

# Group 10 – HA2 Report

## SF2955: Computer Intensive Methods in Mathematical Statistics

David Haapanen

Noah Eriksson Tewolde Berhane

May 21, 2025

### 1 Problem 1: Bayesian Analysis of Coal Mine Disasters

#### 1.1 a) Marginal Posterior Computations

We seek to compute the marginal posterior distributions of  $f(\theta \mid \boldsymbol{\lambda}, \mathbf{t}, \boldsymbol{\tau})$ ,  $f(\boldsymbol{\lambda} \mid \theta, \mathbf{t}, \boldsymbol{\tau})$ , and  $f(\mathbf{t} \mid \theta, \boldsymbol{\lambda}, \boldsymbol{\tau})$  up to an normalizing constant. Here,  $\boldsymbol{\tau} = (\tau_1, \dots, \tau_n)$  denotes the observed time points of the  $n = 191$  disasters,  $\mathbf{t} = (t_1, \dots, t_{d+1})$  contains the fixed end points coupled with all the breakpoints. Furthermore  $\boldsymbol{\lambda} = (\lambda_1, \dots, \lambda_n)$  contains the intensities  $\lambda_i$  for the interval  $[t_i, t_{i+1})$

##### Marginal posterior of $\theta$

Utilizing *Bayes' Theorem* one can rewrite the marginal posterior as proportional to the likelihood times the prior,

$$f(\theta \mid \boldsymbol{\lambda}, \mathbf{t}, \boldsymbol{\tau}) \propto f(\boldsymbol{\tau} \mid \mathbf{t}, \boldsymbol{\lambda}, \theta) \cdot f(\mathbf{t}, \boldsymbol{\lambda}, \theta)$$

where the likelihood can be simplified as  $f(\boldsymbol{\tau} \mid \mathbf{t}, \boldsymbol{\lambda}, \theta) = f(\boldsymbol{\tau} \mid \mathbf{t}, \boldsymbol{\lambda})$ . Furthermore, the breakpoints  $\mathbf{t}$  are independent of the intensities  $\boldsymbol{\lambda}$  and  $\theta$  such that

$$f(\mathbf{t}, \boldsymbol{\lambda}, \theta) = f(\mathbf{t}) \cdot f(\boldsymbol{\lambda}, \theta) = f(\mathbf{t}) \cdot f(\boldsymbol{\lambda} \mid \theta) \cdot f(\theta).$$

Thus we can consider the expression of

$$f(\theta \mid \boldsymbol{\lambda}, \mathbf{t}, \boldsymbol{\tau}) \propto f(\boldsymbol{\tau} \mid \mathbf{t}, \boldsymbol{\lambda}) \cdot f(\mathbf{t}) \cdot f(\boldsymbol{\lambda} \mid \theta) \cdot f(\theta)$$

Here only  $f(\boldsymbol{\lambda} \mid \theta)$  and  $f(\theta)$  are the terms that will depend on  $\theta$  since the prior,  $f(\boldsymbol{\tau} \mid \mathbf{t}, \boldsymbol{\lambda})$ , for given values of  $\boldsymbol{\lambda}$  and  $\mathbf{t}$  doesn't vary w.r.t.  $\theta$ . Thus discarding these terms we have that

$$f(\theta \mid \boldsymbol{\lambda}, \mathbf{t}, \boldsymbol{\tau}) \propto f(\boldsymbol{\lambda} \mid \theta) \cdot f(\theta)$$

Knowing that the prior for each intensity is  $\lambda_i \mid \theta \sim \Gamma(2, \theta)$  and  $\theta \sim \Gamma(2, \vartheta)$ , whilst all  $\lambda_i$  are independent we finally have that

$$f(\theta \mid \boldsymbol{\lambda}, \mathbf{t}, \boldsymbol{\tau}) \propto f(\boldsymbol{\lambda} \mid \theta) \cdot f(\theta) \propto \theta^{(2d+2)-1} \exp \left\{ -\theta \cdot \left( \vartheta + \sum_{i=1}^d \lambda_i \right) \right\}$$

which upon inspection is a  $\Gamma(2d + 2, \vartheta + \sum_{i=1}^d \lambda_i)$  distribution.

##### Marginal posterior of $\boldsymbol{\lambda}$

As described from before the marginal posterior can be rewritten with the likelihood times the prior using *Bayes' Theorem*, furthermore the simplifications mentioned from the previous section still holds. Thus, it holds that

$$f(\boldsymbol{\lambda} \mid \theta, \mathbf{t}, \boldsymbol{\tau}) \propto f(\boldsymbol{\tau} \mid \mathbf{t}, \boldsymbol{\lambda}) \cdot f(\mathbf{t}) \cdot f(\boldsymbol{\lambda} \mid \theta) \cdot f(\theta)$$

where only  $f(\boldsymbol{\tau} \mid \mathbf{t}, \boldsymbol{\lambda})$  and  $f(\boldsymbol{\lambda} \mid \theta)$  will depend on  $\boldsymbol{\lambda}$ . Hence we can conclude that

$$f(\boldsymbol{\lambda} \mid \theta, \mathbf{t}, \boldsymbol{\tau}) \propto f(\boldsymbol{\tau} \mid \mathbf{t}, \boldsymbol{\lambda}) \cdot f(\boldsymbol{\lambda} \mid \theta)$$

Each  $\lambda_i$  is independent from each other and follow with a  $\sim \Gamma(2, \theta)$  distribution. This coupled with that  $f(\boldsymbol{\tau} \mid \boldsymbol{\lambda}, \mathbf{t}) \propto \prod_{i=1}^d \lambda_i^{n_i(\boldsymbol{\tau})} \exp(-\lambda_i(t_{i+1} - t_i))$  gives the final expression of

$$f(\boldsymbol{\lambda} \mid \theta, \mathbf{t}, \boldsymbol{\tau}) \propto f(\boldsymbol{\tau} \mid \mathbf{t}, \boldsymbol{\lambda}) \cdot f(\boldsymbol{\lambda} \mid \theta) \propto \prod_{i=1}^d \lambda_i^{n_i(\boldsymbol{\tau})+1} \exp(-\lambda_i(\theta + t_{i+1} - t_i))$$

The independence between the  $\lambda_i$  terms further yields that  $\lambda_i \sim \Gamma(n_i(\boldsymbol{\tau}) + 2, \theta + t_{i+1} - t_i)$

### Marginal posterior of $\mathbf{t}$

Following the same steps from before we know that

$$f(\mathbf{t} \mid \theta, \boldsymbol{\lambda}, \boldsymbol{\tau}) \propto f(\boldsymbol{\tau} \mid \mathbf{t}, \boldsymbol{\lambda}) \cdot f(\mathbf{t}) \cdot f(\boldsymbol{\lambda} \mid \theta) \cdot f(\theta)$$

where only the terms  $f(\boldsymbol{\tau} \mid \mathbf{t}, \boldsymbol{\lambda})$  and  $f(\mathbf{t})$  depend on  $\mathbf{t}$  such that

$$\begin{aligned} f(\mathbf{t} \mid \theta, \boldsymbol{\lambda}, \boldsymbol{\tau}) &\propto f(\boldsymbol{\tau} \mid \mathbf{t}, \boldsymbol{\lambda}) \cdot f(\mathbf{t}) \propto \left( \prod_{i=1}^d \lambda_i^{n_i(\boldsymbol{\tau})} \exp(-\lambda_i(t_{i+1} - t_i)) \right) \left( \prod_{i=1}^d (t_{i+1} - t_i) \right) \\ &\propto \{\lambda_i \text{ are considered as constans}\} \propto \exp\left(-\sum_{i=1}^d \lambda_i(t_{i+1} - t_i)\right) \cdot \prod_{i=1}^d (t_{i+1} - t_i) \end{aligned}$$

for  $t_1 < t_2 < \dots < t_d < t_{d+1}$ . This has however no standard or clear distribution.

## 1.2 b) Hybrid MCMC Algorithm

Gibbs sampling (GS) was implemented combined with a Metropolis-Hastings (MH) step with the chosen proposal distribution of a random walk (provided in the assignment). A pseudo-code for the full algorithm, with  $\vartheta = 1$ ,  $\rho = 1$  and  $d - 1$  breakpoints, can be seen below

---

### Algorithm 1 Hybrid MCMC

---

```

Parameters:  $N = 10^5$ ,  $\vartheta = 1$ ,  $\rho = 1$ 
for  $i = 1$  to  $N - 1$  do
    Draw  $\theta^{i+1} \sim f(\theta \mid \boldsymbol{\lambda}^i, \mathbf{t}^i, \boldsymbol{\tau})$ 
    Draw  $\boldsymbol{\lambda}^{i+1} \sim f(\boldsymbol{\lambda} \mid \theta^{i+1}, \mathbf{t}^i, \boldsymbol{\tau})$ 
    for  $k = 2$  to  $d$  do
        Draw  $\varepsilon \sim \text{Unif}(-R, R)$  with  $R = \rho(t_{k-1}^i - t_{k+1}^i)$ 
        Propose  $t_k^* = t_k^i + \varepsilon$ 
        if  $t_{k-1} < t_k^* < t_{k+1}$  then
            Compute  $\alpha = \min\left(1, \frac{f(t_k^*)}{f(t_k)} = \frac{(t_k^* - t_{k-1})(t_{k+1} - t_k^*)}{(t_k - t_{k-1})(t_{k+1} - t_k)}\right)$ 
            if  $\text{Unif}(0,1) \leq \alpha$  then
                 $t_k^{i+1} \leftarrow t_k^*$ 
            else
                 $t_k^{i+1} \leftarrow t_k^i$ 
            end if
        end if
    end for
end for

```

---

Histograms of the distributions from sampled  $\theta$ ,  $\lambda$  and  $t$  were furthermore obtained when letting the number of breakpoints be 1, 2, 3 and 4 which can be seen in the different rows in fig. 1.

### 1.3 c) Behavior of MCMC with varying amount of breakpoints

The distribution of  $\theta$ , illustrated in fig. 1, doesn't change for the varying amount of breakpoints having mean  $\approx 1.33$ . The distributions of  $t$  are also similar in their shape but have different means for the varying amount of break points. This is not strange since one introduces new equidistant break points for the different cases. Furthermore the varying histograms of the  $\lambda$  distributions, adding more breakpoints produces increasingly overlapping intensity distributions. Overlapping peaks for the different intensities indicate that they may have the same intensity hence introducing extra break points may fail to capture any genuinely distinct shifts in disaster rates, and are thus redundant. This is clearly seen in fig. 1k as  $\lambda_3$  and  $\lambda_4$  have distributions that overlap. On the other hand having at least one break point can be strongly argued for as the peaks between the distribution of  $\lambda_1$  and  $\lambda_2$  are very distinct as illustrated in fig. 1b.

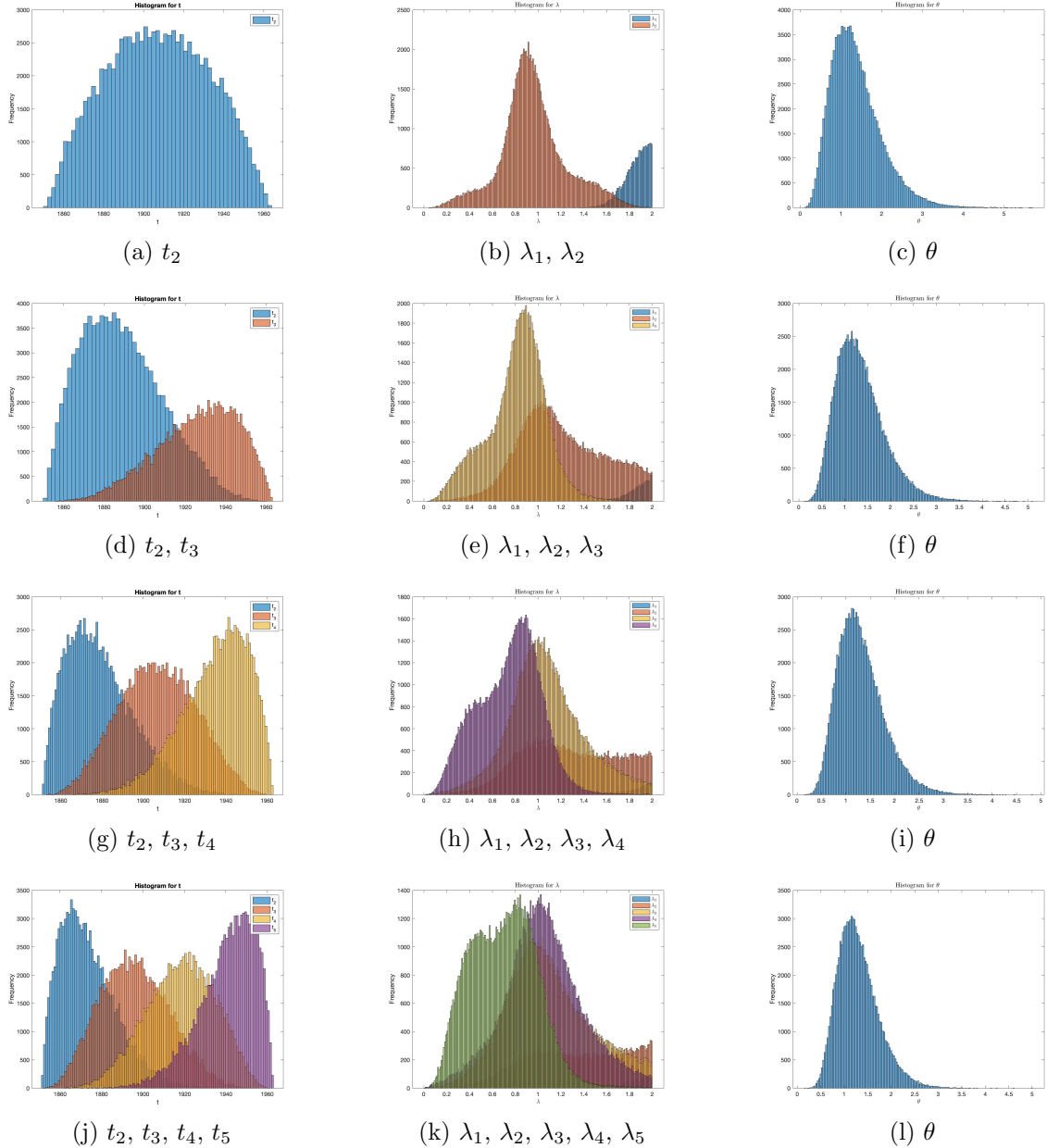


Figure 1: Distributions of  $\theta$ ,  $\lambda$ , and  $t$  using 1, 2, 3 and 4 breakpoints, with  $\vartheta = 1$  and  $\rho = 1$ .

#### 1.4 d) Sensitivity to choice of $\vartheta$

Histograms of  $\mathbf{t}$ ,  $\boldsymbol{\lambda}$  and  $\theta$  for  $\rho = 1$ ,  $\vartheta = 1$ , and 2 breakpoints can be seen as the second row of fig. 1. To investigate the sensitivity to the hyperparameter  $\vartheta$ , the MCMC algorithm was applied with  $\vartheta = 10$  and  $\vartheta = 100$ . The resulting histograms are displayed in fig. 2. When comparing these to the corresponding baseline histograms mentioned in the begining of this section, we observe minor differences in the posterior distributions of  $\mathbf{t}$  and  $\boldsymbol{\lambda}$  whilst the variation in distribution of  $\theta$  changes drastically. Since  $\theta \sim \Gamma(2, \vartheta)$  the change in distribution of  $\theta$  is expected. The result indicate that the distributions of  $\mathbf{t}$  and  $\boldsymbol{\lambda}$  are non-sensitive to changes in  $\theta$ , and therefore  $\vartheta$ . This further implies that the prior distribution  $f(\theta)$  doesn't affect the distributions of  $\mathbf{t}$  and  $\boldsymbol{\lambda}$ .

#### 1.5 e) Sensitivity to choice of $\rho$

Similarly from the previous question the MCMC algorithm was applied for  $\vartheta = 1$  and 2 breakpoints when letting  $\rho = 2$  and  $\rho = 3$  to produce the histograms seen in fig. 3. In comparison to the second row of fig. 1 one can see the magnitude of  $t_2$  and  $t_3$  varies for the different cases while the prior distribution of  $\boldsymbol{\lambda}$  and  $\theta$  is unchanged. This suggests that only the prior distribution of  $\mathbf{t}$  is sensitive to the value of  $\rho$ . Furthermore the mixing of the MH step was observed by editing the algorithm in section 1.2 such that the acceptance rate in the MH step could be evaluated. When varying  $\rho$  one could see that the acceptance rate would also change. An indication for good mixing is when the acceptance rate is  $\approx 30\%$  which could be continuously observed for when  $\rho = 1.25$  (even for the varying amount of breakpoints). Thus the mixing of the MH step is sensitive to  $\rho$  which aligns with the sensitivity of the prior distribution of  $\mathbf{t}$ . That is the draws are made using MH sampling hence candidates are either more likely to be rejected or accepted dependent on  $\rho$ , which correspondingly will affect the histograms of  $\mathbf{t}$ .

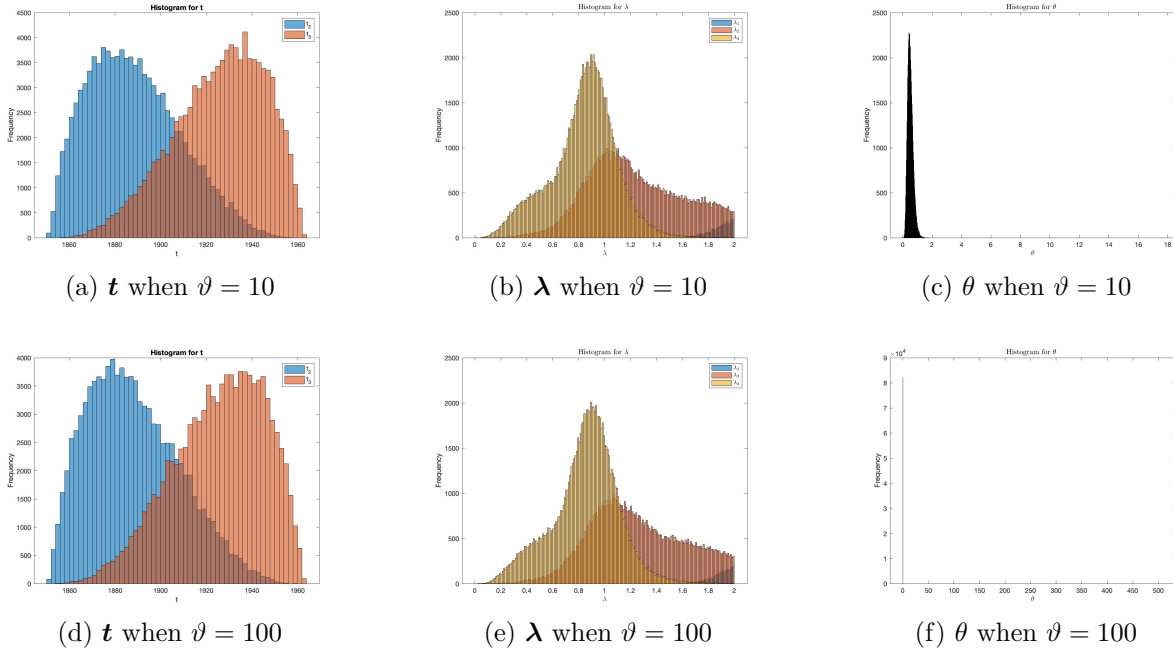


Figure 2: Distributions of  $\mathbf{t}$ ,  $\boldsymbol{\lambda}$ , and  $\theta$  for  $\rho = 1$ , two breakpoints, and varying  $\vartheta$ .

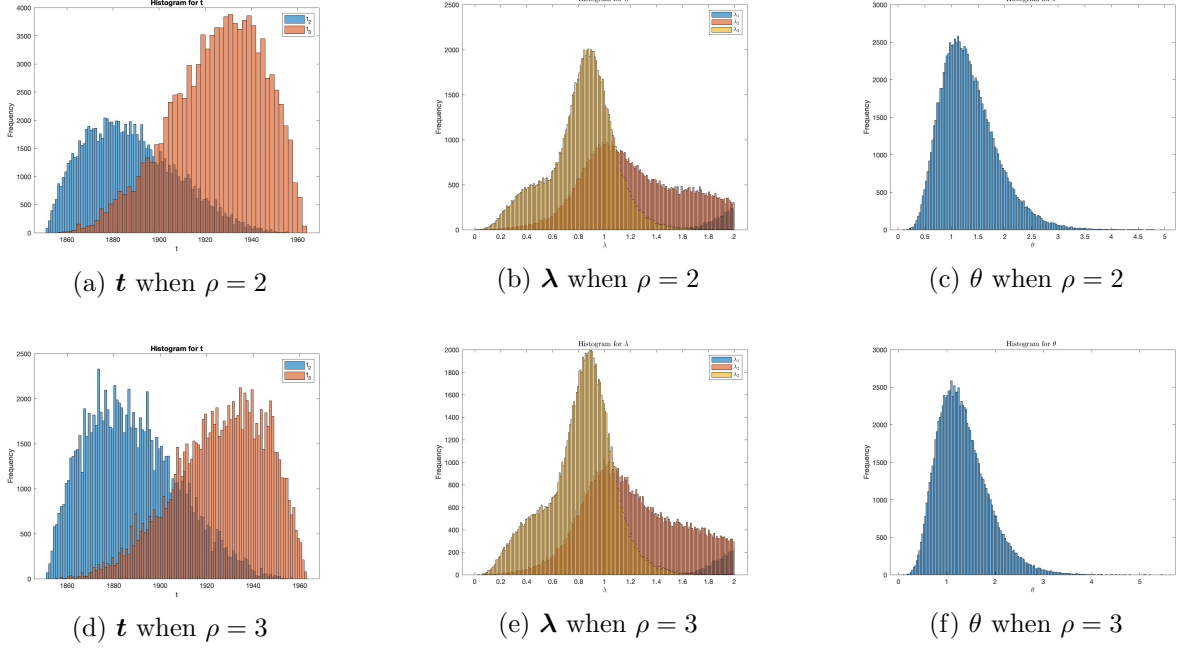


Figure 3: Distributions of  $\mathbf{t}$ ,  $\lambda$ , and  $\theta$  for  $\vartheta = 1$  and 2 breakpoints under varying  $\rho$ .

## 2 Problem 2: Sampling from a circle-shaped Posterior using Hamiltonian Monte Carlo

### 2.1 a) Posterior $\ln f(\theta | \mathbf{y})$ and its gradient $\nabla_{\theta} \ln f(\theta | \mathbf{y})$

Using Baye's theorem we have that  $f(\theta | \mathbf{y}) \propto f(\mathbf{y} | \theta) \cdot f(\theta)$  which upon taking the natural logarithm results in the expression of

$$\ln f(\theta | \mathbf{y}) \propto \ln f(\mathbf{y} | \theta) + \ln f(\theta)$$

Furthermore using that each  $y_i | \theta \sim \mathcal{N}(\theta_1^2 + \theta_2^2 = \|\theta\|^2, \sigma^2)$  and that each  $y_i$  are i.i.d one can get the following expressions

$$f(\mathbf{y} | \theta) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y_i - \|\theta\|^2)^2}{2\sigma^2}\right) \Rightarrow \ln f(\mathbf{y} | \theta) \propto -\sum_{i=1}^n \frac{(y_i - \|\theta\|^2)^2}{2\sigma^2}$$

$$\frac{\partial}{\partial \theta_j} \ln f(\mathbf{y} | \theta) = 2 \sum_{i=1}^n \frac{(y_i - \|\theta\|^2)}{\sigma^2} \cdot \theta_j$$

It is also given that  $f(\theta) \sim \mathcal{N}(\mathbf{0}, \Sigma)$ , where using that  $\Sigma^{-1}$  is symmetric yields in the following expressions

$$f(\theta) = \frac{1}{2\pi|\Sigma|^{1/2}} \exp\left(-\frac{1}{2}\theta^\top \Sigma^{-1} \theta\right) \Rightarrow \ln f(\theta) \propto -\frac{1}{2}\theta^\top \Sigma^{-1} \theta$$

$$\frac{\partial}{\partial \theta} \ln f(\theta) \propto \frac{\partial}{\partial \theta} \left(-\frac{1}{2}\theta^\top \Sigma^{-1} \theta\right) \propto -\frac{2}{2} \cdot \Sigma^{-1} \theta \propto -\Sigma^{-1} \theta$$

Thus the resulting expression of  $\ln f(\theta | \mathbf{y})$  and its gradient becomes

$$\ln f(\theta | \mathbf{y}) \propto \ln f(\mathbf{y} | \theta) + \ln f(\theta) \propto -\sum_{i=1}^n \frac{(y_i - \|\theta\|^2)^2}{2\sigma^2} - \frac{1}{2}\theta^\top \Sigma^{-1} \theta$$

$$\nabla_{\theta} \ln f(\theta | \mathbf{y}) \propto \frac{2}{\sigma^2} \sum_{i=1}^n (y_i - \|\theta\|^2) \cdot \theta - \Sigma^{-1} \theta$$

## 2.2 b) Hamiltonian Monte Carlo Algorithm Design

Using the previously determined posterior of  $\ln f(\boldsymbol{\theta} \mid \mathbf{y})$  and  $\nabla_{\boldsymbol{\theta}} \ln f(\boldsymbol{\theta} \mid \mathbf{y})$  an algorithm for sampling from  $f(\boldsymbol{\theta} \mid \mathbf{y})$  was designed using the Hamiltonian Monte Carlo (HMC) method. Similar to the MH algorithm an auxiliary proposal is used for drawing candidate samples. The HMC proposal is based on Hamiltonian dynamics, involving the gradient of the desired distribution to propose more informed moves leading to faster mixing. The core of HMC is the *Leapfrog* integrator, which is a volume-preserving involution, yielding proposals that are accepted with high probability. The integrator solves Hamilton's equations numerically to propose new  $(\boldsymbol{\theta}, v)$  pairs that explore the target distribution efficiently. Given an initial momentum  $v$  a *Leapfrog* integrator utilized for time length  $T = \varepsilon \cdot L$ , using  $L$  iterations of time-step  $\varepsilon$  can be described as:

$$\begin{aligned} v_{m+1/2} &\leftarrow v_m - \frac{\varepsilon}{2} \nabla U(\boldsymbol{\theta}_m), \\ \boldsymbol{\theta}_{m+1} &\leftarrow \boldsymbol{\theta}_m + \varepsilon v_{m+1/2}, \\ v_{m+1} &\leftarrow v_{m+1/2} - \frac{\varepsilon}{2} \nabla U(\boldsymbol{\theta}_{m+1}), \end{aligned}$$

for  $m \in \{0, 1, \dots, L-1\}$  before applying the sign-shift operator  $\sigma(\boldsymbol{\theta}, v) = (\boldsymbol{\theta}, -v)$ . The newly computed  $\boldsymbol{\theta}$  is then either accepted or rejected in a Metropolis step. A pseudo code of the HMC algorithm can be seen below.

---

### Algorithm 2 Hamiltonian Monte Carlo (HMC)

---

**Parameters:**  $\varepsilon, L, N$ , observed data  $\mathbf{y}$ ,  $\sigma, \boldsymbol{\Sigma}$ , initial value  $\boldsymbol{\theta}^{(0)}$

**for**  $i = 2$  to  $N$  **do**

    Sample momentum:  $\mathbf{v}^{(i)} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_{2 \times 2})$

    Set  $\boldsymbol{\theta}_{\text{current}} \leftarrow \boldsymbol{\theta}^{(i)}$

    Perform leapfrog integration using step size  $\varepsilon$  and  $L$  iterations

$[\boldsymbol{\theta}^*, \mathbf{v}^*] \leftarrow \text{Leapfrog}(\boldsymbol{\theta}_{\text{current}}, \mathbf{v}^{(i)}, \varepsilon, L)$

    Compute:  $\alpha = \min(1, \exp(\ln f(\boldsymbol{\theta}^* \mid \mathbf{y}) - \ln f(\boldsymbol{\theta}_{\text{current}} \mid \mathbf{y}) - \frac{1}{2}\|\mathbf{v}^*\|^2 + \frac{1}{2}\|\mathbf{v}^{(i)}\|^2))$

    Sample  $U \sim \text{Unif}(0, 1)$

**if**  $U \leq \alpha$  **then**

$\boldsymbol{\theta}^{(i+1)} \leftarrow \boldsymbol{\theta}^*$

**else**

$\boldsymbol{\theta}^{(i+1)} \leftarrow \boldsymbol{\theta}^{(i)}$

**end if**

**end for**

**Output:** Samples  $\{\boldsymbol{\theta}^{(i)}\}_{i=1}^N$

---

## 2.3 c) Comparison between MH and HMC sampling

Using the outline of the HMC algorithm provided in section 2.2 with  $N = 100\,000$  the distribution of  $f(\boldsymbol{\theta} \mid \mathbf{y})$  was sampled using the given parameters  $\sigma$  and  $\boldsymbol{\Sigma}$  combined with the set values for  $\varepsilon$  and  $L$ . The values for these were set by using various combinations and assessing autocorrelation plots of  $\theta_1$  and  $\theta_2$  and acceptance rates of the HMC algorithm.

- Smaller values of  $L$  could lead to higher acceptance rates but continuously showed autocorrelation plots with slower decay. Since the autocorrelation plots showed more correlation between samples the higher acceptance rate was not seen as beneficial.
- Larger values of  $\varepsilon$  made the assessed distribution of  $f(\boldsymbol{\theta} \mid \mathbf{y})$  look less like the provided plot from the assignment. The autocorrelation functions would also display values that were further away from zero.
- Smaller values of  $\varepsilon$  would make the computational time of the algorithm longer but show

less correlation between samples in the autocorrelation plots.

Based on the observation the values were set to  $\varepsilon = 0.09$  and  $L = 25$ , where it was found that these values continuously led to a acceptance rate of  $\approx 77\%$  while the samples of  $\theta$  seem to have sufficiently low correlation between them. The sampled distribution of  $f(\theta | \mathbf{y})$  and the autocorrelation plots of  $\theta_1$  and  $\theta_2$  can be seen in fig. 4.

$N = 100\,000$  samples from the distribution of  $f(\theta | \mathbf{y})$  was also retrieved using a MH algorithm with a bivariate random walk with the candidate having the distribution of  $\theta^* \sim \mathcal{N}(\theta, \zeta^2 \mathbf{I}_{2 \times 2})$ . The tuning of  $\zeta$  was done by trying various values and assessing the acceptance rate of the algorithm.

- The acceptance rate was  $\approx 32\%$  around  $\zeta = 0.35$
- The autocorrelation plots showed very strong correlation between samples for the varied values of  $\zeta$ , but slightly lower around the chosen  $\zeta = 0.35$ .
- The sampled distribution of  $f(\theta | \mathbf{y})$  was generally very unstable and not smooth looking compared to the plot provided in the assignment.

Autocorrelation plots of  $\theta_1$  and  $\theta_2$  and the sampled distribution of  $f(\theta | \mathbf{y})$  using the chosen parameter value  $\zeta = 0.35$  can be seen in fig. 5

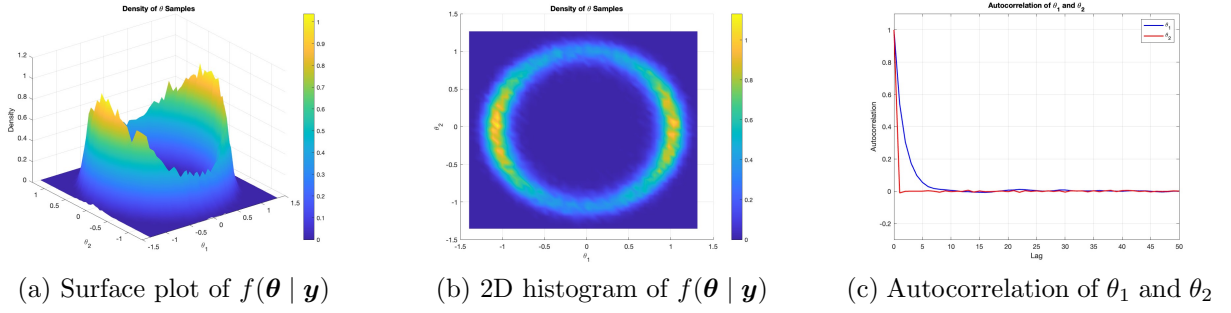


Figure 4: Distribution and autocorrelation plots using HMC algorithm with  $\varepsilon = 0.09$ ,  $L = 25$ .

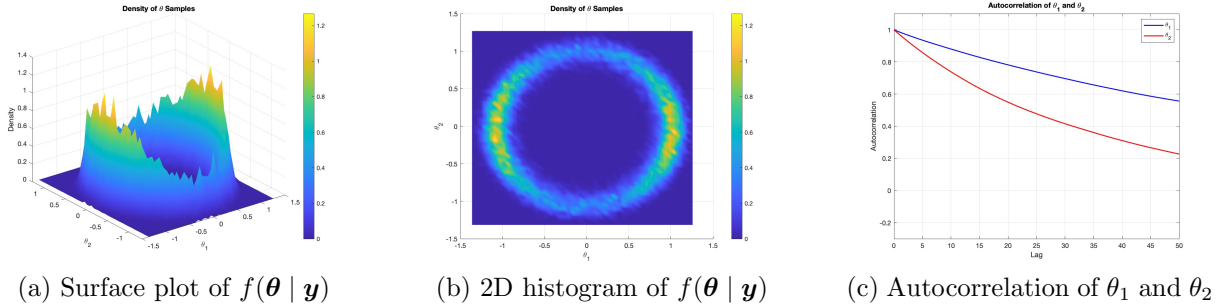


Figure 5: Distribution and autocorrelation plots using MH algorithm with  $\zeta = 0.35$ .

After drawing  $N = 100,000$  samples using the MH and HMC algorithm to estimate the distribution of  $f(\theta | \mathbf{y})$ , it is clear that HMC provides a more accurate estimate of the target distribution. The HMC algorithm samples seen in fig. 4b show a circular shape which has a more even edge while the samples from the MH are mostly drawn from the right or left focal points seen in fig. 5b. One can also see that the distribution in fig. 4a resembles the provided posterior distribution from the assignment more in comparison to fig. 5a. The difference is also supported by the autocorrelation plots. In fig. 4c, the autocorrelation decreases rapidly, indicating that the HMC samples are less correlated and mix more efficiently. In contrast, fig. 5c shows a much slower decay in autocorrelation, meaning the MH samples are more dependent.