# Computational Statistics Coursework 3

Lorenz Wolf (CID 01351308)

June 16, 2021

## Question 1

In this question we consider the following density on $\mathbb{R}^2$

$$f(x,y) = k \exp\left\{-\frac{x^2}{100} - \left(y + \frac{3}{100}x^2 - 3\right)^2\right\}$$

where $k$ is a normalising constant. First we will construct and implement an MCMC algorithm to sample from this distribution. We will then use this sampler to estimate the expected values of X and Y with joint distribution f(x,y). Finally, a sampling algorithm with adaptive proposal distribution as described by Haario et al. in [1] is implemented and then compared to the MCMC sampler.

### a)

To sample from f(x,y) we propose a Metropolis Hastings algorithm with a additive random walk specified by a two dimensional normal distribution as outlined in Algorithm 1. We define

- $\Sigma$, a valid variance-covariance matrix of a 2-dimensional random variable

- nsteps, the number of iterations

- f(x,y), the target density up to a normalising constant

- $\mathcal{N}_2(\mu, \Sigma)$, the 2-d normal distribution with mean $\mu$ and variance covariance matrix $\Sigma$

---

**Algorithm 1:** MCMC algorithm with target f(x,y)

---

**Inputs:** $X_0 = (x_0, y_0)$, $\Sigma$, nsteps, f

**for** $i=1\ldots nsteps$ **do**

    Propose $Y \sim \mathcal{N}_2(X_{i-1}, \Sigma)$

    Simulate $U \sim \text{Uniform}(0,1)$

    **if** $U \le \alpha = min\left(1, \frac{f(Y)}{f(X_{i-1})}\right)$ **then**

        $X_i = Y$

    **else**

        $X_i = X_{i-1}$

    **end**

**end**

---

In Figure 1 the unnormalised density is shown. It can be observed that X and Y have very different scales, i.e. the variance in X is much larger than in Y. Thus, running the algorithm for a simple standard proposal with $\Sigma$ such as

$$\Sigma = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$$

will lead to poor mixing of the chain.

Since the mode of the distribution can be analytically computed to be $p^* = (x,y) = (0,3)$, we will use this as starting values for all future runs of the algorithm. To improve the mixing we need to take the different scales and dependency of X and Y into account by choosing $\Sigma$ appropriately. In order to choose these values iteratively use Algorithm 1 (starting with $\Sigma = I_1$) to sample from the target distribution, then compute sample covariance, and standard deviations to obtain a better estimate for $\Sigma$. This new estimate is used to specify the proposal in the next iteration. Afterwards we try several values for $\Sigma$ around this estimate and obtain that the effective sample size is maximised for

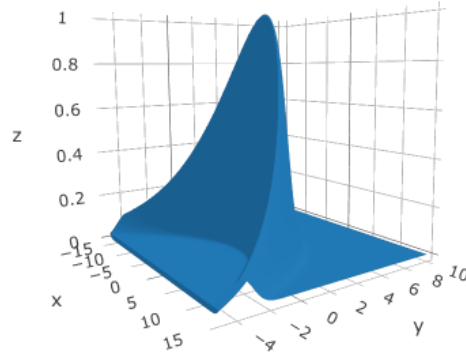$$\hat{\Sigma} = \begin{pmatrix} 7^2 & 0.22 \\ 0.22 & 2.1^2 \end{pmatrix}.$$

Figure 1: Joint density of X and Y.

To assess the convergence we run the algorithm for a sample size of $5 \times 10^4$ and $\Sigma = \hat{\Sigma}$ with 6 different starting values. The traceplots presented in Figures 2 and 3 show good convergence with occasional values of Y large and negative. Furthermore, we compute the Gelman Rubin diagnostic and obtain a point estimate and upper confident interval of (1,1.01) and (1,1.01) for X and Y respectively. Additionally, the gelman plots shown in Figures 4 and 5 do not indicate problems with convergence.
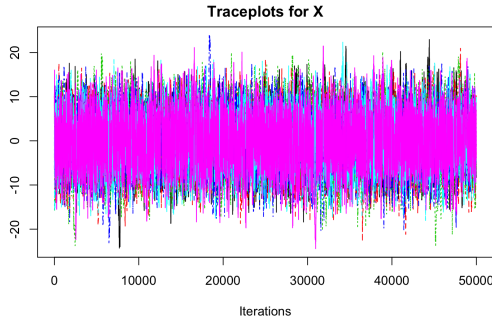


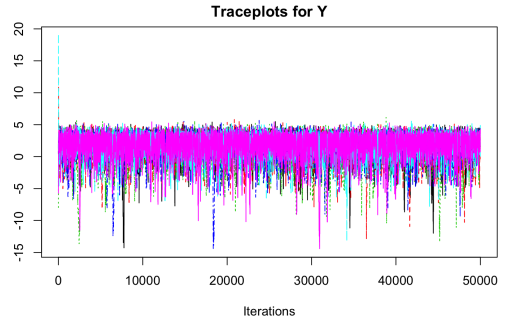Figure 2: Traceplots of X for chain length $5 \times 10^4$ with 6 different starting values



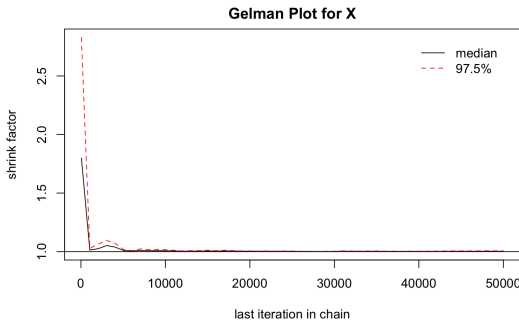Figure 3: Traceplots of Y for chain length $5 \times 10^4$ with 6 different starting values



Figure 4: Gelman plot for X for chain lengths $5 \times 10^4$ with 6 different starting values.
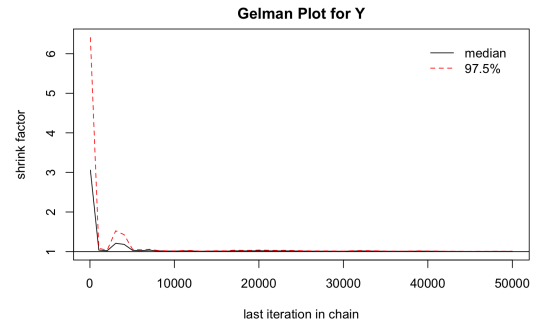


Figure 5: Gelman plot for Y for chain lengths $5 \times 10^4$ with 6 different starting values.

With $\Sigma = \hat{\Sigma}$ and starting value $p^*$ the algorithm is run to obtain a sample of size $10^5$. The first $10^4$ values of this sample are omitted as burn-in to eliminate the influence of the starting value. In Figures 6 and 7 the acf's are presented. They show a reasonable decrease indicating that mixing is ocurring. We note that there is still significant dependence in the chain. However, having tried several different inputs for $\Sigma$ no significant improvement has been found. Plotting the joint density in Figure 8 reveals the 'banana'-like shape required according to the plot of the real density in Figure 1. However, in the extreme regions the symmetry of f in x is not observed.
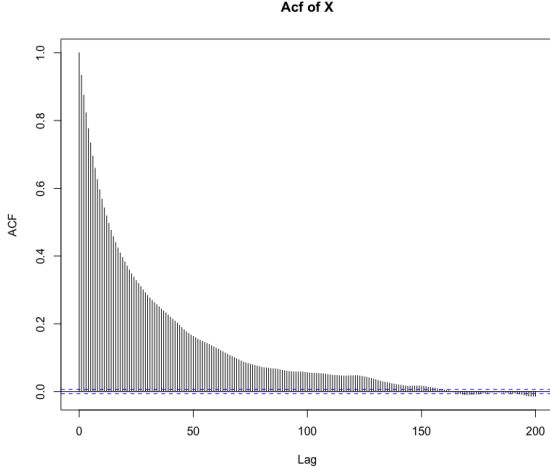


Figure 6: Acf of Y for sample size of $9 \times 10^4$ after burn-in with starting value $(0,3)$, $\Sigma = \hat{\Sigma}$ and burn-in of $10^4$.
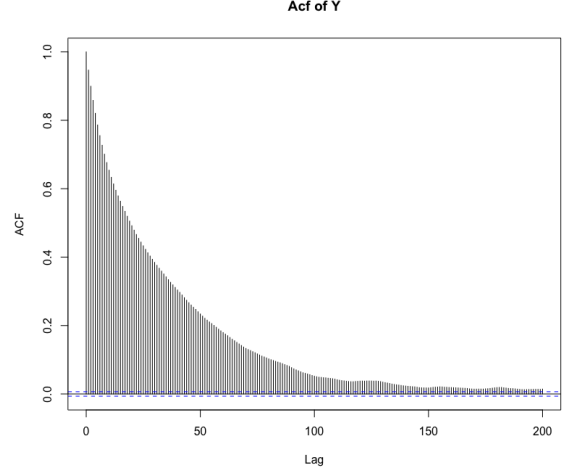
Figure 7: Acf of Y for sample size of $9 \times 10^4$ after burn-in with starting value $(0,3)$, $\Sigma = \hat{\Sigma}$ and burn-in of $10^4$.
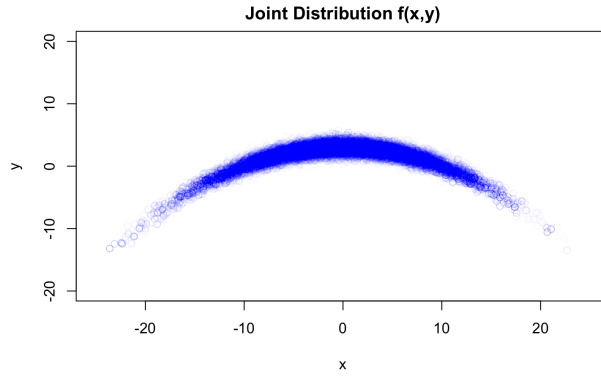


Figure 8: Joint density of X and Y for sample size of $9 \times 10^4$ after burn-in with starting value $(0,3)$, $\Sigma = \hat{\Sigma}$ and burn-in of $10^4$.

## b)

Using the sampler discussed in part a) with starting value $X_0 = (0,3)$ and $\Sigma = \hat{\Sigma}$ as specified in part a) we obtain a sample of size $9 \times 10^4$, after having discarded the first $10^4$ points as burn-in. In the following computations of the means we use this sample. Due to the relatively high dependency in the resulting chain we suggest to use thinning. However, since the acf decreases slowly an appropriate step size for thinning would need to be large. As this would leave us with only few observations (limitations due to time and computational power to sample a longer chain) we omit thinning at this point and note that for better estimates a longer chain

3

should obtained and thinned appropriately. The estimates obtained are

$$\hat{E}(X) = \frac{1}{9 \times 10^4} \sum_{t=1}^{9 \times 10^4} X_i \approx -0.2817$$

$$\hat{E}(Y) = \frac{1}{9 \times 10^4} \sum_{t=1}^{9 \times 10^4} Y_i \approx 1.5461$$

Using numerical integration the true values are computed to be $\mathbb{E}(X) = 0$ and $\mathbb{E}(Y) \approx 1.614$. Thus, the estimates (especially for X) are not very accurate with errors of 0.28 and 0.068 respectively.

## c)

In this section we implement a random walk Metropolis algorithm with an adaptive proposal distribution as suggested in [1]. The key difference to the sampler in part a) is that the proposed adaptive algorithm adjusts the proposal distribution based on the samples so far obtained.

In particular assume the points $X_1, \ldots, X_k$ have been already sampled. Then the proposal Y is sampled from the distribution $q_k(.|X_1, \ldots, X_k)$ which depends on the last H observations $X_{k-H+1}, \ldots, X_k$, where H is one of the input parameters to be specified. As proposal $q_k(.|X_1, \ldots, X_k)$ a 2 dimensional Gaussian distribution with mean $(X_k)$ and a covariance dependent on the sampled points is implemented. With at least H points sampled at time t define $K \in \mathbb{R}^{H \times d}$ to be the matrix with row i equal to $X_{t-H+i}$ and define $\tilde{K} = K - \mathbb{E}(K)$ the corresponding centered matrix. Here $\mathbb{E}(K)$ is estimated by the mean of the H values. Then the proposal is taken to be

$$q_t(.|X_1, \ldots, X_t) \sim X_t + \frac{c_d}{\sqrt{H-1}} \tilde{K}^T \mathcal{N}(0, I_H), \tag{1}$$

where $I_H$ is the H dimensional identity matrix and $c_d = \frac{2.4}{\sqrt{d}}$.

To initialise the algorithm the first H points are sampled with the MCMC sampler from part a). Afterwards the proposal is updated as specified above and then further updates are performed according to the update frequency U. For an outline of the algorithm refer to Algorithm 2.

---

**Algorithm 2:** Adaptive Proposal Metropolis algorithm with target f(x,y)

---

**Inputs:** $X_0 = (x_0, y_0)$ (starting value), $\Sigma$ (inital covariance), H (History length), U (Update frequency), nsteps (Number of iterations), f (Target up to proportionality)

**for** $i=1 \ldots H$ **do**
    Propose $Y \sim \mathcal{N}_2(X_{i-1}, \Sigma)$
    Simulate $V \sim \text{Uniform}(0, 1)$
    **if** $V \leq \alpha = min\left(1, \frac{f(Y)}{f(X_{i-1})}\right)$ **then**
        $X_i = Y$
    **else**
        $X_i = X_{i-1}$
    **end**
**end**
$K = (X_1, \ldots, X_H)^T$
$\tilde{K} = K - E(K)$
**for** $i=H+1 \ldots nsteps+H$ **do**
    Propose $Y \sim X_t + \frac{c_2}{\sqrt{H-1}} \tilde{K}^T \mathcal{N}(0, I_H)$
    Simulate $V \sim \text{Uniform}(0, 1)$
    **if** $V \leq \alpha = min\left(1, \frac{f(Y)}{f(X_{i-1})}\right)$ **then**
        $X_i = Y$
    **else**
        $X_i = X_{i-1}$
    **end**
    **if** $i = 0 \ mod \ U$ **then**
        $K = (X_{i-H+1}, \ldots, X_i)^T$
        $\tilde{K} = K - E(K)$
    **end**
**end**
**return** $(X_{H+1}, \ldots, X_{H+nsteps})$

---

To assess the convergence, the algorithm is run for 6 different starting values with U=H=200 (chosen according to [1]). The traceplots of the chains of length $5 \times 10^4$ for X and Y are shown in Figures 9 and 10 respectively and do not reveal any issues with convergence. Furthermore, the Gelman Rubin diagnostic and an upper confident interval are computed to be (1,1) for X and (1,1) for Y. Additionally, the Gelman plots are checked to look fine.
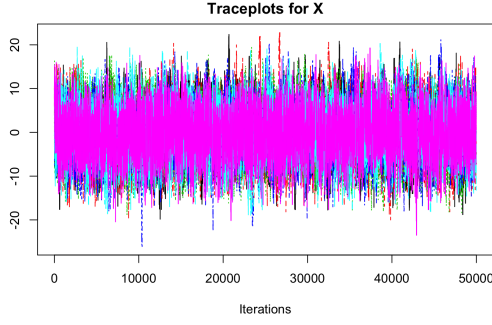


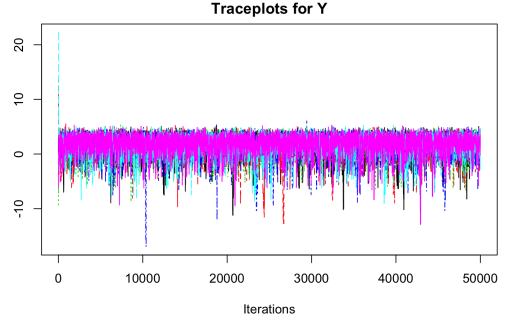Figure 9: Traceplots of X for chain length $3 \times 10^4$ with 6 different starting values



Figure 10: Traceplots of Y for chain length $3 \times 10^4$ with 6 different starting values

With U=H=200, starting value (0,3) and $\Sigma$ as above we obtain a sample of size $9 \times 10^4$ after a burn-in of $10^4$. The Acf's of X and Y for this sample are presented in Figures 11 and 12 respectively and indicate no problems. The joint density has the 'banana'-like shape as required.
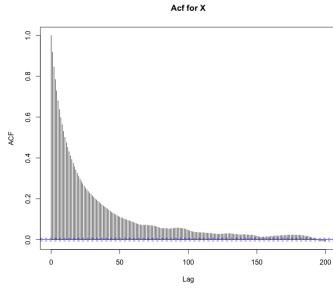


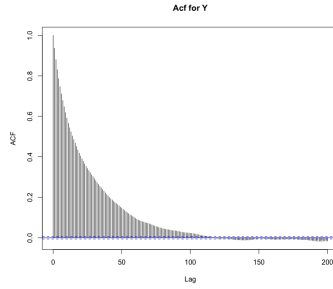Figure 11: Acf of X for adaptive proposal sampler.



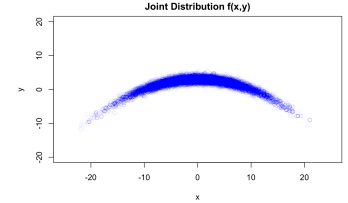Figure 12: Acf of Y for adaptive proposal sampler.



Figure 13: Joint density of X and Y

## Advantages compared to sampler a)

Haario et al describe the specification of a proposal as the potential bottleneck of the Metropolis algorithm when we have small knowledge of the target distribution. The proposal distribution influences the mixing of the chain and the convergence. Furthermore, it should be easy to sample from for computational efficiency. The adaptive proposal addresses the problem of choosing a proposal by automatically adjusting the proposal distribution based on the samples already obtained. The adaptive proposal proves to be especially strong when the target distribution is non-linear or curved and there is correlation between the components. In the case of correlated components Haario et al. have found the adaptive proposal to be more effective compared to a Metropolis algorithm with not well tuned proposal and not much less effective compared to a Metropolis algorithm with optimally selected proposal. Compared to the sampler in part a) with an adaptive proposal we do not have to be concerned with the choice of $\Sigma$ which can prove to be difficult. The implementation of the sampler in part c) shows increased computational speed and smaller upper confident intervals for the Gelman Rubin diagnostic. Furthermore, acf's for the adaptive proposal sampler shown in Figures 11 and 12 decrease quicker than those for the sampler in part a), which indicates better mixing and smaller dependence in the chain. Additionally we obtain increased effective sample sizes compared to the sampler in part a) as shown in Table 1. The resulting estimates of the means are also more accurate. Comparing the joint densities shown in Figures 8 and 13 it can be observed that the properties of the true joint density (e.g. symmetry) for extreme x-values are more accurate for the sample obtained with the adaptive proposal sampler. All these are desirable properties of a sampler which are obtained, without choosing a specific proposal. A disadvantage that is not observed in our case is that the adaptive proposal sampler is an approximate method which can lead to a bias

in the simulation [1].

For a more detailed comparison of the samplers the bias and standard errors as well as root mean squared errors of several estimates could be computed by repeated sampling with the different samplers.

|  |  | ESS | $\hat{\mathbb{E}}(.)$ | Time |
|---|---|---|---|---|
| Sampler a) | X | 1883.4 | -0.2817 | 1.34s |
|  | Y | 1409.6 | 1.5461 |  |
| Sampler b) | X | 2374.1 | 0.021 | 10.22s |
|  | Y | 1989.8 | 1.622 |  |

Table 1: Properties based on samples of size $10^5$ and starting values $X_0 = c(0,3)$ and $\Sigma = \hat{\Sigma}$. For the adaptive proposal we take U=H=200. A burn-in of $10^4$ points is chosen. Time reported is the average time of 10 iterations taken to sample a chain of length $5 \times 10^4$

# Question 2

In Question 2 we consider the density

$$g(\theta) = k \left( \frac{4}{10} \exp\left\{-(8-\theta)^2\right\} + \frac{6}{10} \exp\left\{-(30-\theta)^2\right\} \right) \propto \exp\left\{-H(\theta)\right\},$$

where k is the normalising constant. In Figure 16 we can observe the two separated modes of g($\theta$) (Tempered distribution with T=1).

## a)

First consider an additive random walk Metropolis Hastings algorithm with normally distributed noise (mean 0 and variance $\sigma^2$) to sample from g(.). To demonstrate that such an algorithm will not explore the entire state space we consider two different scenarios.

First, consider such an algorithm with small *sigma* (e.g. around 1). With such a $\sigma$, the acceptance rate will be good enough for the states to change frequently. However, the variance of the proposal is to small for the chain to jump between the different modes of g(.) (except in some very rare cases). Thus, depending on the starting value a chain converges to one of the modes and does not explore the entire state space. This scenario can be observed in Figure 14.

Secondly consider such an algorithm with large $\sigma$. Due to the large variance the proposals are more likely to be far away from the current state. Thus the chain might jump between the different modes. However, such proposals are less likely to be accepted, so the chain jumps fewer times. This causes very high dependence and slow to no convergence of the chain. This scenario can be observed in Figure 15.

We conclude that the entire state space can not be explored properly with a simple Metropolis Hastings algorithm.
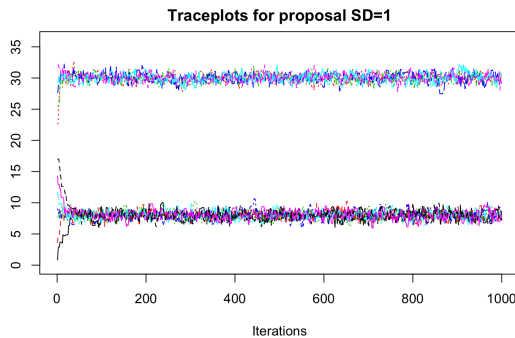


Figure 14: Traceplots for chains of length $10^3$ for several starting values with proposal standard deviation equal to 1.
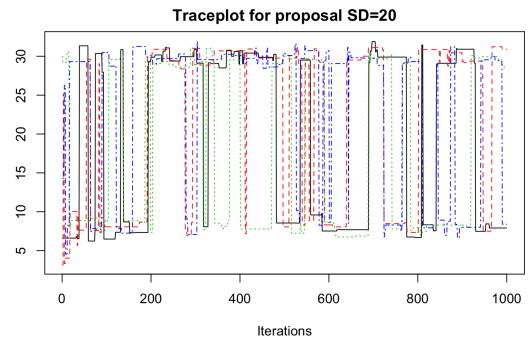
Figure 15: Traceplots for chains of length $10^3$ for several starting values with proposal standard deviation equal to 20.

6

**b)**

To sample from this density a parallel tempering algorithm is implemented. We define the tempered distributions with temperatures $T_i$, i=1,...,M, to be

$$g_{T_i}(\theta) = k \left( \frac{4}{10} \exp\left\{-(8-\theta)^2\right\} + \frac{6}{10} \exp\left\{-(30-\theta)^2\right\} \right)^{1/T_i} \frac{1}{Z_{T_i}} \propto \exp\left\{ \frac{-H(\theta)}{T_i} \right\},$$

where $Z_{T_i} = \int g_{T_i}(\theta)d\theta$ and $T_1 < T_2 < \cdots < T_M$.

One iteration of the implemented parallel tempering algorithm consists of a normal additive random walk proposal Metropolis Hastings step to update each of the M chains, which upon setting the variance of the proposal distribution for each chain fixes the transition kernels $P_1, \ldots, P_M$ (for details refer to Algorithm 3). After each chain is updated, swaps between the chains are proposed. For each chain l=1,...,M we define the between-chain exchange probabilities according to which swaps are proposed to be

$$q_B(l,m) = P[\text{ Swap chain } l \text{ with chain } m \mid \text{ chain } l \text{ selected }]$$

$$= \begin{cases} \frac{1}{2} & 2 \le I < M, m = l-1, l+1 \\ 1 & l = 1 \text{ and } m = 2, \text{ or } l = M \text{ and } m = M-1 \\ 0 & \text{otherwise} \end{cases}$$

Then the between-chain exchange that proposes to swap $x_l$, the current state of chain l, with $x_m$, that of chain m, is accepted with probability

$$\alpha = \min\left\{ 1, \frac{\pi_l(x_m)\,\pi_m(x_l)\,q_B(m,l)}{\pi_l(x_l)\,\pi_m(x_m)\,q_B(l,m)} \right\}$$

For the tempered distributions

$$\frac{\pi_l(x_m)\,\pi_m(x_l)}{\pi_l(x_l)\,\pi_m(x_m)} = \exp\left\{ (H(x_l) - H(x_m)) \left( \frac{1}{T_l} - \frac{1}{T_m} \right) \right\}$$

The samples collected for $x_1$ across iterations represent a sample from the marginal $\pi_1 \equiv \pi$, the target distribution. In this particular implementation outlined in Algorithm 3 we propose swaps in a random order each iteration. For cleaner code we define a function called $rq_B(.)$ proposing the swapping partner given chain l is selected according to $q_B(l,.)$. Furthermore, the Metropolis steps are vectorised but will be written as a for loop in the pseudo code.

---
**Algorithm 3:** Parallel tempering algorithm for $g(\theta)$

---
**Inputs: $\mathbf{X_0} = (\mathbf{X_0^{(1)}}, \ldots, \mathbf{X_0^{(M)}})$** (vector of starting values), **Ts** (vector of temperatures),
$\qquad \sigma = (\sigma^{(1)}, \ldots, \sigma^{(M)})$ (proposal standard deviations for each temperature), **nsteps** (Number of
$\qquad$ iterations), **g** (Target up to proportionality)

**for** $i=1\ldots nsteps$ **do**
$\quad$ **for** $l=1\ldots M$ **do**
$\quad\quad$ Propose $Y \sim \mathcal{N}(X_{i-1}^{(l)}, \sigma^{(l)})$
$\quad\quad$ Simulate $V \sim \text{Uniform}(0,1)$
$\quad\quad$ **if** $V \le \alpha = min\left(1, \frac{g(Y)}{g(X_{i-1}^{(l)})}\right)$ **then**
$\quad\quad\quad$ $X_i^{(l)} = Y$
$\quad\quad$ **else**
$\quad\quad\quad$ $X_i^{(l)} = X_{i-1}^{(l)}$
$\quad\quad$ **end**
$\quad$ **end**
$\quad$ **for** $l$ *in permutation*$(\{1,...,M\})$ **do**
$\quad\quad$ Propose swapping partner $m \sim rq_B(l,.)$
$\quad\quad$ Simulate $V \sim \text{Uniform}(0,1)$
$\quad\quad$ **if** $V \le \alpha = min\left(1, \exp\left\{(H(x_l) - H(x_m))\left(\frac{1}{T_l} - \frac{1}{T_m}\right)\right\} \frac{q_B(m,l)}{q_B(l,m)}\right)$ **then**
$\quad\quad\quad$ Swap $X_i^{(l)}$ with $X_i^{(m)}$
$\quad\quad$ **end**
$\quad$ **end**
**end**

---

For the following discussion we use temperatures $Ts = (1, 2, 5, 10, 20, 50, 100)$ with transition kernels specified by the standard deviations $\sigma = (0.3, 0.6, 1, 1.6, 4, 9, 14)$. The variances are chosen increasingly as with increasing

temperature the tempered distributions become flatter so that big jumps are still accepted. The swaps will then accelerate the mixing. In Figure 16 the shapes of some tempered distributions can be observed. Figure 17 shows the traceplots for 6 distinct starting values and $10^4$ iterations indicating good convergence. Furthermore the Gelman plot is presented in Figure 18 and the Gelman Rubin diagnostic is computed to be 1 with upper confident interval 1.
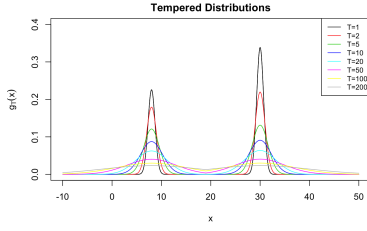


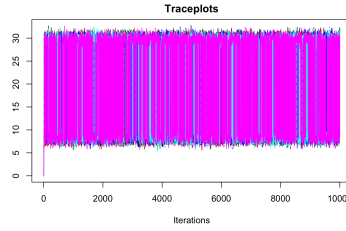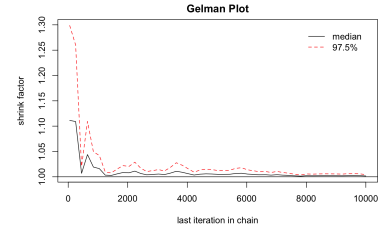Figure 16: Tempered distributions



Figure 17: Traceplots



Figure 18: Gelman plot

Taking starting the starting value $X_0 = 0$ for each chain and with the transition kernels as specified above we obtain a sample of length $9 \times 10^4$ after the first $10^4$ points are discarded as burn-in. Figure 20 shows the acf of this sample. The acf decreases relatively quickly which indicates good mixing and hence good exploration of the state space. The ESS is computed to be 4316. The Histogram presented in Figure 20 matches the shape of that of the true density well. For a better comparison the estimated density based on the sample is plotted together with the true density in Figure 21. The two densities are a very close fit with small discrepancies in the modes. Thus, we can conclude that the sampler is exploring the full state space and correctly sampling from the target distribution.



Figure 19: Acf of sample



Figure 20: Histogram of sample



Figure 21: Estimated vs true density

## c)

In this part we estimate $P(\theta > 20)$ where $\theta \sim g(.)$. The true value is numerically computed to be 0.6 with absolute error smaller $1.6 \times 10^{-5}$. Based on the sample obtained in the way as described above we compute the estimate

$$\hat{P}_{20} = \frac{1}{9 \times 10^4} \sum_{i=1}^{9 \times 10^4} I(X_i > 20) = 0.6005$$

which is accurate with an error of $5 \times 10^{-4}$.

## References

[1] Haario & Saksman & Tamminen. "Adaptive proposal for random walk Metropolis algorithm." In: *Computational Statistics* 14 (1999), pp. 375–395.

# Appendix

```r
1   # plot density
2   f <- function(x,y) exp(-x^2/100 - (y+3/100*x^2-3)^2)
3
4   x_vec <- seq(-15,15,0.1)
5   y_vec <- seq(-5,10,0.1)
6   grid <- expand.grid(x=x_vec,y=y_vec)
7   f_points <- f(grid$x, grid$y)
8
9   dat <- data.frame(x = c(grid$x), y=c(grid$y), f=c(f_points))
10
11  fig <- plot_ly(x = dat$x, y=dat$y, z=dat$f, type = 'mesh3d')
12  fig
13  # simple metropolis hastings
14  MetrHastw <- function(x0, sigmapropx, sigmapropy, nsteps, f){
15    X <- matrix(NA,nrow=2, ncol=nsteps+1)
16    X[,1] <- x0
17
18    for (i in 2:(nsteps+1)){
19      x <- rnorm(1,mean=X[1,(i-1)], sd=sigmapropx)
20      y <- rnorm(1,mean=X[2,(i-1)], sd=sigmapropy)
21      if (runif(1)<=min(exp(log(f(x,y))-log(f(X[1,i-1], X[2,i-1]))), 1)) X[,i] <- c(x↩
         ,y)
22      else X[,i] <- X[,i-1]
23    }
24    X[,-1]
25  }
26  # algorithm taking into account scale and dependency
27  # by specifying the covariance matrix of the proposal
28  MetrHastw_cov <- function(x0, sigmapropx, sigmapropy, cova, nsteps, f){
29    X <- matrix(NA,nrow=2, ncol=nsteps+1)
30    X[,1] <- x0
31
32    for (i in 2:(nsteps+1)){
33      prop <- rmvnorm(1,mean=X[,i-1], sigma=matrix(c(sigmapropx^2,cova,cova,↩
         sigmapropy^2), ncol=2))
34      x <- prop[1]
35      y <- prop[2]
36      if (runif(1)<=min(exp(log(f(x,y))-log(f(X[1,i-1], X[2,i-1]))), 1)) X[,i] <- c(x↩
         ,y)
37      else X[,i] <- X[,i-1]
38    }
39    X[,-1]
40  }
41
42
43  # tuning parameters
44  X_mh_adj <- MetrHastw_cov(c(0,3), 6.884244,2.122618, 0.2264226,10^5,f)
45  cov <- cov(X_mh_adj[1,], X_mh_adj[2,])
46  sd_x <- sd(X_mh[1,])
47  sd_y <- sd(X_mh[2,])
48  cov <- 0
49  sd_x <- 1
50  sd_y <- 1
51
52  for (i in 1:5){
53
```

```r
54     print(c(sd_x, sd_y,cov))
55     Dummy <- MetrHastw_cov(c(0,3), sd_x,sd_y, cov,10^5,f)
56     # update parameters
57     cov <- cov(Dummy[1,10^4:10^5], Dummy[2,10^4:10^5])
58     sd_x <- sd(Dummy[1,10^4:10^5])
59     sd_y <- sd(Dummy[2,10^4:10^5])
60
61   }
62   # Then trying several values around this estimate
63   # multiple chain convergence diagnostics
64   set.seed(0)
65   N <- 5*10^4 # chain length
66   X0s <- replicate(6, rnorm(2,c(0,3), 8))
67   samplesX <- matrix(rep(0,dim(X0s)[2]*N), nrow=dim(X0s)[2])
68   samplesY <- matrix(rep(0,dim(X0s)[2]*N), nrow=dim(X0s)[2])
69   for (i in 1:dim(X0s)[2]){
70     sample <- MetrHastw_cov(X0s[,i], 7 ,2.1, 0.22,N,f)
71     samplesX[i,] <- sample[1,]
72     samplesY[i,] <- sample[2,]
73   }
74
75
76   chainsX <- lapply(1:dim(X0s)[2], function(i) mcmc(samplesX[i,]))
77   chainsY <- lapply(1:dim(X0s)[2], function(i) mcmc(samplesY[i,]))
78   X0s
79   traceplot(chainsX, main='Traceplots for X')
80   traceplot(chainsY, main='Traceplots for Y')
81   gelman.diag(mcmc.list(chainsX))
82   gelman.diag(mcmc.list(chainsY))
83   gelman.plot(mcmc.list(chainsX), main='Gelman Plot for X')
84   gelman.plot(mcmc.list(chainsY), main='Gelman Plot for Y')
85
86   # sample
87   set.seed(0)
88   X_mh_final <- MetrHastw_cov(c(0,3), 7 ,2.1, 0.22,10^5,f)
89   plot(X_mh_final[1,], type='l')
90   plot(X_mh_final[2,], type='l')
91
92   plot(X_mh_final[1,(500:10^5)], X_mh_final[2,(500:10^5)], col=rgb(red=0, green=0, ↵
         blue=1, alpha=0.01), xlab='x', ylab='y', main='Joint Distribution f(x,y)', xlim↵
         =c(-25,25), ylim=c(-20,20))
93   acf(X_mh_final[1,10^4:10^5], lag.max=200, main='Acf of X')
94   acf(X_mh_final[2,10^4:10^5], lag.max=200, main='Acf of Y')
95
96   #mean estimates
97   mean(X_mh_final[1,(10^4+1):10^5])
98   mean(X_mh_final[2,(10^4+1):10^5])
99
100  #numerical integration for finding means
101  IntY <- function(x) sapply(x, function(b) integrate(function(y) f(b,y),-Inf,Inf)$↵
         value)
102  IntX <- function(y) sapply(y, function(b) integrate(function(x) f(x,b),-Inf,Inf)$↵
         value)
103  k <- 1/integrate(function(x) IntY(x), -Inf, Inf)$value
104
105  k*integrate(function(x) x*IntY(x), -Inf, Inf)$value
106  k*integrate(function(y) y*IntX(y), -Inf, Inf)$value
107
108  # adaptive proposal algorithm
109  c_2 <- 2.4/sqrt(2)
110  AP <- function(X0=c(0,3), H, U, nsteps, f){
```

```r
111
112    X <- matrix(NA,nrow=2, ncol=nsteps+H)
113    X[,1] <- X0
114    X[,2:(H)] <- MetrHastw_cov(X0, 7, 2, 0.22, H-1, f)
115
116    Kt <- X[,1:H]
117    Kt_centered <- Kt - apply(Kt, 1, mean)
118
119    for (i in (H+1):(nsteps+H)){
120
121       Y <- X[,i-1] +  c_2/sqrt(H-1)* Kt_centered  %*% rnorm(H)
122
123       if (runif(1)<=min(f(Y[1], Y[2])/f(X[1,i-1], X[2,i-1]), 1)) X[,i] <- Y
124       else X[,i] <- X[,i-1]
125
126       if ((i%%U)==0){
127          Kt <- X[,(i-H+1):(i)]
128          Kt_centered <- Kt - apply(Kt, 1, mean)
129       }
130    }
131    X[,(1+H):(H+nsteps)]
132 }
133
134 # ap multiplc chain conv diagnostics
135 set.seed(0)
136 N <- 5*10^4 # chain length
137 X0s <- replicate(6, rnorm(2,c(0,3), 8))
138 samplesX_ap <- matrix(rep(0,dim(X0s)[2]*N), nrow=dim(X0s)[2])
139 samplesY_ap <- matrix(rep(0,dim(X0s)[2]*N), nrow=dim(X0s)[2])
140 for (i in 1:dim(X0s)[2]){
141    sample_ap <- AP(X0s[,i],H=200, U=200, nsteps=N, f)
142    samplesX_ap[i,] <- sample_ap[1,]
143    samplesY_ap[i,] <- sample_ap[2,]
144 }
145
146
147 chainsX_ap <- lapply(1:dim(X0s)[2], function(i) mcmc(samplesX_ap[i,]))
148 chainsY_ap <- lapply(1:dim(X0s)[2], function(i) mcmc(samplesY_ap[i,]))
149 X0s
150 traceplot(chainsX_ap, main='Traceplots for X')
151 traceplot(chainsY_ap, main='Traceplots for Y')
152 gelman.diag(mcmc.list(chainsX_ap))
153 gelman.diag(mcmc.list(chainsY_ap))
154 gelman.plot(mcmc.list(chainsX_ap), main='Gelman Plot for X')
155 gelman.plot(mcmc.list(chainsY_ap), main='Gelman Plot for Y')
156
157 # sample
158 set.seed(0)
159 X_AP2 <- AP(c(0,3),H=200, U=200, nsteps=10^5, f)
160 plot(X_AP2[1,(10^4:10^5)], X_AP2[2,(10^4:10^5)], col=rgb(red=0, green=0, blue=1, ←
        alpha=0.01), xlab='x', ylab='y', main='Joint Distribution f(x,y)', xlim=c←
        (-25,25), ylim=c(-20,20))
161 acf(X_AP2[1,(10^4:10^5)], lag.max=200, main='Acf for X')
162 acf(X_AP2[2,(10^4:10^5)], lag.max=200, main='Acf for Y')
163 plot(X_AP2[1,(10^4:10^5)], type='l')
164 plot(X_AP2[2,(10^4:10^5)], type='l')
165
166
167
168 #effective sample size and estimates
169 mean(X_AP2[1,10^4:10^5])
```

```
170   mean(X_AP2[2,10^4:10^5])
171
172   ESS_ap2 <- c(effectiveSize(mcmc(X_AP2[1,10^4:10^5])), effectiveSize(mcmc(X_AP2↩
          [2,10^4:10^5])))
173   ESS_ap2
174   ESS_mh <- c(effectiveSize(mcmc(X_mh_final[1,10^4:10^5])), effectiveSize(mcmc(X_mh_↩
          final[2,10^4:10^5])))
175   ESS_mh
176   # time algorithms
177   system.time(replicate(10, AP(c(0,3),H=200, U=200, nsteps=5*10^4, f)))
178   system.time(replicate(10,MetrHastw_cov(c(0,3), 7 ,2.1, 0.22,5*10^4,f)))
```

<div align="center">Listing 2: Code for question 2</div>

```
 1   a <- 8
 2   g <- function(theta) 4/10*exp(-(a-theta)^2) + 6/10 * exp(-(30-theta)^2)
 3
 4   x <- seq(0,38,0.01)
 5   plot(x,g(x), type='l', main='g(.) up to proportionality')
 6
 7   MetrHastw1 <- function(X0, sigmaprop, nsteps, f){
 8     X <- numeric(nsteps+1)
 9     X[1] <- X0
10
11     for (i in 2:(nsteps+1)){
12       Y <- rnorm(1,mean=X[i-1], sd=sigmaprop)
13       if (runif(1)<=min(f(Y)/f(X[i-1]), 1)) X[i] <- Y
14       else X[i] <- X[i-1]
15     }
16     X[-1]
17   }
18   # not exploring entire state space
19   X0s <- seq(0,35,3)
20   chains <- lapply(1:length(X0s), function(i) mcmc(MetrHastw1(X0s[i], 1, 10^3, g)))
21   traceplot(chains, ylim=c(0,36), main='Traceplots for proposal SD=1')
22
23   X0s <- c(1,17,30,40)
24   chains <- lapply(1:length(X0s), function(i) mcmc(MetrHastw1(X0s[i], 20, 10^3, g)))
25   traceplot(chains, main='Traceplot for proposal SD=20')
26
27   # parallel tempering algorithm
28   g_temp <- function(theta,T) g(theta)^(1/T)
29   H <- function(x) -log(g(x))
30   # swapping probabilities
31   q_B <- function(l,m,M) {
32     if(2<=l & l<M & (m==(l-1) | m==l+1)) return(0.5)
33     if((l==1 & m==2) | (l==M & m==M-1)) return(1)
34     else return(0)
35   }
36   rq_B <- function(l,M){
37     if(l==1) return(l+1)
38     if(l==M) return(l-1)
39     if(2<=l & l<M) return(ifelse(runif(1)<1/2, l-1, l+1))
40   }
41
42   tempered <- function(Ts,x0s, sigmaprops, nsteps){
43     M <- length(Ts)
44     X <- matrix(rep(0,M*(nsteps+1)), nrow=M, ncol=nsteps+1)
45     X[,1] <- x0s
46     accepted_mh <- 0
```

```r
47      accepted_sw <- 0
48    # prob accept swap
49    for (i in 2:nsteps+1){
50      #proposal
51      prop <- X[,i-1] + rnorm(M,mean=rep(0,M), sd=sigmaprops)
52      U <- runif(M)
53      ind <- c(U<= sapply(1:M, function(j) exp(log(g_temp(prop[j],Ts[j]))- log(g_temp↩
          (X[j,i-1],Ts[j])))))
54      # mh update for each density
55      X[ind,i] <- prop[ind]
56      X[!ind,i] <- X[!ind, i-1]
57      accepted_mh <- accepted_mh + sum(ind)
58      # between moves
59      for(l in sample(1:M)){
60        m <- rq_B(l,M)
61        if (runif(1)<=min(1, exp((H(X[l,i])-H(X[m,i]))*(1/Ts[l]-1/Ts[m])))) {
62          X[c(l,m),i] <- X[c(m,l),i]
63          accepted_sw <- accepted_sw + 1}
64      }
65    }
66    print(accepted_mh/(nsteps*M))
67    print(accepted_sw/(nsteps*M))
68    return(X[,-1])
69 }
70
71 #plot tempered distributions
72 Ts <- c(1,2,5,10,20, 50, 100, 200)
73 x <- seq(-10,50,0.01)
74 k <- 1/integrate(g,lower=-Inf,upper=Inf)$value
75 plot(x,g(x)*k, type='l', main='Tempered Distributions', xlim=c(-10,50), ylim=c↩
      (0,.4), ylab=expression(g[T](x)))
76 for(i in 1:length(Ts)) lines(x, g_temp(x,Ts[i])/integrate(function(y) g_temp(y,Ts[i↩
      ]), lower=-Inf, upper=Inf)$value, col=i)
77 legend('topright', legend=c("T=1", "T=2","T=5","T=10","T=20","T=50","T=100","T=200"↩
      ),
78        col=1:(length(Ts)), lty=1, cex=0.8)
79
80 # multiple chains convergence diagnostics
81 X0s <- replicate(6, rnorm(length(Tsq), 0 ,10))
82 chains <- lapply(1:6, function(i) mcmc(tempered(Tsq, X0s[,i], c(0.3, 0.6, 1, 1.6, ↩
      4, 9,14), 10^4)[1,]))
83 gelman.diag(mcmc.list(chains))
84 traceplot(chains, main='Traceplots')
85 gelman.plot(mcmc.list(chains), main='Gelman Plot')
86 # sample
87 set.seed(0)
88 Tsq <- c(1,2, 5, 10,20, 50, 100)
89 X_tempq <- tempered(Tsq, rep(0,length(Tsq)), c(0.3, 0.6, 1, 1.6, 4, 9,14), 10^5)
90 hist(X_tempq[1,10^3:10^5], breaks=150, main='Histogram of the sample', xlab='x')
91 plot(X_tempq[1,10^3:10^5], type='l', main='Traceplot', ylab=expression(X[t]))
92 acf(X_tempq[1,10^3:10^5], lag.max=100, main='ACF')
93
94 #plot true density against estimated density
95 x <- seq(0,35,0.1)
96 plot(density(X_tempq[1,10^3:10^5], adjust=0.1 ),ylim=c(0,1), main='True density and↩
      estimated density')
97 lines(x, g(x)*k, col='red')
98 legend('topright', legend=c("Estimated", "True"),
99        col=c('black', 'red'), lty=1, cex=0.8)
100
101 # estimate and numerical integration
```

```r
102   mean(X_tempq[1,10^4:10^5]>20)
103   integrate(function(x) k*g(x), lower=20, upper=Inf)
```