$

- In this case, the candidate $X'_{k:k+L} \sim \pi(x_{k:k+L} | x_{k-1}, x_{k+L+1}, \theta)$ is accepted with probability

$$\begin{aligned}
 &\min \left(1, \frac{\pi(x'_{k:k+L} | y_{k:k+L}, x_{k-1}, x_{k+L+1}, \theta) \pi(x_{k:k+L} | x_{k-1}, x_{k+L+1}, \theta)}{\pi(x_{k:k+L} | y_{k:k+L}, x_{k-1}, x_{k+L+1}, \theta) \pi(x'_{k:k+L} | x_{k-1}, x_{k+L+1}, \theta)} \right) \\
 &= \min \left(1, \frac{\prod_{i=k}^{k+L} g_{\theta}(y_i | x'_i)}{\prod_{i=k}^{k+L} g_{\theta}(y_i | x_i)} \right)
 \end{aligned}$$

- Simple but one cannot expect it to be too efficient when the observations are very informative compared to the prior.

- Consider the case where

$$X_k = AX_{k-1} + BV_k, V_k \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, I).$$

- Consider the case where

$$X_k = AX_{k-1} + BV_k, \quad V_k \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, I).$$

- Particular cases include

$$X_k = X_{k-1} + \sigma V_k, \quad \text{where } V_k \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, 1),$$

$$X_k = \begin{pmatrix} \alpha_k \\ \alpha_{k-1} \end{pmatrix} = \begin{pmatrix} 2 & -1 \\ 1 & 0 \end{pmatrix} X_{k-1} + \begin{pmatrix} \sigma \\ 0 \end{pmatrix} V_k, \quad V_k \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, 1)$$

- In this case, it is simple to see that $\pi(x_{k:k+L} | x_{k-1}, x_{k+1}, \theta)$ is a Gaussian distribution.

- In this case, it is simple to see that $\pi(x_{k:k+L} | x_{k-1}, x_{k+1}, \theta)$ is a Gaussian distribution.
- In (Knorr-Held, 1999), one samples from this distribution by computing directly the parameters of this joint distribution: complexity $O(L^2)$.

- In this case, it is simple to see that $\pi(x_{k:k+L} | x_{k-1}, x_{k+1}, \theta)$ is a Gaussian distribution.
- In (Knorr-Held, 1999), one samples from this distribution by computing directly the parameters of this joint distribution: complexity $O(L^2)$.
- We can derive a simpler method of complexity $O(L)$ based on the following decomposition (omitting θ in the notation)

$$\begin{aligned} \pi(x_{k:k+L} | x_{k-1}, x_{k+L+1}) &= \prod_{i=k}^{k+L} \pi(x_i | x_{k-1}, x_{k+L+1}, x_{i+1}) . \\ &= \prod_{i=k}^{k+L} \pi(x_i | x_{k-1}, x_{i+1}) \end{aligned}$$

- Moreover it is easy to establish the expression for $\pi(x_i | x_{k-1}, x_{i+1})$

$$\pi(x_i | x_{k-1}, x_{i+1}) \propto \pi(x_i | x_{k-1}) f(x_{i+1} | x_i)$$

as

$$\pi(x_i | x_{k-1}) = \int \pi(x_{k:i} | x_{k-1}) dx_{k:i-1} = \mathcal{N}(x_i; \mu_i(x_{k-1}), \Sigma_i)$$

with, for $X_n = AX_{n-1} + BV_n$, $\mu_{k-1}(x_{k-1}) = x_{k-1}$, $\Sigma_{k-1} = 0$ and for $i \geq k$

$$\begin{aligned} \mu_i(x_{k-1}) &= A\mu_{i-1}(x_{k-1}), \\ \Sigma_i &= A\Sigma_{i-1}A^T + \Sigma \text{ with } \Sigma = BB^T. \end{aligned}$$

- Moreover it is easy to establish the expression for $\pi(x_i | x_{k-1}, x_{i+1})$

$$\pi(x_i | x_{k-1}, x_{i+1}) \propto \pi(x_i | x_{k-1}) f(x_{i+1} | x_i)$$

as

$$\pi(x_i | x_{k-1}) = \int \pi(x_{k:i} | x_{k-1}) dx_{k:i-1} = \mathcal{N}(x_i; \mu_i(x_{k-1}), \Sigma_i)$$

with, for $X_n = AX_{n-1} + BV_n$, $\mu_{k-1}(x_{k-1}) = x_{k-1}$, $\Sigma_{k-1} = 0$ and for $i \geq k$

$$\begin{aligned} \mu_i(x_{k-1}) &= A\mu_{i-1}(x_{k-1}), \\ \Sigma_i &= A\Sigma_{i-1}A^T + \Sigma \text{ with } \Sigma = BB^T. \end{aligned}$$

- To obtain $\pi(x_i | x_{k-1}, x_{i+1})$, we combine the prior $\pi(x_i | x_{k-1})$ with the “likelihood” $f(x_{i+1} | x_i)$.

- We have $\pi(x_i | x_{k-1}) = \mathcal{N}(x_i; \mu_i(x_{k-1}), \Sigma_i)$ and $f(x_{i+1} | x_i) = \mathcal{N}(x_{i+1}; Ax_i, \Sigma)$ then

$$\pi(x_i | x_{k-1}, x_{i+1}) = \mathcal{N}(x_i; \mu_i(x_{k-1}, x_{i+1}), \tilde{\Sigma}_i)$$

where

$$\begin{aligned}\tilde{\Sigma}_i &= \left(\Sigma_i^{-1} + A^T \Sigma^{-1} A \right)^{-1}, \\ \mu_i(x_{k-1}, x_{i+1}) &= \tilde{\Sigma}_i \left(A^T \Sigma^{-1} x_{i+1} + \Sigma_i^{-1} \mu_i(x_{k-1}) \right).\end{aligned}$$

- We have $\pi(x_i | x_{k-1}) = \mathcal{N}(x_i; \mu_i(x_{k-1}), \Sigma_i)$ and $f(x_{i+1} | x_i) = \mathcal{N}(x_{i+1}; Ax_i, \Sigma)$ then

$$\pi(x_i | x_{k-1}, x_{i+1}) = \mathcal{N}(x_i; \mu_i(x_{k-1}, x_{i+1}), \tilde{\Sigma}_i)$$

where

$$\begin{aligned} \tilde{\Sigma}_i &= \left(\Sigma_i^{-1} + A^T \Sigma^{-1} A \right)^{-1}, \\ \mu_i(x_{k-1}, x_{i+1}) &= \tilde{\Sigma}_i \left(A^T \Sigma^{-1} x_{i+1} + \Sigma_i^{-1} \mu_i(x_{k-1}) \right). \end{aligned}$$

- To sample a realization of $\pi(x_{k:k+L} | x_{k-1}, x_{k+L+1})$, first compute $\mu_i(x_{k-1}), \Sigma_i$ for $i = k, \dots, k+L$ using a forward recursion. Then sample backward $X_{k+L} \sim \pi(\cdot | x_{k-1}, x_{k+L+1})$, $X_{k+L-1} \sim \pi(\cdot | x_{k-1}, X_{k+L})$, \dots , $X_k \sim \pi(\cdot | x_{k-1}, X_{k+1})$.

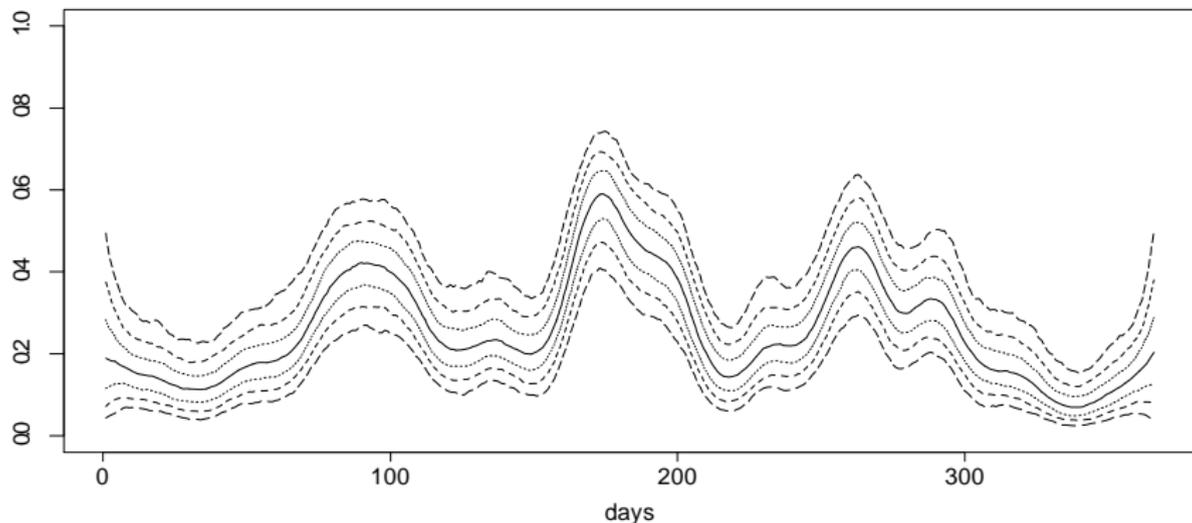


Figure: Number of occurrences of rainfall in Tokyo for each day during 1983-1984 reproduced as relative frequencies between 0, 0.5 and 1 ($n = 366$)

- Consider the following model

$$X_k = \begin{pmatrix} \alpha_k \\ \alpha_{k-1} \end{pmatrix} = \begin{pmatrix} 2 & -1 \\ 1 & 0 \end{pmatrix} X_{k-1} + \begin{pmatrix} \sigma \\ 0 \end{pmatrix} V_k, \quad V_k \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, 1)$$

and

$$Y_k | X_k \sim \begin{cases} B(2, \pi_k) & k \neq 60, \\ B(1, \pi_k) & k = 60 \text{ (February 29)} \end{cases},$$

where

$$\pi_k = \frac{\exp(\alpha_k)}{1 + \exp(\alpha_k)}.$$

- Consider the following model

$$X_k = \begin{pmatrix} \alpha_k \\ \alpha_{k-1} \end{pmatrix} = \begin{pmatrix} 2 & -1 \\ 1 & 0 \end{pmatrix} X_{k-1} + \begin{pmatrix} \sigma \\ 0 \end{pmatrix} V_k, \quad V_k \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, 1)$$

and

$$Y_k | X_k \sim \begin{cases} B(2, \pi_k) & k \neq 60, \\ B(1, \pi_k) & k = 60 \text{ (February 29)} \end{cases},$$

where

$$\pi_k = \frac{\exp(\alpha_k)}{1 + \exp(\alpha_k)}.$$

- We also use for $\sigma^2 \sim \mathcal{IG}(\frac{\nu_0}{2}, \frac{\gamma_0}{2})$.

- We use the block sampling strategies discussed before where candidates are sampled according to $\pi(x_{k:k+L} | x_{k-1}, x_{k+L+1})$ and accepted with proba

$$\min \left(1, \frac{\prod_{i=k}^{k+L} g(y_i | x'_i)}{\prod_{i=k}^{k+L} g(y_i | x_i)} \right).$$

- We use the block sampling strategies discussed before where candidates are sampled according to $\pi(x_{k:k+L} | x_{k-1}, x_{k+L+1})$ and accepted with proba

$$\min \left(1, \frac{\prod_{i=k}^{k+L} g(y_i | x'_i)}{\prod_{i=k}^{k+L} g(y_i | x_i)} \right).$$

- The parameter σ^2 is updated through a simple Gibbs step

$$\begin{aligned} \sigma^2 &\sim \pi(\sigma^2 | x_{1:n}, y_{1:n}) = \pi(\sigma^2 | x_{1:n}) \\ &= \mathcal{IG} \left(\frac{\nu_0 + n - 1}{2}, \frac{\gamma_0 + \sum_{k=2}^n (\alpha_k - 2\alpha_{k-1} + \alpha_{k-2})^2}{2} \right) \end{aligned}$$

- We use the block sampling strategies discussed before where candidates are sampled according to $\pi(x_{k:k+L} | x_{k-1}, x_{k+L+1})$ and accepted with proba

$$\min \left(1, \frac{\prod_{i=k}^{k+L} g(y_i | x'_i)}{\prod_{i=k}^{k+L} g(y_i | x_i)} \right).$$

- The parameter σ^2 is updated through a simple Gibbs step

$$\begin{aligned} \sigma^2 &\sim \pi(\sigma^2 | x_{1:n}, y_{1:n}) = \pi(\sigma^2 | x_{1:n}) \\ &= \mathcal{IG} \left(\frac{\nu_0 + n - 1}{2}, \frac{\gamma_0 + \sum_{k=2}^n (\alpha_k - 2\alpha_{k-1} + \alpha_{k-2})^2}{2} \right) \end{aligned}$$

- For block size $L = 1, 5, 20$ and 40 , we compute the average trajectories of 100 parallel chains after 10, 50, 100 and 500 iterations with initialization $x_k = 0$ for all $k, \sigma^2 = 0.1$.

After 10 Iterations

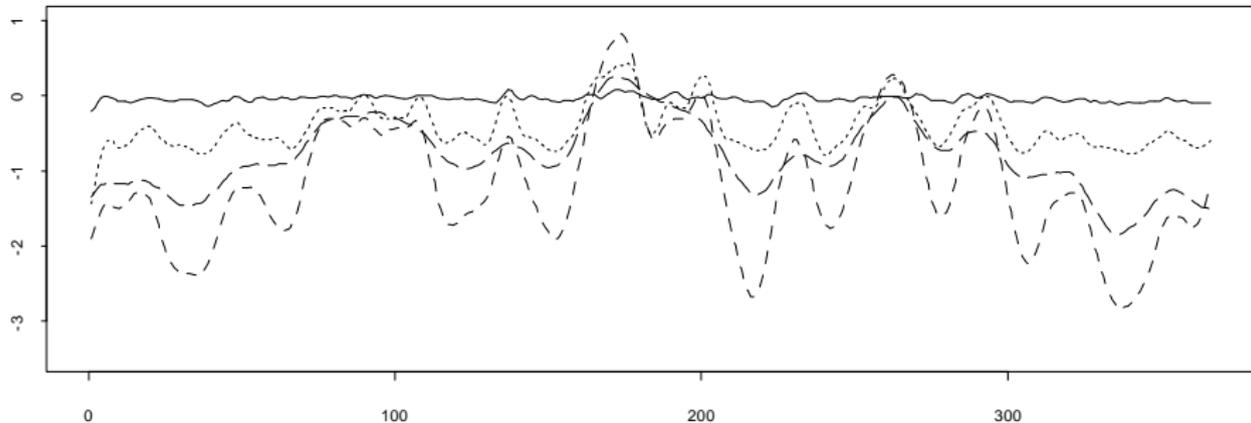


Figure: Average trajectories over 100 chains for $L = 1, 5, 20$ and 40 from top to bottom.

After 50 Iterations

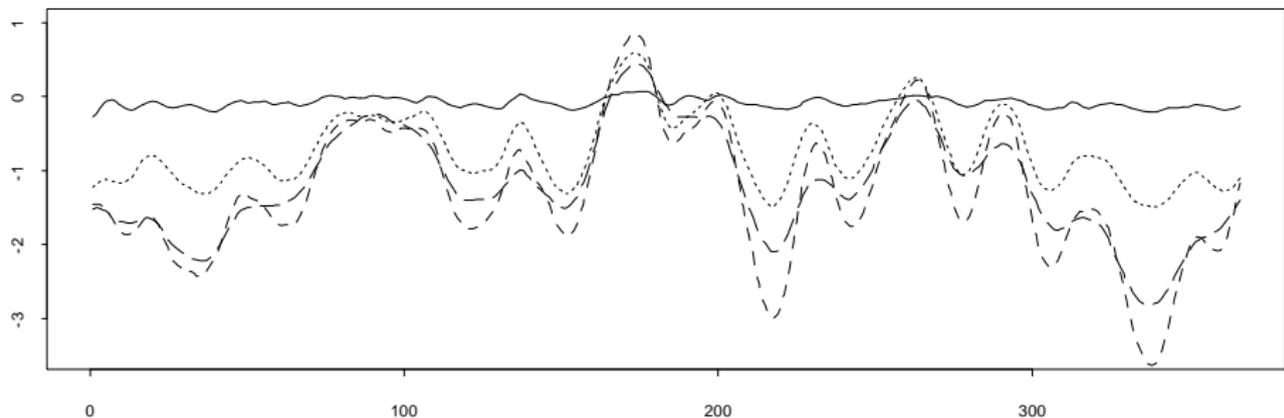


Figure: Average trajectories over 100 chains for $L = 1, 5, 20$ and 40 from top to bottom.

After 100 Iterations

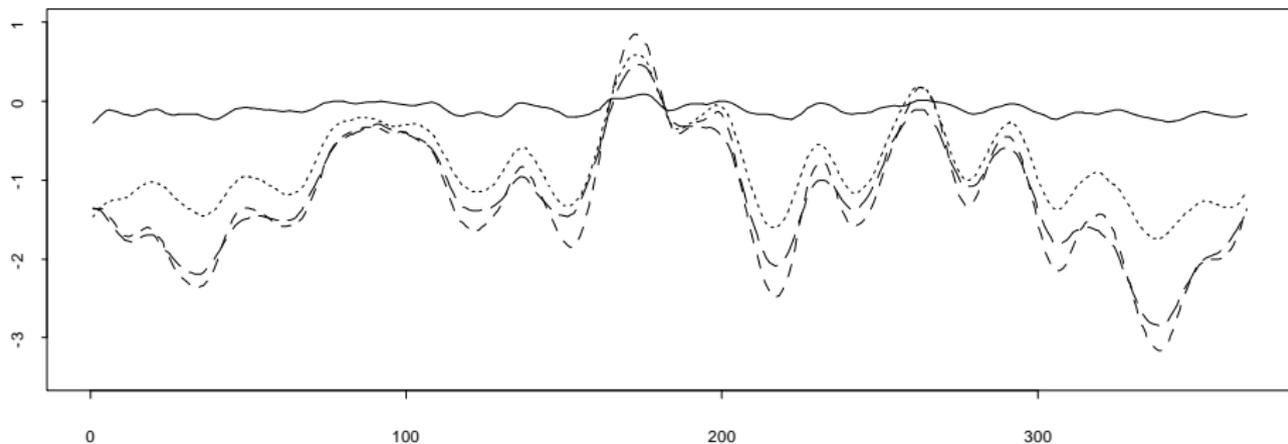


Figure: Average trajectories over 100 chains for $L = 1, 5, 20$ and 40 from top to bottom.

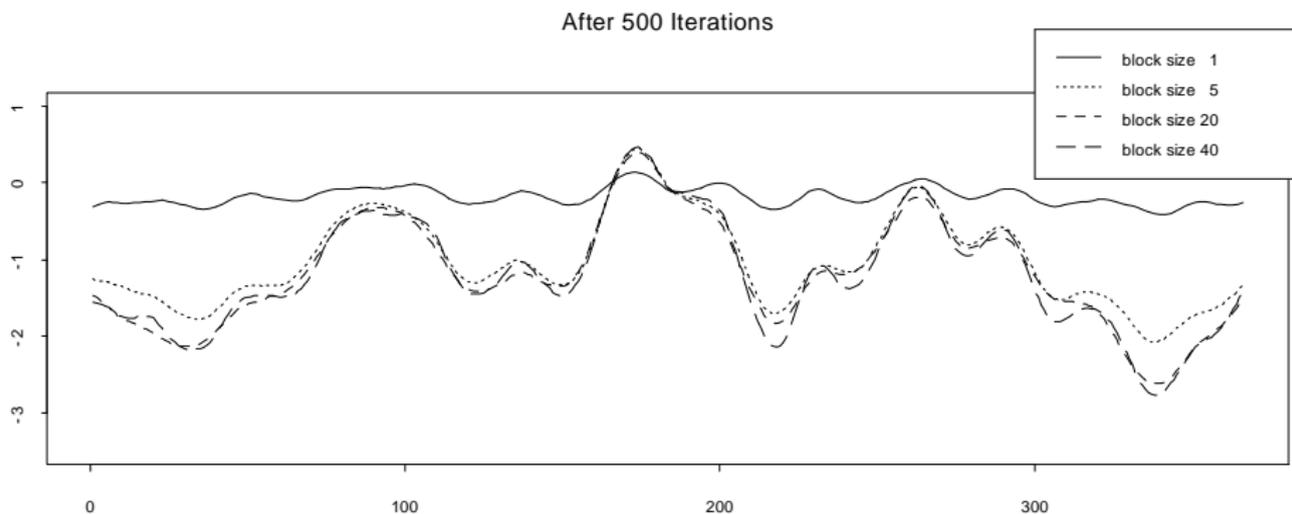


Figure: Average trajectories over 100 chains for $L = 1, 5, 20$ and 40 from top to bottom.

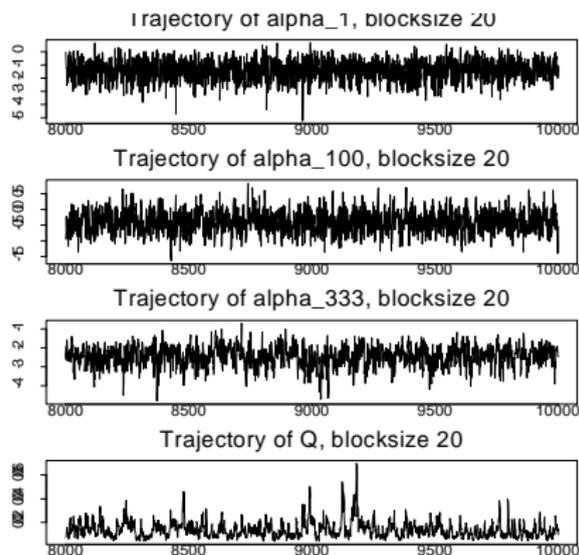
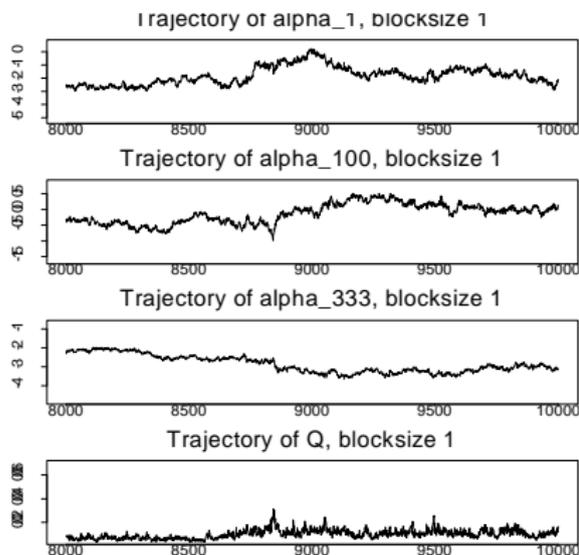


Figure: Traces of α_1 , α_{100} , α_{333} and σ^2 for $L = 1$ (left) and $L = 20$ (right).

- This (naive!) block sampling strategy performs well here because the likelihood of the observations is fairly flat.

- This (naive!) block sampling strategy performs well here because the likelihood of the observations is fairly flat.
- For a linear Gaussian observation equation, Knorr-Held compares this strategy to a direct Gibbs sampling implementation. As expected, the conditional proposal strategy is competitive when the observations are not very informative compared to the prior.

- This (naive!) block sampling strategy performs well here because the likelihood of the observations is fairly flat.
- For a linear Gaussian observation equation, Knorr-Held compares this strategy to a direct Gibbs sampling implementation. As expected, the conditional proposal strategy is competitive when the observations are not very informative compared to the prior.
- For more complex problems, such strategies are inefficient and we will need to use the observations to build the proposal.

- (Pitt & Shephard, 1999) propose a more efficient strategy... also more computationally intensive.

- (Pitt & Shephard, 1999) propose a more efficient strategy... also more computationally intensive.
- Consider the log full conditional distribution

$$\begin{aligned}
 & \log \pi (x_{k:k+L} | y_{k:k+L}, x_{k-1}, x_{k+L+1}) \\
 &= \sum_{i=k}^{k+L} \log g (y_i | x_i) + \sum_{i=k}^{k+L+1} \log f (x_{i+1} | x_i) \\
 &\equiv \sum_{i=k}^{k+L} \log g (y_i | x_i) - \frac{1}{2} \sum_{i=k}^{k+L+1} (x_{i+1} - Ax_i)^T \Sigma^{-1} (x_{i+1} - Ax_i)
 \end{aligned}$$

which is not quadratic in x_i hence $\pi (x_{k:k+L} | y_{k:k+L}, x_{k-1}, x_{k+1})$ is not Gaussian.

- (Pitt & Shephard, 1999) propose a more efficient strategy... also more computationally intensive.
- Consider the log full conditional distribution

$$\begin{aligned}
 & \log \pi (x_{k:k+L} | y_{k:k+L}, x_{k-1}, x_{k+L+1}) \\
 &= \sum_{i=k}^{k+L} \log g (y_i | x_i) + \sum_{i=k}^{k+L+1} \log f (x_{i+1} | x_i) \\
 &\equiv \sum_{i=k}^{k+L} \log g (y_i | x_i) - \frac{1}{2} \sum_{i=k}^{k+L+1} (x_{i+1} - Ax_i)^\top \Sigma^{-1} (x_{i+1} - Ax_i)
 \end{aligned}$$

which is not quadratic in x_i hence $\pi (x_{k:k+L} | y_{k:k+L}, x_{k-1}, x_{k+1})$ is not Gaussian.

- The idea is to expand the log-likelihood part around some point estimates

$$\begin{aligned}
 \log g (y_i | x_i) &\simeq \log g (y_i | \hat{x}_i) + \nabla \log g (y_i | \hat{x}_i) \cdot (x_i - \hat{x}_i) \\
 &\quad + \frac{1}{2} (x_i - \hat{x}_i)^\top \nabla^2 \log g (y_i | \hat{x}_i) (x_i - \hat{x}_i)
 \end{aligned}$$

- By doing this, we have a Gaussian approximation of the log-likelihood and then we obtain a Gaussian proposal

$$q(x_{1:n}, x'_{k:k+L}) = q(x_{-(k:k+L)}, x'_{k:k+L})$$

$$\begin{aligned} \log q(x_{-(k:k+L)}, x'_{k:k+L}) &\equiv \sum_{i=k}^{k+L} \nabla \log g(y_i | \hat{x}_i) \cdot (x_i - \hat{x}_i) \\ &+ \frac{1}{2} (x_i - \hat{x}_i)^\top \nabla^2 \log g(y_i | \hat{x}_i) (x_i - \hat{x}_i) \\ &- \frac{1}{2} \sum_{i=k}^{k+L+1} (x_{i+1} - Ax_i)^\top \Sigma^{-1} (x_{i+1} - Ax_i) \end{aligned}$$

- By doing this, we have a Gaussian approximation of the log-likelihood and then we obtain a Gaussian proposal

$$q(x_{1:n}, x'_{k:k+L}) = q(x_{-(k:k+L)}, x'_{k:k+L})$$

$$\begin{aligned} \log q(x_{-(k:k+L)}, x'_{k:k+L}) &\equiv \sum_{i=k}^{k+L} \nabla \log g(y_i | \hat{x}_i) \cdot (x_i - \hat{x}_i) \\ &+ \frac{1}{2} (x_i - \hat{x}_i)^\top \nabla^2 \log g(y_i | \hat{x}_i) (x_i - \hat{x}_i) \\ &- \frac{1}{2} \sum_{i=k}^{k+L+1} (x_{i+1} - Ax_i)^\top \Sigma^{-1} (x_{i+1} - Ax_i) \end{aligned}$$

- (Pitt & Shepard, 1999) propose to select

$$\hat{x}_{k:k+1} = \arg \max \pi(x_{k:k+L} | y_{k:k+L}, x_{k-1}, x_{k+L+1})$$

and a scheme to sample from $q(x_{-(k:k+L)}, x'_{k:k+L})$ which is of complexity $O(L)$.

- This algorithm is applied to the SV model where

$$X_k = \phi X_{k-1} + \sigma V_k, \quad V_k \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, 1)$$

$$Y_k = \beta \exp(X_k/2) W_k, \quad W_k \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, 1).$$

- This algorithm is applied to the SV model where

$$X_k = \phi X_{k-1} + \sigma V_k, \quad V_k \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, 1)$$

$$Y_k = \beta \exp(X_k/2) W_k, \quad W_k \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, 1).$$

- Prior are set to $\phi \sim \mathcal{U}[-1, 1]$, $\sigma^2 \sim \mathcal{IG}(\frac{\nu_\sigma}{2}, \frac{\gamma_\sigma}{2})$ and $\beta \sim \mathcal{IG}(\frac{\nu_\beta}{2}, \frac{\gamma_\beta}{2})$.

- This algorithm is applied to the SV model where

$$X_k = \phi X_{k-1} + \sigma V_k, \quad V_k \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, 1)$$

$$Y_k = \beta \exp(X_k/2) W_k, \quad W_k \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, 1).$$

- Prior are set to $\phi \sim \mathcal{U}[-1, 1]$, $\sigma^2 \sim \mathcal{IG}(\frac{\nu_\sigma}{2}, \frac{\gamma_\sigma}{2})$ and $\beta \sim \mathcal{IG}(\frac{\nu_\beta}{2}, \frac{\gamma_\beta}{2})$.
- Full conditional distributions of the parameters given $x_{1:n}, y_{1:n}$ are standard.

- This algorithm is applied to the SV model where

$$X_k = \phi X_{k-1} + \sigma V_k, \quad V_k \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, 1)$$

$$Y_k = \beta \exp(X_k/2) W_k, \quad W_k \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, 1).$$

- Prior are set to $\phi \sim \mathcal{U}[-1, 1]$, $\sigma^2 \sim \mathcal{IG}(\frac{\nu_\sigma}{2}, \frac{\gamma_\sigma}{2})$ and $\beta \sim \mathcal{IG}(\frac{\nu_\beta}{2}, \frac{\gamma_\beta}{2})$.
- Full conditional distributions of the parameters given $x_{1:n}, y_{1:n}$ are standard.
- Compared to standard single move strategies, the authors report significant improvement.

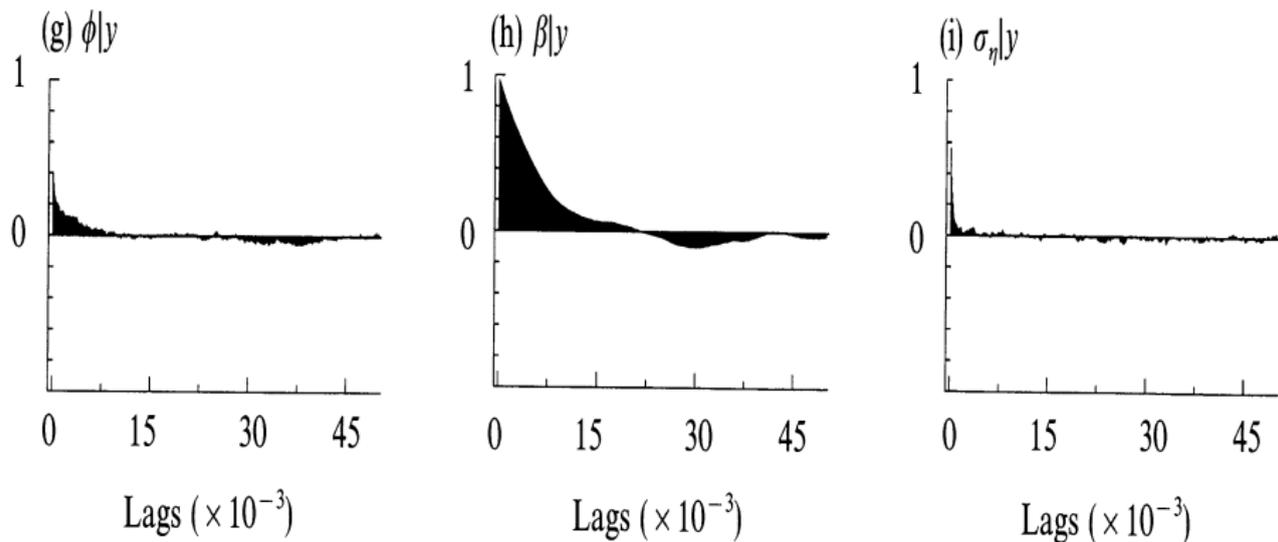


Figure: Autocorrelation plots for (ϕ, σ^2, β) with $L = 1$

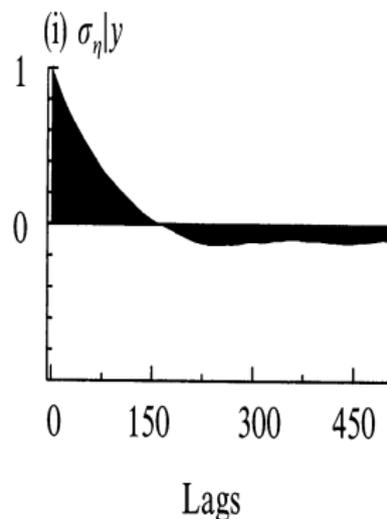
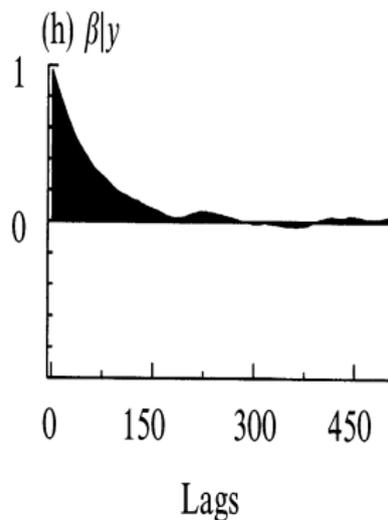
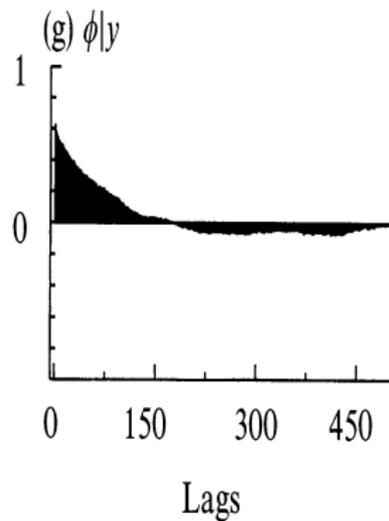


Figure: Autocorrelation plots for (ϕ, σ^2, β) with $L = 50$ on average