

Convergence Diagnostic for Markov Chain Monte Carlo Methods

I. Abstract

This project aims to provide researchers utilizing MCMC methods with a comprehensive guide on the practical application of diagnostic tools to ensure dependable results. Two examples are presented to introduce the use of convergence diagnostic methods and compare the performance of different testing methods in different situations. To assess their convergence, code is written for five methods: Graphical observation, Effective sample size, Geweke, Heidelberg and Welch, Gelman-Rubin and kernel density methods to assess their convergence. In particular, the team focused on the Geweke and ESS calculations employed in prevalent R packages such as “coda” and “mcmcse”, while also including additional features such as considering batch means, spectral variance estimator method [2], and multivariate cases. Additionally, the team modified the Geweke code to consider batch means [3] and multivariate cases in an innovative manner.

II. Introduction

The Markov Chain Monte Carlo (MCMC) method is a widely used tool for statistical inference and simulation-based estimation. The idea of MCMC is that if simulating from a target density π is difficult, so that ordinary Monte Carlo methods based on independent and identically distributed (i.i.d.) samples cannot be used for inferences on π , it may be possible to construct a Markov chain $\{X_n\}_{(n \geq 0)}$ with stationary density to form Monte Carlo estimators [3].

It is particularly useful for problems where the likelihood function is intractable or computationally expensive. In such cases, MCMC methods allow users to simulate samples from the posterior distribution using only the likelihood function and prior distributions. However, MCMC simulations can suffer from slow or even failed convergence to the target distribution. This is because the algorithm works by constructing a Markov chain that produces a sequence of samples, and the convergence of the chain to the target distribution can be slow or uncertain.

To overcome these issues, it is crucial to use convergence diagnostics to assess the reliability of MCMC simulations. Convergence diagnostics are statistical tools that help evaluate whether the MCMC chain has converged to the target distribution. The diagnostic tools indicate whether the MCMC simulations have produced enough independent samples from the target distribution. The diagnostic tools can help identify issues with the MCMC simulations and provide guidance on improving the simulations. This paper aims to provide a comprehensive review of commonly used diagnostic tools for evaluating the convergence of MCMC simulations. The diagnostic tools covered in this paper include the Geweke statistic, Heidelberger and Welch test, effective sample size, Gelman-Rubin diagnostic, Kernel density-based method, and graphical outputs.

III. MCMC Diagnostics

The first code, titled “1. Generate two chains for testing”, is used to implement a Markov Chain Monte Carlo (MCMC) simulation using the Metropolis-Hastings (MH) algorithm. In this specific case, two chains are generated to test different proposal distributions. The first chain uses an exponential target distribution with $\lambda = 1$ of $\pi(x) = \exp(-x)$ and a proposal distribution of $P(x'|x) = 0.5\exp(-0.5x)$. The second chain uses the same target distribution of $\pi(x) = \exp(-x)$ but with a proposal distribution of $P(x'|x) = 5\exp(-5x)$. The purpose of the simulation is to generate a sequence of random samples from the target distribution that can be used for statistical inference or simulation-based estimation. In addition, the generated samples from the two chains will be used to assess the convergence of MCMC simulations by applying the aforementioned diagnostic methods. The code includes a burn-in period to allow the chains to converge to the target distribution, and the final output is a matrix of the generated samples after discarding the burn-in period if specified.

A. Effective Sample Size

The second code, titled “2. Effective sample size”, contains several functions that are used to calculate effective sample size (ESS) (which is defined in [3]) for Markov Chain Monte Carlo (MCMC) simulations. ESS is a measure of the number of independent samples in a MCMC sample, which is important for estimating the precision of the sample mean and other summary statistics. The “min.ESS” function computes the minimum ESS required for a given parameter space dimensionality p , a significance level α , and a desired error margin ε .

$$\widehat{mESS} \geq \frac{2^{2/p} \pi}{(p\Gamma(p/2))^{2/p}} \frac{\chi_{1-\alpha,p}^2}{\varepsilon^2}$$

where ε is the desired level of precision for the volume of the $100(1 - \alpha)\%$ asymptotic confidence region, and $\chi_{1-\alpha,p}^2$ is the $(1 - \alpha)$ quantile of $\chi_{1-\alpha,p}^2$.

The “get.ESS” function is used to estimate the ESS of a given MCMC chain X , using one of several methods depending on the value of the argument ESS. If ESS = 1, the function calculates ESS using the following formula (the most common definition (Robert and Casella 2004) [2]):

$$ESS = \frac{n}{1 + 2 \sum_{i=1}^{\infty} Corr(g(X_0), g(X_i))}$$

where g is a real valued function.

If ESS is 'coda', then it uses the `mcse.multi` function from the “coda” package. If ESS is “batchmeans”, then it uses the `mcse_batchmeans` function defined in the code [2]. (These two are equivalent while method “batchmeans” will use the local “mcse_batchmeans” function). This two method also considered multivariate settings, that is, when $p \geq 1$, Vats et al. (2019) define multivariate ESS (mESS) as:

$$mESS = n \left(\frac{\Lambda_g}{\Sigma_g} \right)^{1/p}$$

where Λ_g is the population covariance matrix and Σ_g is defined as an estimate of the spectral density at frequency zero (introduced later).

The “mcse_batchmeans” function is used to estimate the Monte Carlo standard error (MCSE) of a given MCMC chain x using the batch means method. This involves dividing the chain into batches, computing the mean of each batch, and then estimating the standard deviation of the batch means. The MCSE is then estimated as the standard deviation of the batch means divided by the square root of the number of batches. In the formula above, Λ_g is the square of MCSE.

The spectrum0.ar function is used to compute the spectral density [2] and autoregressive (AR) model order for each column of a given matrix \mathbf{X} (Σ_g in the above formula). The function fits each column to a linear regression model and checks if its residuals are zero. If the residuals are zero, then the column is regarded as a constant sequence with a spectral density of zero and an AR model order of zero. Otherwise, an AR model is used to fit the column, and its spectral density and AR model order is computed.

Spectral estimation values are used to estimate the variances before and after a given point, and Geweke's statistics are then calculated by taking the weighted average of these two variances. The zero-frequency spectral estimate is used as the variance estimate value for Geweke's statistics calculation, as it can reflect the overall variance in time series data and remove the influence of autoregressive coefficients from spectral estimation. However, the AR method is not always more accurate than directly calculating the sample variance, especially when there are outliers or anomalies in the data that may affect the stationarity of the time series.

B. Geweke

The third code, titled “3. Geweke”, is used to test the convergence of a Markov Chain using Geweke's test. Geweke's test calculates the standardized difference between the means of the first a_n samples (SSA) and the last b_n samples of the chain (SSB), and tests whether this value is significantly different from zero using a two-sided t-test [1], where

$$\text{where } SSA = \frac{1}{A-1} \sum_{t=1}^A (X_t - \bar{X}_{1:A})^2, \quad SSB = \frac{1}{n-B} \sum_{t=B}^{n-1} (X_t - \bar{X}_{B:n-1})^2$$

$$z = \frac{\bar{X}_{1:A} - \bar{X}_{B:n-1}}{\sqrt{\frac{SSA}{A} + \frac{SSB}{n-B+1}}}$$

The function “geweke_toy” implements Geweke's test and returns the Geweke statistic and p-value for a given Markov chain. The function takes in the Markov chain (chain), the proportions of the chain to use for the first set of samples (a) and the last set of samples (b), and the method to use for variance calculation (method). Two methods are available for variance calculation: “normal” uses the var function, while “spectral” uses a spectral variance estimator calculated by the spectrum0.ar function.

The function “geweke” is a more complex implementation of Geweke's test that divides the Markov chain into multiple sub-sequences to increase the sample size and improve the accuracy of the test. The function takes in the Markov chain (x), the proportions of the chain to use for the first set of samples (a) and the last set of samples (b), the number of batches [2] to divide the chain into (num_batches), and the method to use for variance calculation (method).

The function initializes vectors to store the Geweke statistics and p-values for each batch, and then calculates these using the same method as “geweke_toy”. It then returns a list of the Geweke statistics and p-values for each batch. Moreover, it calculates the proportion of p-values below 0.05 among all p-values and plots the statistics and p-values, so a high proportion means that p-values calculated in multiple batches tend to reject the null hypothesis: the chain has converged.

A function for multi-dimensional cases using a “geweke_multi” function is also considered, which performs Geweke's convergence diagnostic test on a multi-dimensional input data. The function takes in additional parameters such as num_batches and method to specify the number of batches to use and the method for calculating the test statistic, respectively. The function returns the test statistics and p-values for each dimension of the input data. The code then

generates a random multi-dimensional input data using the ‘mvrnorm’ function from the “MASS” package, and applies the “geweke_multi” function to it.

C. Heidelberg and Welch

The fourth R code, titled “4. Heidelberg and Welch”, is an implementation of the Heidelberg and Welch test for checking the stationarity and adequacy of Markov chains, as described in Heidelberg and Welch's papers from 1981 and 1983. The test consists of two parts: a stationary portion test and a half-width test. The stationary portion test assesses the stationarity of a Markov chain by testing the hypothesis that the chain comes from a covariance stationary process. The half-width test [1] checks whether the Markov chain sample size is adequate to estimate the mean values accurately.

Given: $\{X_t\}$, $S_0 = 0$, $S_n = \sum_{t=1}^n X_t$, and $\bar{X} = \frac{1}{n} \sum_{t=1}^n X_t$. The following sequence can be constructed with s coordinates on values from $\frac{1}{n}, \frac{2}{n}, \dots, 1$:

$$B_n(s) = \frac{S_{[ns]} - [ns]\bar{X}}{n\hat{p}(0)}$$

where $[]$ is the rounding operator, and $\hat{p}(0)$ is an estimate of the spectral density at zero frequency that uses the second half of the sequence.

The statistic used in these procedures is the Cramer–von Mises statistic; that is $\int_0^1 B_n(S)^2 ds$. As $N \rightarrow \infty$, the statistic converges in distribution to a standard Cramer–von Mises distribution.

The “pcramer” function is a helper function that computes the probability of exceeding the given value of a test statistic, based on the Cramer-Lundberg approximation. The Cramer-Lundberg approximation approximates the Cramer-von Mises test and is used to calculate the p-value of the Cramer-von Mises statistic. The form of the Cramer-Lundberg approximation is as follows:

$$p_{value} = \sum_{k=0}^3 z_k \cdot e^{-u_k} \cdot K_{\frac{1}{4}}(u_k)$$

$$\text{where } z_k = \frac{\Gamma(k + 0.5) \cdot \sqrt{4k - 1}}{\Gamma(k + 1) \cdot \pi^{\frac{3}{2}} \cdot \sqrt{q}}, \quad u_k = \frac{(4k + 1)^2}{16q}$$

$$\text{and } K_{\frac{1}{4}}(x) = z_k \cdot e^{-x} \cdot \text{besselK}(x = u_k, \nu = 1/4)$$

This test can be performed repeatedly on the same chain, and it helps identify a time t when the chain has reached stationarity. The whole chain, $\{\theta^t\}$, is first used to construct the Cramer–von Mises statistic. If it passes the test, the conclusion is that the entire chain is stationary. If it fails the test, the initial 10% of the chain is dropped and the test redone by using the remaining 90%. This process is repeated until either a time t is selected, or it reaches a point where there is not enough data remaining to construct a confidence interval (the cutoff proportion is set to be 50%).

The RHW quantifies accuracy of the $1 - \alpha$ level confidence interval of the mean estimate by measuring the ratio between the sample standard error of the mean and the mean itself. In other words, the Markov chain is stopped if the variability of the mean stabilizes with respect to the “mean.val”. The RHW for a confidence interval of level $1 - \alpha$ is:

$$RHW = \frac{z_{(1-\alpha/2)} \cdot (\hat{s}_n/n)^{1/2}}{\hat{\theta}}$$

The “Heidelberger” function takes as input a matrix of data x and optional parameters “eps” and “p-value.” It returns a matrix containing the test results for each column of x , including the starting position of the stationary portion, the p-value of the test, the half-width of the stationary portion, and whether the half-width test passed or failed.

D. Gelman-Rubin

The fifth code, titled “5. Gelman-Rubin”, implements the Gelman-Rubin diagnostic to verify if parallel chains with dispersed initial values converge to the same target distribution. This method is useful for detecting multi-modal posterior distribution and to identify the need to run a longer chain. The function “gelman_toy” calculates potential scale reduction factors (PSRF) by estimating the between-chain variance, within-chain variance, and the PSRF \widehat{R}_c .

There are two main problems with the toy model. The first small problem is that the improved version of Brooks and Gelman (1997) is not considered when calculating the PSRF statistics.

This improved version considers the possible existence of multicollinearity between samples and considers the effective degrees of freedom at the end of the calculation. The second problem is that the function only considers chains with the same target distribution. That is, even if multiple Markov chains are input, they are considered to have the same target distribution.

Next, ways to solve the above problems are considered. In Brooks and Gelman's (1998) modified version of PSRF, \hat{d} represents the effective number of degrees of freedom for the model parameters, calculated as: $\frac{2\hat{V}^2}{\widehat{Var}(\hat{V})}$. At the same time, when calculating the variance, the covariance between different variables needs to be calculated (because the target distribution may be different at this time), so cov.wb needs to calculate the covariance and use if statements to judge whether it is multivariate. The specific calculation method is below [1]:

$$B = \frac{n}{M-1} \sum_{m=1}^M (\bar{X}_m - \bar{X})^2, \text{ where } \bar{X}_m = \frac{1}{n} \sum_{t=1}^n X_m^t, \bar{X} = \frac{1}{M} \sum_{m=1}^M \bar{X}_m$$

$$W = \frac{1}{M} \sum_{m=1}^M s_m^2, \text{ where } s_m^2 = \frac{1}{n-1} \sum_{t=1}^n (X_m^t - \bar{X}_m)^2$$

The posterior marginal variance, $var(\theta|y)$, is a weighted average of W and B. The estimate of the variance:

$$\hat{V} = \frac{n-1}{n} W + \frac{M+1}{nM} B$$

A refined version of PSRF is calculated, as suggested by Brooks and Gelman (1997), as:

$$\widehat{R}_c = \sqrt{\frac{\hat{d}+3}{\hat{d}+1} \cdot \frac{\hat{V}}{W}} = \sqrt{\frac{\hat{d}+3}{\hat{d}+1} \left(\frac{n-1}{n} + \frac{M+1}{nM} \frac{B}{W} \right)}$$

$$\text{where } \hat{d} = \frac{2\hat{V}^2}{\widehat{Var}(\hat{V})}$$

$$\text{where } \widehat{Var}(\bar{V}) = \left(\frac{n-1}{n}\right)^2 \frac{1}{M} \widehat{Var}(s_m^2) + \left(\frac{M+1}{nM}\right)^2 \frac{2}{M-1} B^2 \\ + 2 \frac{(M+1)(n-1)}{n^2 M} \frac{n}{M} (\widehat{cov}(s_m^2, (\bar{X}_m)^2) - 2\bar{X}_m \widehat{cov}(s_m^2, \bar{X}_m))$$

The refined function "gelman" solves the problems in the toy example by considering the adequate number of degrees of freedom for the model parameters (using \hat{d} to modify the variance estimation bias and reflect the degree of autocorrelation of MCMC simulation results). The advanced function also considers when the target distributions are different from each other.

When the “autoburnin” parameter is set to TRUE by default, the "gelman" function automatically removes the first half of the MCMC chain (the burn-in period). In practical experience, it has been repeatedly verified that when the option "autoburnin" is set to TRUE, non-convergent results tend to be produced by the function (the potential scale reduction factors increase, as do the point estimates and upper confidence intervals). An investigation is needed to determine whether setting “autoburnin” to TRUE by default is appropriate.

A large PSRF indicates that the between-chain variance is substantially greater than the within-chain variance, so a more extended simulation is needed. If the PSRF is close to 1, it can be concluded that each M chain has stabilized, and they will likely reach the target distribution.

It is best to choose different initial values for all M chains. The initial values should be as dispersed from each other as possible so that the Markov chains can fully explore different parts of the distribution before they converge on the target. Similar initial values can be risky because all chains can get stuck in a local maximum; that is something this convergence test cannot detect. If initial values for all the chains are not supplied, then the procedures generate them.

E. Kernel density-based methods [3]

There are MCMC diagnostics which compute distance between the kernel density estimates of two chains or two parts of a single chain and conclude convergence when the divergence is close to zero. These tools are intended to assess the convergence of the whole distributions [3].

Let $\{X_{ij}: i = 1, 2; j = 1, 2, \dots, n\}$ be the n observations obtained from each of the two Markov chains initialized from two points well separated with respect to the target density π . The adaptive kernel density estimates of observations obtained from the two chains are denoted by p_{1n} and p_{2n} respectively. The Kullback-Leibler (KL) divergence between p_{in} and p_{jn} is denoted by $KL(p_{in}|p_{jn}), i \neq j, i, j = 1, 2$, that is,

$$KL(p_{in}|p_{jn}) = \int p_{in}(x) \log \frac{p_{in}(x)}{p_{jn}(x)} dx.$$

The KL divergence has some important properties:

1. Non-negativity: The KL divergence is always non-negative ($KL(P \parallel Q) \geq 0$), and is zero if and only if P and Q are the same distribution.
2. Not symmetric: The KL divergence of Q from P is not the same as the KL divergence of P from Q .

This maximum KL divergence gives an indication of how much the chains differ in their distributions. If the chains are all converging to the same underlying distribution (as they should in a well-behaved MCMC), this maximum KL divergence should be small. If the chains are not converging to the same distribution, the maximum KL divergence will be larger.

Specifically in the Tool1 (different chains), the statistic represents the maximum symmetric KL divergence among the adaptive kernel density estimates of multiple MCMC chains. Initially, the symmetric KL divergence for each pair of chains is computed. Since the KL divergence is a measure of the difference between two probability distributions, whereas the symmetric KL divergence is the average of the KL divergences in both directions. The results for each pair of chains are stored in a matrix. The maximum symmetric KL divergence is then selected from this matrix to form the T1 statistic. This process provides a measure of the similarity between the MCMC chains, assisting in the assessment of whether the chains have adequately mixed, that is, whether they have reached their stationary distributions.

For multimodal target distributions, if all chains are stuck at the same mode, then empirical convergence diagnostics based solely on MCMC samples may falsely treat the target density as unimodal and are prone to failure. In such situations, Dixit and Roy (2017) propose another tool (Tool2) that makes use of the KL divergence between the kernel density estimate of MCMC samples and the target density (generally known up to the unknown normalizing constant) to detect divergence.

More specifically, T_2 represents the percentage of the target distribution that the Markov chain has yet to reach. If T_2 is close to 0, this indicates that the chain has reached its target distribution. In particular, let $\pi(x) = f(x)/c$, where $c = \int f(x)dx$ is the unknown normalizing constant. Dixit and Roy (2017)'s Tool2 is given by

$$T_2^* = \frac{|\hat{c} - c^*|}{c^*}$$

where \hat{c} is a Monte Carlo estimate, as described in section 3.3 of Dixit and Roy (2017) [4], of the unknown normalizing constant (c), based on the KL divergence between the adaptive kernel density estimate of the chain and π , and c^* is an estimate of c obtained by numerical integration.

Dixit and Roy (2017) discuss that T_2^* can be interpreted as the percentage of the target distribution not yet captured by the Markov chain. Using this interpretation, they advocate that if $T_2^* > 0.05$, then the Markov chain has not yet captured the target distribution adequately. Since the calculation involves numerical integration, it cannot be used in high-dimensional examples.

F. Graphical Analysis

The sixth and final code, titled “6. Graphical methods”, performs a Markov Chain Monte Carlo (MCMC) simulation using the Metropolis-Hastings algorithm to estimate the posterior distribution of a target distribution with a known density function $\pi(x)$. The final MCMC chain is then plotted using the ggplot2 package, and the variance and summary statistics are also computed. The “mcmcse” package is used to estimate the effective sample size of the chain, and the acf() function is used to plot the autocorrelation of the chain.

IV. Comparison

Table 1: Convergence Diagnostic in the Bayesian Procedures (MCMC)[1]

Name	Description	Interpretation of the Test
Gelman-Rubin	Uses parallel chains with dispersed initial values to test whether they all converge to the same target distribution. Failure could indicate the presence of a multi-mode posterior distribution (different chains converge to different local modes) or the need to run a longer chain (burn-in is yet to be completed).	One-sided test based on a variance ratio test statistic. Large \widehat{R}_c values indicate rejection.
Geweke	Tests whether the mean estimates have converged by comparing means from the early and latter part of the Markov chain.	Two-sided test based on a z-score statistic. Large absolute z values indicate rejection.
Heidelberger-Welch (stationarity test)	Tests whether the Markov chain is a covariance (or weakly) stationary process. Failure could indicate that a longer Markov chain is needed.	One-sided test based on a Cramer–von Mises statistic. Small p-values indicate rejection.
Heidelberger-Welch (half-width test)	Reports whether the sample size is adequate to meet the required accuracy for the mean estimate. Failure could indicate that a longer Markov chain is needed.	If a relative half-width statistic is greater than a predetermined accuracy measure, this indicates rejection.
Effective Sample Size	Relates to autocorrelation; measures mixing of the Markov chain.	Large discrepancy between the effective sample size and the simulation sample size indicates poor mixing.
Kernel density-based methods	Tool1 consists of calculating KL divergence for each pair of chains, and possibly a symmetrical KL divergence test that takes the average of KL divergence in two directions. Tool2 measures the percentage of the target distribution that the MCMC chain has not yet covered. Having good performance in multimodal.	The small KL divergence value suggests that the chains have a high degree of similarity, implying that they may have reached a stationary distribution. On the contrary, a larger value indicates potential divergence and the need for further iterations.

V. Examples

A. An exponential target distribution

The first chain with the target distribution is $\pi(x) = \exp(-x)$, and the proposal density function of $P(x'|x) = 0.5\exp(-0.5x)$

The second chain with the target distribution is $\pi(x) = \exp(-x)$, and the proposal density function of $P(x'|x) = 5\exp(-5x)$.

Each chain is run with 323,700 iterations and then convergence diagnostics are performed.

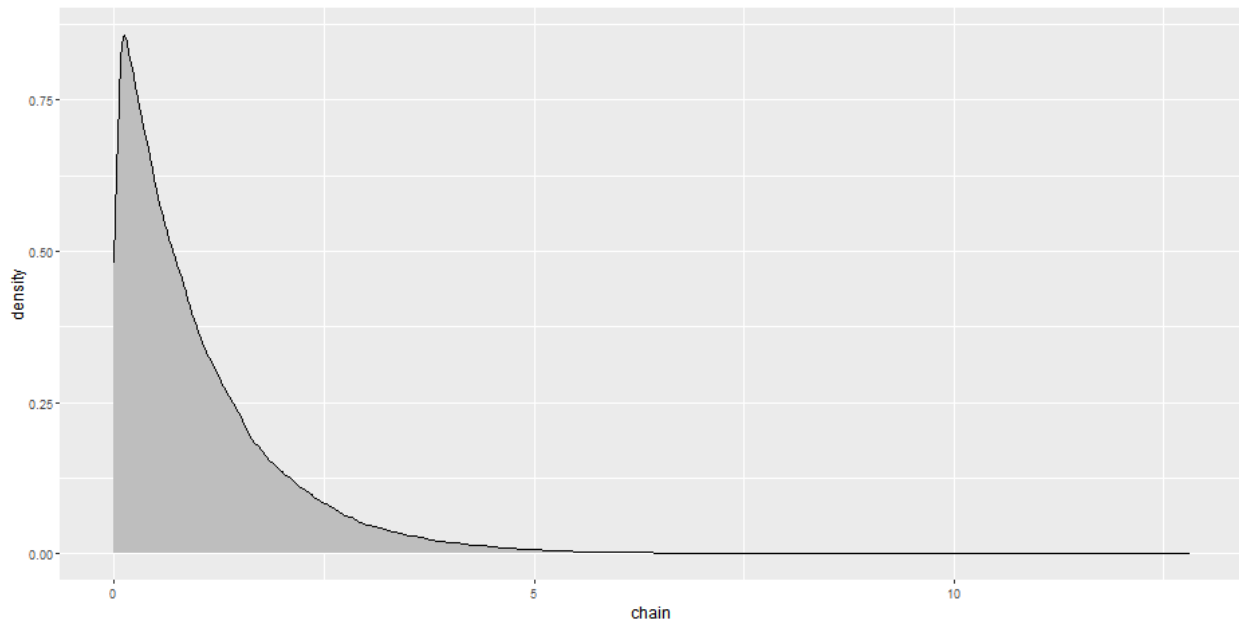
Notably, due to improper proposal distribution selection, the second Markov chain tends not to converge, which can be seen from the graphical representation in Figure 6, where the ACF function values remain high even at lag=50, while the first chain quickly converges.

Therefore, an appropriate convergence criterion would lead to accepting the null hypothesis of convergence for the first chain and rejecting the null hypothesis for the second chain.

1. Graphical Output Results

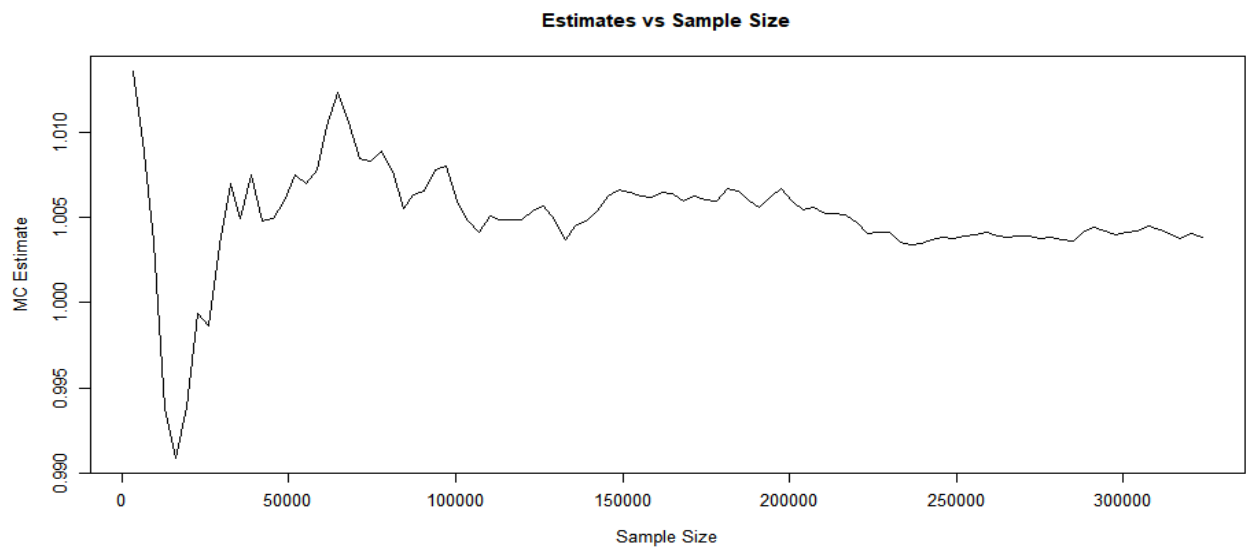
First chain:

Figure 1: Empirical pdf simulation of the first chain ($\pi(x) = \exp(-x)$)



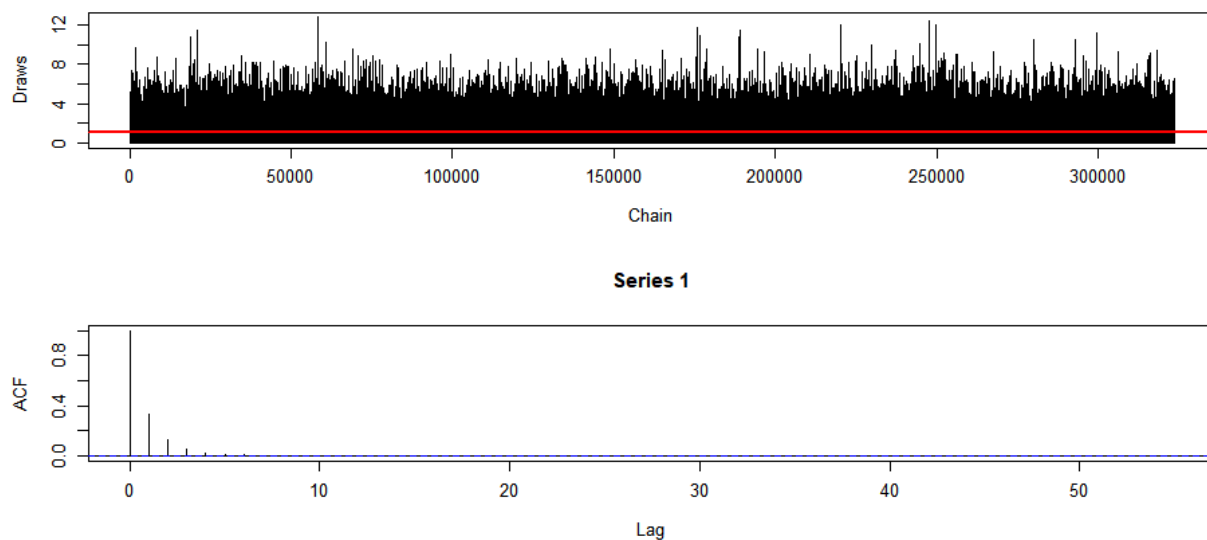
The function “estvssamp” plots the Monte Carlo estimates versus the sample size for a component of the MCMC output, indicating whether the Monte Carlo estimate has stabilized:

Figure 2: “estvssamp” plots of the first chain



“ts” is used to represent a set of data arranged in chronological order, “acf” will show the autocorrelation function image:

Figure 3: Time series plot and autocorrelation function plot of the first chain



From the first figure, it can be seen that the empirical pdf fits very well, and the simulated histogram is very close to the *Exponential*(1) distribution.

In the second figure, the mean of the Markov chain stabilizes around 1.005 after 100,000 iterations.

The time series plot in the third figure indicates that there is no apparent stagnation, indicating that the chain does not have periods of low acceptance rates. The ACF plot also shows that the autocorrelation coefficients of the chain are close to 0 after lag=5, indicating that the obtained samples are almost independent and identically distributed (i.i.d.).

Second chain:

Figure 4: Empirical pdf simulation of the second chain ($\pi(x) = \exp(-x)$)

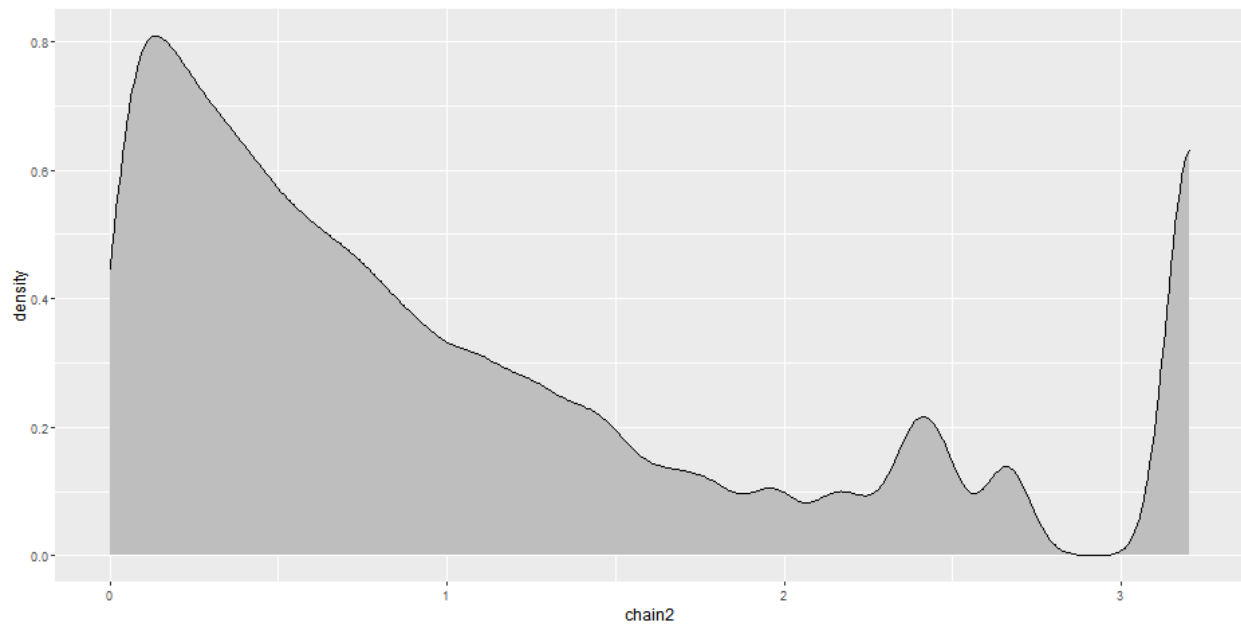


Figure 5: “estvssamp” plots of the second chain

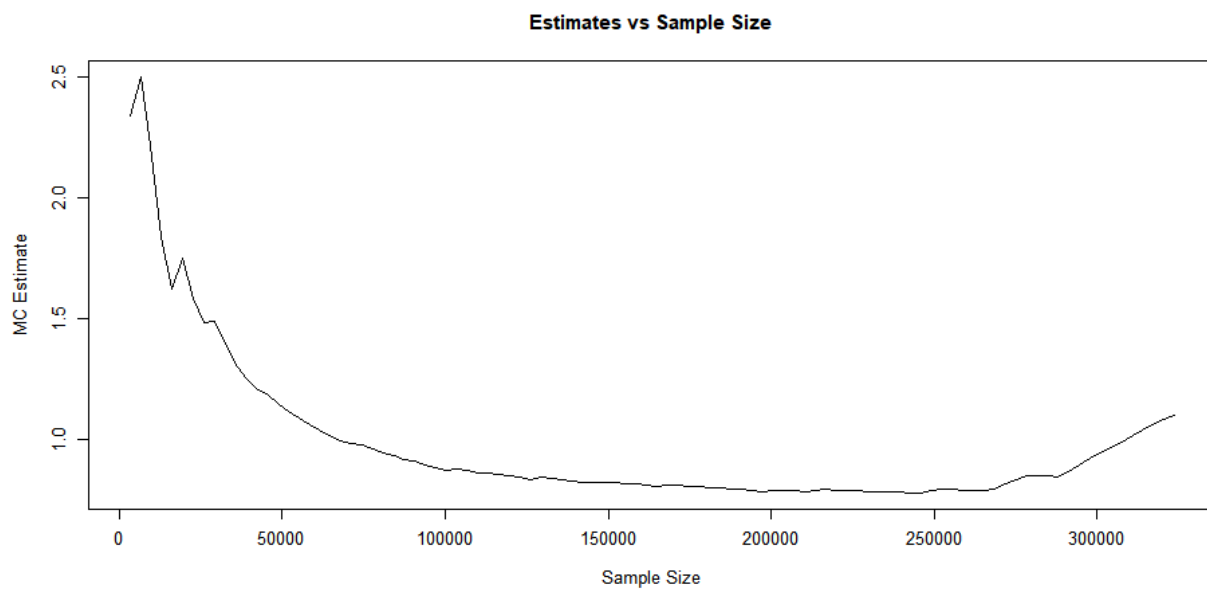
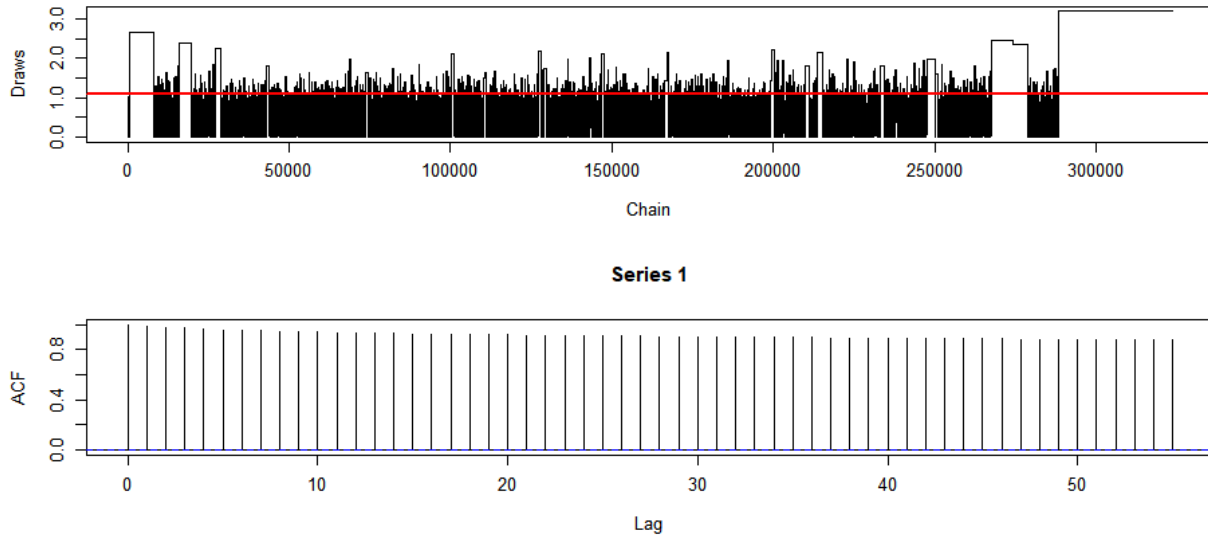


Figure 6: Time series plot and autocorrelation function plot of the second chain



From the fourth figure, it can be seen that the empirical pdf fits very poorly, and the simulated histogram is significantly different from the *Exponential*(1) distribution. In fact, the simulated density even has high-density regions in the end part that violate the target function.

Some clues from the fifth figure can be seen: the mean has a clear increasing trend before and after 300,000 iterations, the whole plot keeps changing substantially at the same time.

In the time series plot of the sixth figure, it can be seen that the chain has many obvious stagnations, indicating that there are many periods of low acceptance rates. There are even periods of complete stagnation before and after 300,000 iterations. The ACF plot also shows that the autocorrelation coefficients of the chain are still large after lag=50, indicating that the obtained samples are not independent and identically distributed (i.i.d.).

In summary, we can conclude from the graphical analysis that the first chain has converged while the second chain has not yet converged.

2. ESS Results

Using the minESS function, the **minimum ESS of chain 1 and chain 2 is 153,658.4** for given dimensionality 1 of the parameter space, significance level $\alpha = 0.05$, and the desired error margin $\epsilon = 0.05$.

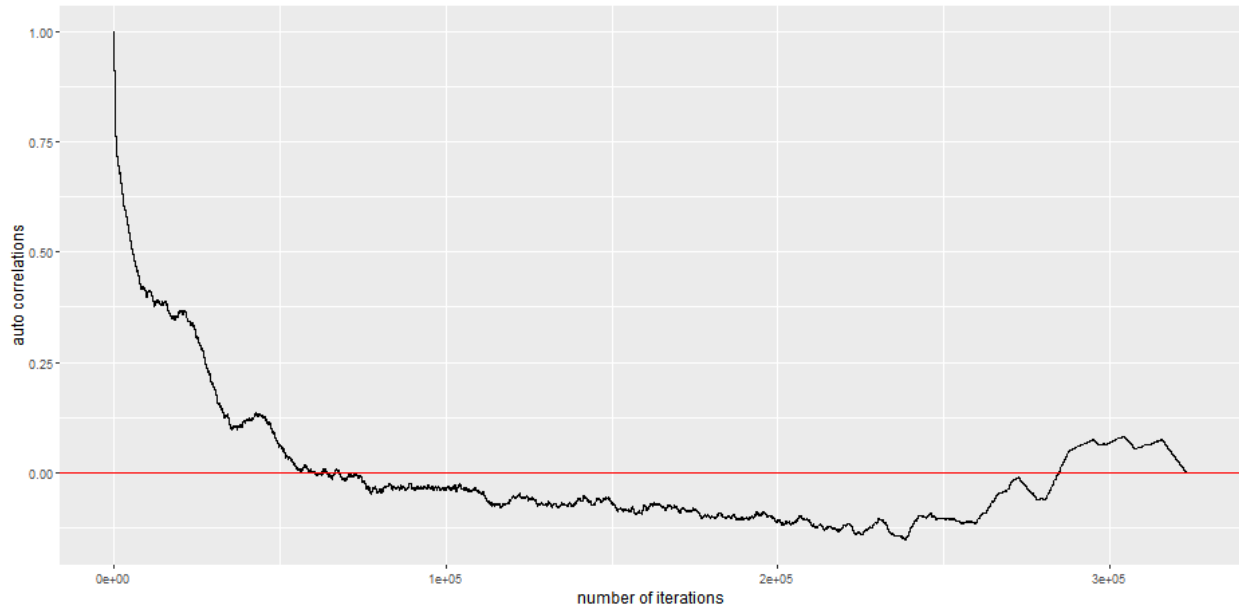
Table 2: Result of different chains using different ESS calculation methods

ESS Method	Calculation	Description	ESS of the first chain	ESS of the second chain
Default	$M / (1 + 2 * \text{sum}(\rho))$	Basic definition	161,849.3	161,849
Coda	$M * (\lambda^2 / \sigma_a^2)$	This method is the default in Coda	154,088.2	84.15389
Batch Means	$M * (\lambda^2 / \sigma_a^2)$	Same as above but with local batch_means function	155,330.8	3,557.791
Regularized	$\sigma^2_{sq} = \lambda^2_{sq} + 2 * \text{sum}(\rho)$	Large sample variance of the sample mean using	162,173	162,220.6

The unreasonable results are marked in red. It can be seen that the ESS values calculated by the four methods for the first Markov chain are consistent with the fact that the chain has converged, but the ESS values calculated by the first and fourth methods for the second chain do not match the actual situation. Why is that?

Due to the high autocorrelation of the second chain, the ESS values calculated by the first and fourth methods appear inflated. This is also evident from the ACF plot, which shows that the ACF value of the second chain is still very high at lag = 50, indicating strong autocorrelation. When calculating the sum of autocorrelation coefficients, a larger value should have been obtained. However, the algorithm returned 0.5, the result $(323700 / (1 + 2 * 0.5)) = 161,850$ obtained by the first method, indicating a much lower autocorrelation. The plot below may provide some clues as to why this is the case.

Figure 7: Autocorrelation verses numbers of iterations of the second chain



From the figure above it can be seen that the autocorrelation coefficients of the second chain remained high during the first tens of thousands of iterations, indicating poor convergence (high autocorrelation means that adjacent samples are highly correlated, indicating that the samples obtained are not i.i.d.). In fact, a sum of 14,526.14 can be obtained for the first 50,000 terms. However, as the number of iterations increases, the autocorrelation coefficients gradually change from positive to negative. Surprisingly, these negative autocorrelation coefficients cancel out the initial positive ones. The cumulative value of these coefficients at iteration 323,299 is 0.5, leading to an erroneous conclusion.

Therefore, these two methods may be unreliable in cases where the autocorrelation is high.

3. Geweke Results

Firstly, the results of the toy function are presented, and it is found that in the case where it is known that the first chain has converged while the second chain has not yet converged, both methods for computing the variance give incorrect conclusions for the toy function:

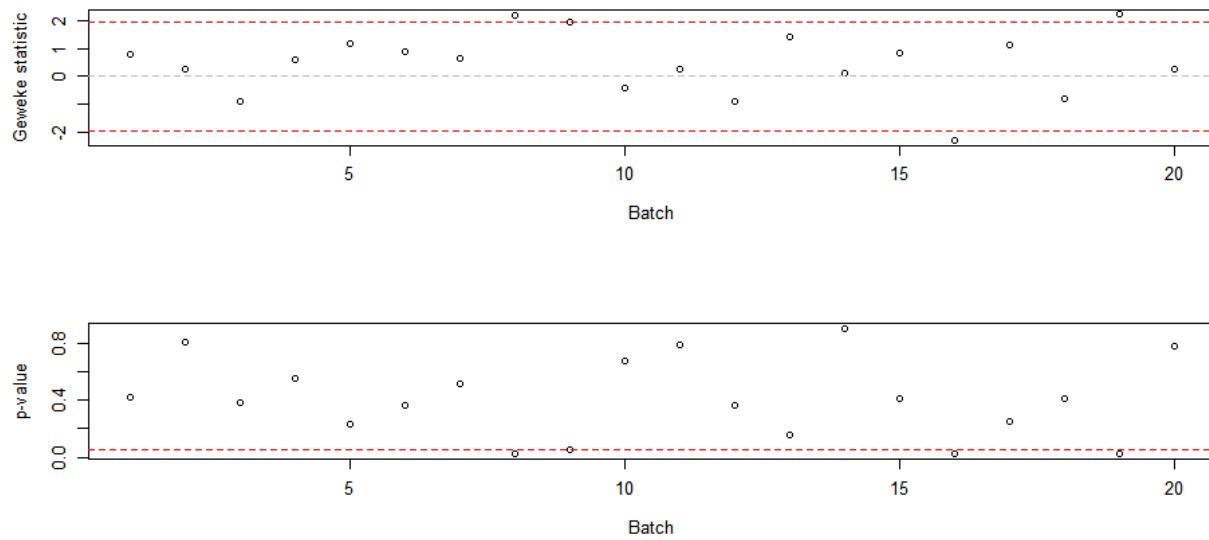
Table 3: Result using different variance calculation methods with the toy function

Geweke Method	Calculation	Description	P-value of the first chain	P-value of the second chain
Normal	$M / (1 + 2 * \text{sum}(\rho))$	Calculate the variance directly using the var() function	0.338618	0.455918
Spectral	$M * (\lambda^2 / \sigma^2)$	Calculate the variance using the spectral variance estimator	0.5085427	0.9845621

Now, the improved Geweke function will be demonstrated to see if it has actually improved. Of course, since the refined version returns results for multiple batches, the two functions cannot be directly compared side by side.

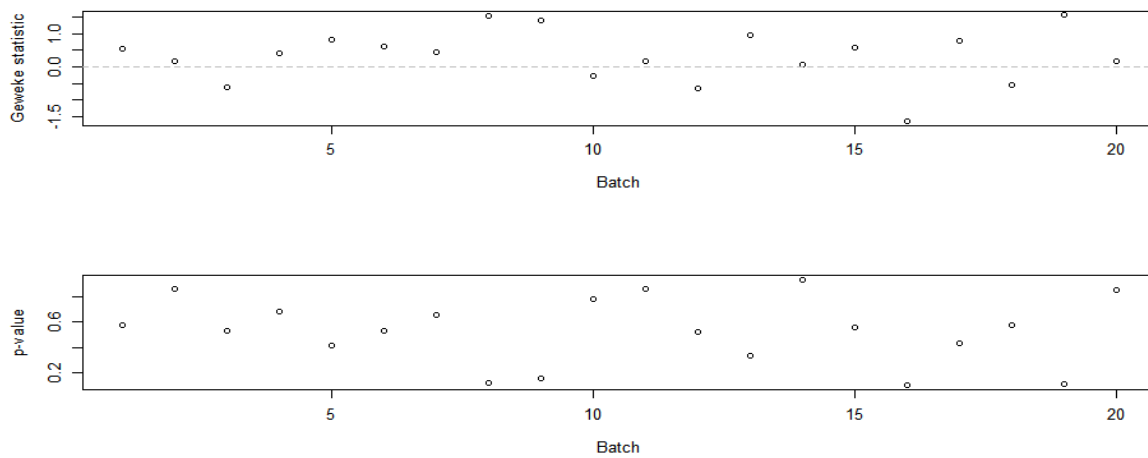
First chain:

Figure 8: Geweke Statistic and corresponding P-value using normal variance function of the first chain



Proportion of $p_value < 0.05 = 20\%$

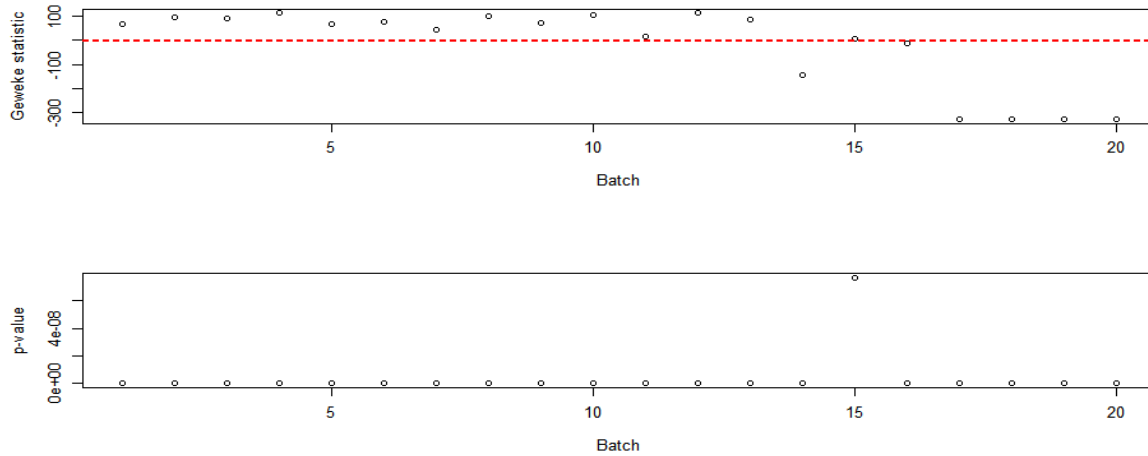
Figure 9: Geweke Statistic and corresponding P-value using spectral variance estimator of the first chain



Proportion of $p_value < 0.05 = 0\%$

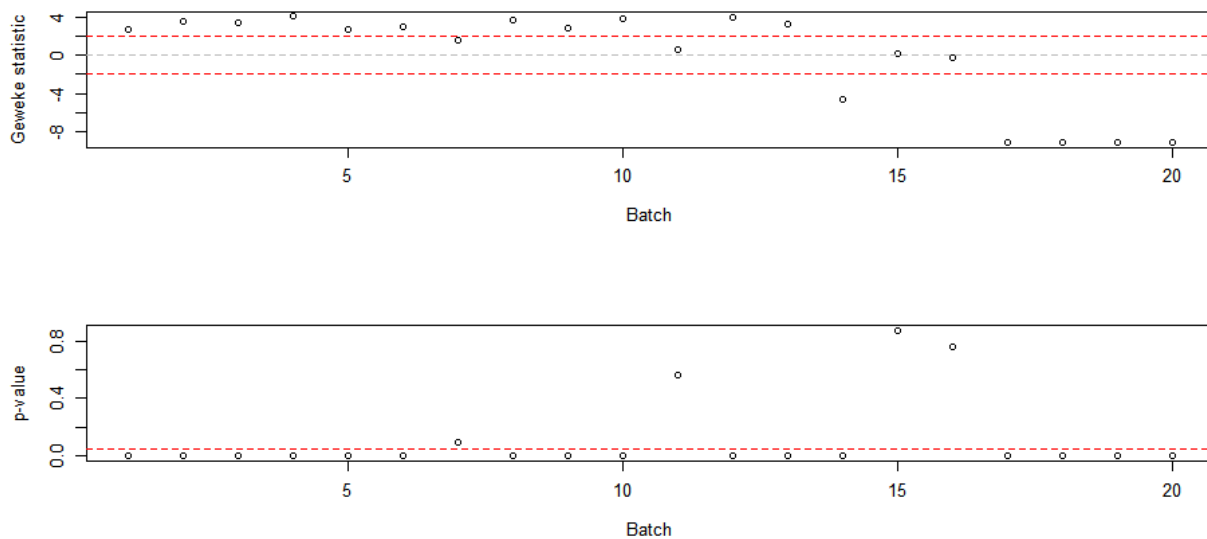
Second chain:

Figure 10: Geweke Statistic and corresponding P-value using normal variance function of the second chain



Proportion of $p_value < 0.05 = 100\%$

Figure 11: Geweke Statistic and corresponding P-value using spectral variance estimator of the second chain



Proportion of $p_value < 0.05 = 80\%$

For the first chain, both methods for computing the variance give results consistent with the truth and the conclusion mentioned above.

From the above figures, it seems that using the `var()` function directly to compute the variance gives a "more accurate" conclusion, but is this really the case?

In fact, it can be seen from Figure 10 that the absolute values of most Geweke statistics are very large, to a degree that seems to be beyond common sense. Further investigation revealed that directly calculating the variance using the `var()` function results in a very small variance (in fact, the second chain did not exhibit a large variance in the mean trajectory plot), but using the Spectral Density Estimate at Zero Frequency to estimate the variance yields a very large estimate (sometimes greater than 1000).

Therefore, although the absolute values of the computed test statistics may all be greater than 1.96, the absolute value of the test statistic will be large due to the small variance (a part of the denominator) computed using the `var()` function, while the Spectral Density Estimator will be relatively small (but still greater than 1.96, and the two-sided p-value will be less than 0.05).

Investigation into this issue needs to be done further in future research. However, in terms of the conclusion of this research, the improved version of Geweke has better generality and more robustness compared to `Geweke_toy`.

4. Heidelberger and Welch Results

Table 4: Result using of the Heidelberger and Welch

Chain	Stationarity test	Starting point	P-value	Halfwidth test	Mean	Halfwidth
First chain	passed	1	0.641	passed	1.004	0.00498
Second chain	failed	NA	5.85e-06	NA	NA	NA

Firstly, based on the results, the running result of the Heidelberger and Welch function is consistent with the actual situation.

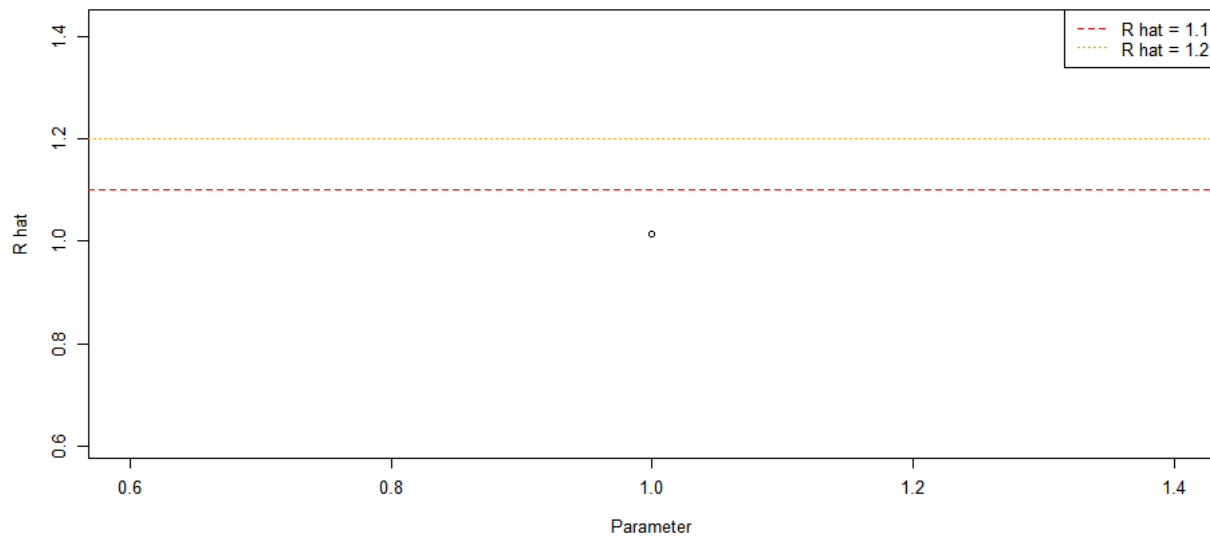
It is found that the first chain converged while the second chain did not converge. The starting point of the first chain is 1, which means that the entire chain can be used for testing to obtain a conclusion of convergence (without the need for burn-in, so that we can obtain more independent samples). Additionally, the first chain also passed the half-width test. In contrast, the second chain obtained a completely opposite conclusion.

Therefore, the performance of Heidelberger and Welch is satisfactory, and there is no need for further improvement.

5. Gelman-Rubin Results

Firstly, the first and second chains are merged and the `gelman_toy` function is run, obtaining a Potential Scale Reduction Factor of 1.014203. The result can be seen in the following figure:

Figure 12: Result of the `gelman_toy` function



As was said earlier in the third chapter, there are two main problems with the toy model: the improved version of PSRF and different target distribution cases, these two problems are solved in the next improved code.

Table 5: Result using of the Gelman-Rubin

Chain	Point est.	Upper C.I.
PSRF	1.065043	1.249431
MPSRF	NULL	NULL

From the table above, the improved PSRF has increased, indicating that the function takes into account the high autocorrelation of the second chain and provides some penalty for it. However,

due to the sufficient convergence of the first chain, the PSRF is still below the threshold of 1.1. The improved function also provides an upper C.I. value that exceeds the threshold of 1.2, indicating that the improved function is more effective than the toy function since the second chain has not converged yet.

However, since this function can only provide a conclusion of rejecting convergence for all chains or not rejecting convergence for all chains, the two chains cannot perfectly demonstrate the value of Gelman-Rubin. Additionally, since both chains have the same target function, MPSRF does not exist. But other examples have been used to demonstrate the reliability of the function.

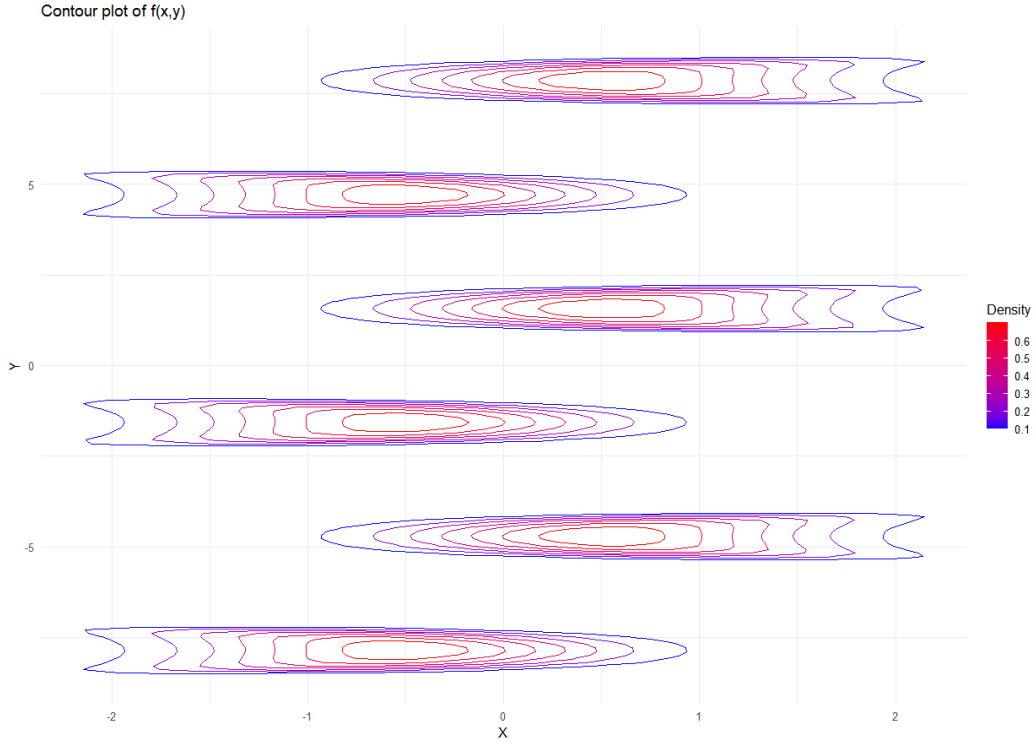
B. A sixmodal target distribution

This example is proposed by Leman et al. (2009) [5] where the target density is as follows:

$$\pi(x, y) \propto e^{-x^2/2} \times e^{-(\csc y^5 - x)^2/2}$$

The contour plot of the target distribution is as follows:

Figure 13: Contour plot of the target distribution in the sixmodal example

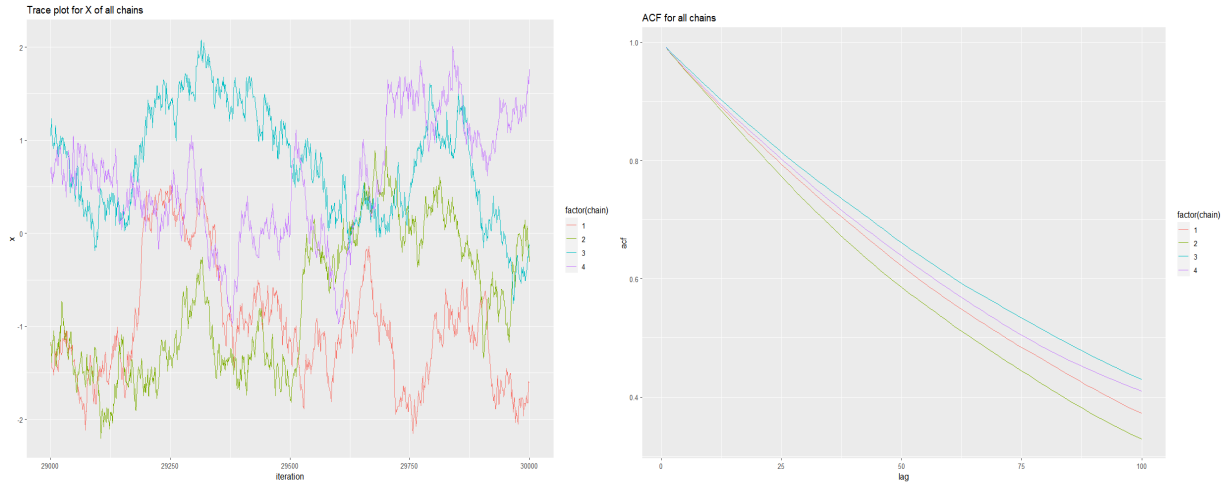


The plots of the joint distributions clearly shows that the target distribution is multimodal in nature. In the following sections, the advantages of using KL divergence's convergence checking tool in multimodal chains will be demonstrated. Since the focus is on Kernel density-based methods, simple results presentations of convergence checking tools such as Gelman-Rubin will be presented.

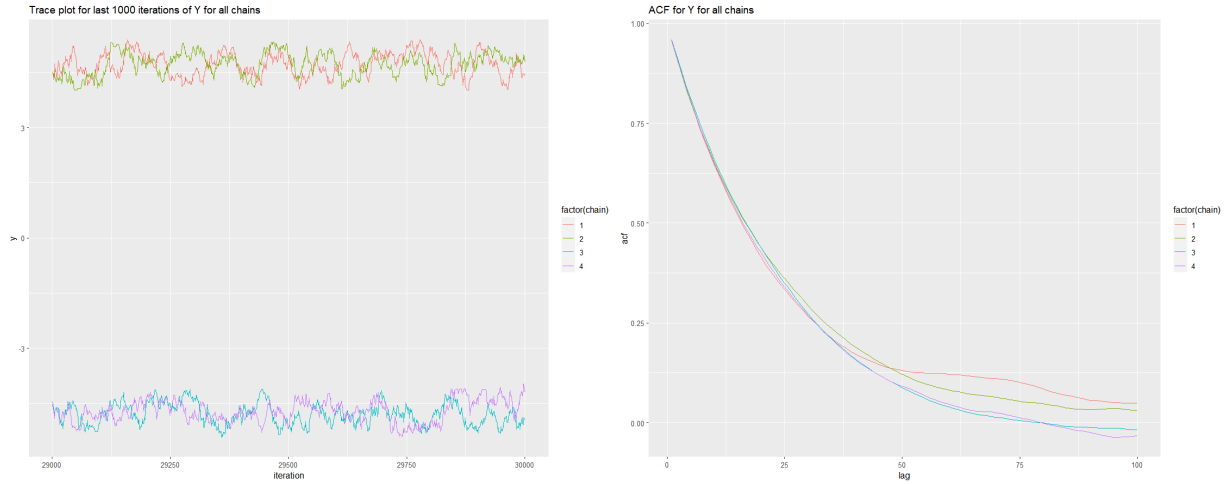
Case 1: In this case, four chains are run, with chains 1 and 2 starting at a particular mode (0, 5), while chains 3 and 4 are started at some other mode (-1, -5). Each of the four chains is run for 30,000 iterations.

1. Graphical Output Results

Figures 14: Trace and autocorrelation function plots of the X marginal of the four chains for the sixmodal example in Case 1



Figures 15: Trace and autocorrelation function plots of the Y marginal of the four chains for the sixmodal example in Case 1.



Trace plots of the last one thousand iterations of the four parallel X and Y marginal chains are given in the left panel of Figures 14 and 15 respectively. Trace plots show the divergence of the Markov chains. High ACF values can also be seen from the autocorrelation plots of the marginal chains in Figures 14 and 15. Therefore, it can be inferred from the figure that the four chains have not yet converged.

2. ESS Results

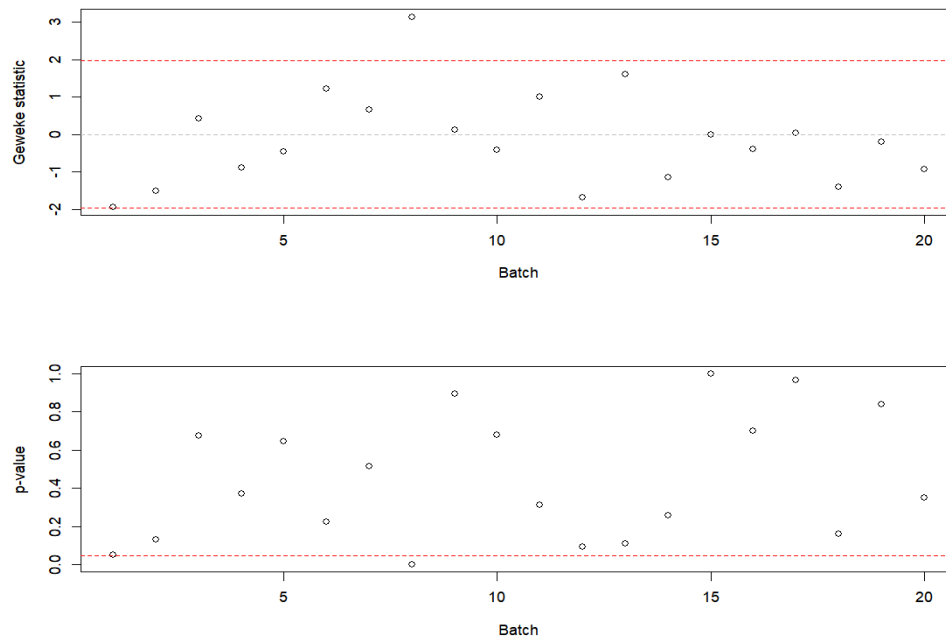
Table 5: ESS Result of the 4 chains in the case 1 using Batch Means

X	Y
544.2388	2584.0141

The minimum ESS for these four two-dimensional chains is 7529, obviously neither X nor Y has reached this value. Therefore, we can draw the conclusion that the chain has not yet converged from the conclusion.

3. Geweke Results

Figures 16: Geweke result using multi-dimensional Geweke in the case 1



Proportion of $p_value < 0.05 = 5\%$

Only one in 20 p-values below 0.05 considers that these four chains (X and Y of the other four chains are similar) have not yet converged, but in fact we know that these four chains have not yet converged, so the Geweke test gives a wrong conclusion.

4. Heidelberger and Welch Results

Table 6: Result using of the Heidelberger and Welch in the case 1

Chain	Stationarity test	Starting point	P-value	Halfwidth test	Mean	Halfwidth
Chain 1 \$ X	passed	3001	0.0533	failed	-0.615	0.1268
Chain 1 \$ Y	passed	1	0.0720	passed	4.713	0.0257
Chain 2 \$ X	passed	1	0.912	failed	-0.653	0.1268
Chain 2 \$ Y	passed	1	0.204	passed	4.724	0.0231
Chain 3 \$ X	passed	1	0.189	failed	0.61	0.124
Chain 3 \$ Y	passed	1	0.382	passed	-4.71	0.024
Chain 4 \$ X	passed	1	0.327	failed	0.573	0.1417
Chain 4 \$ Y	passed	1	0.705	passed	-4.724	0.0242

All chains pass the Heidelberger and Welch (1983) test for stationarity so Heidelberger and Welch (1983)'s diagnostics fail to detect the non-convergence of the chains to the target distribution.

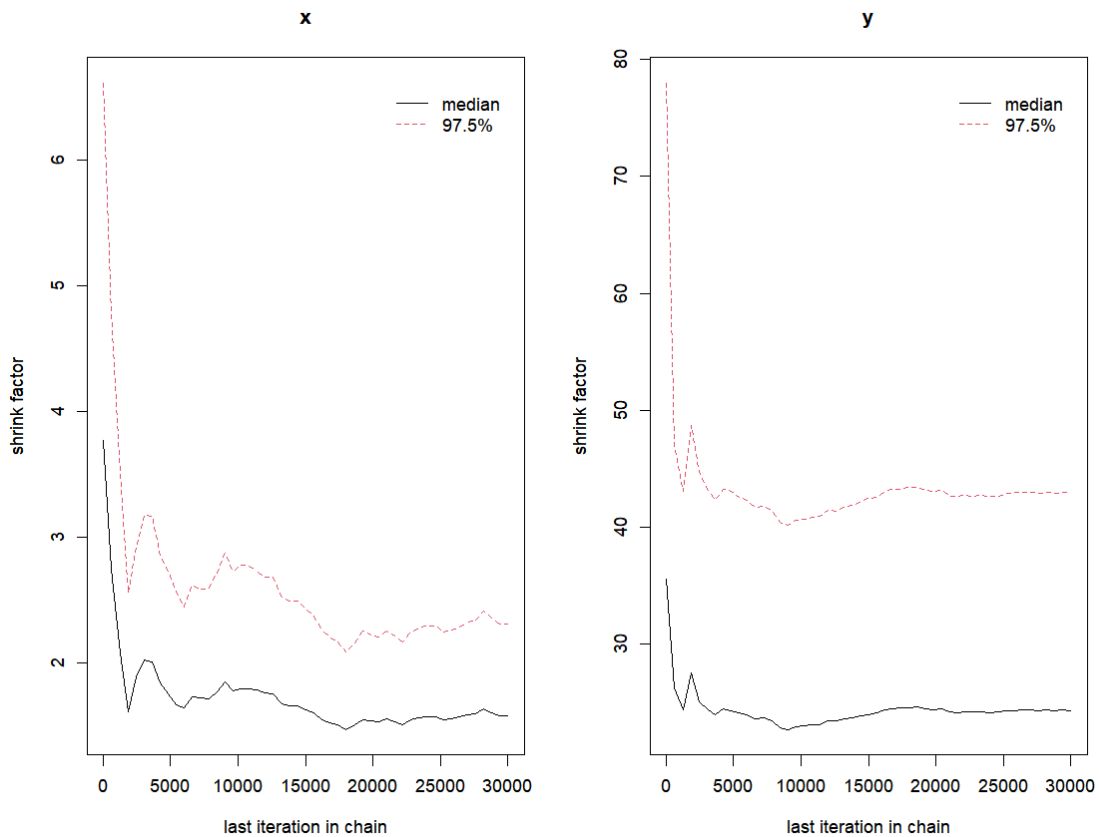
According to the Halfwidth test, the Markov chain can be stopped if the variability of the chain stabilizes with respect to the mean. Figure 15 clearly illustrates that each pair of chains, iterating from the same initial point, consistently explores and stays within a single mode without crossing over to other modes. This remarkable behavior is the key reason behind the stabilization of the chain's variability concerning the mean.

5. Gelman-Rubin Results

Table 7: Result using of the Gelman-Rubin (4 chains) in case 1

Variable	Point est.	Upper C.I.
X	1.58	2.31
Y	24.34	42.95
Multivariate PSRF	21.8	

Figures 17: Iterative \hat{R} plot from four parallel chains in case 1



A large PSRF indicates that the between-chain variance is substantially greater than the within-chain variance, so that longer simulation is needed. It is evident that for Y (the two chain pairs), the between-chain variance is much larger than the within-chain variance, resulting in a significantly higher Potential Scale Reduction Factor (PSRF) for Y compared to X (which exhibits relatively good convergence).

6. KL divergence Tool

Next, Dixit and Roy (2017)'s bivariate KL divergence Tool1 are applied on the joint chain. The maximum symmetric KL divergence among the six pairs is 77.09 significantly larger than the cutoff value 0.06.

Table 7: Result of Tool2 in case 1

Chain	Chain 1	Chain 2	Chain 3	Chain 4
T_2^*	0.8315992	0.828666	0.8256517	0.8265348

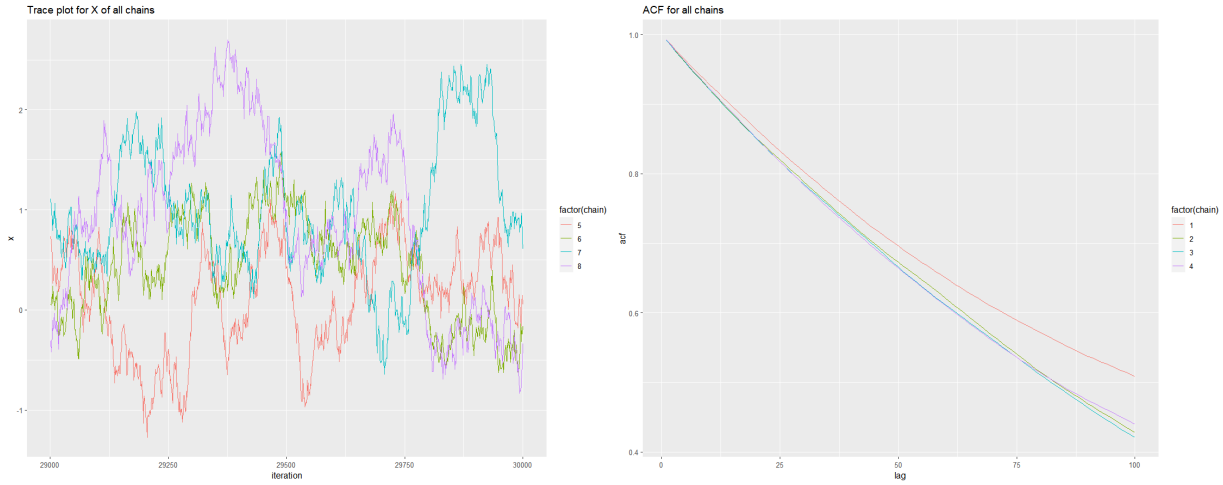
Indeed, Dixit and Roy (2017)'s Tool2 requires just one chain. The computed PSRF values suggest that the four chains exhibit similar behavior, allowing us to select any one of them for further analysis. In this experiment, all four chains were tested to demonstrate the stability of the diagnostic method. All of the T_2^* exceed 0.82 which are significantly greater than zero. This result indicates that all four chains are stuck at the same mode, failing to explore other modes effectively.

So, it seems that even for different starting points (to meet the requirements of other methods such as Gelman-Rubin), there are many tests that are fooled into thinking that the chain has converged - especially for Y (because Y never starts from a modal goes to another modal). Next, we try the second case where the initial points of the four chains are the same.

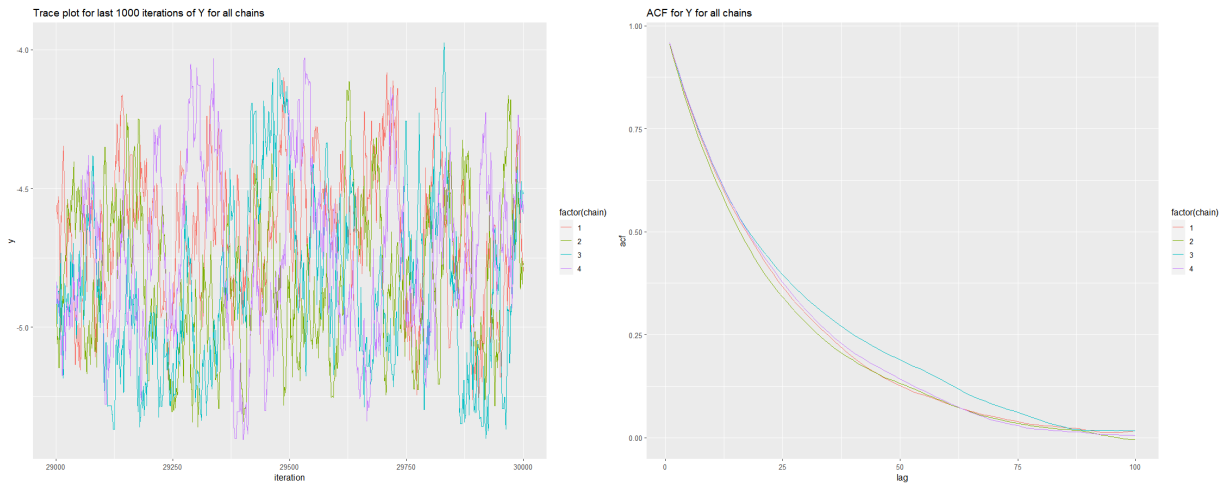
Case 2: In this case four chains are run but all the chains are started at the same local mode (0, - 5). Each of the four chains is run for 30,000 iterations. The trace and autocorrelation plots of the marginal chains are given in Figures 18 and 19. From these plots one may conclude mixing of the Markov chains, although the large autocorrelations result in low ESS for the chains. It is worth noting that the four chains mix very well in the Y direction (compared to case1), it only appears as if Y is constantly fluctuating due to the scale being magnified.

1. Graphical Output Results

Figures 18: Trace and autocorrelation function plots of the X marginal of the four chains for the sixmodal example in Case 2



Figures 19: Trace and autocorrelation function plots of the Y marginal of the four chains for the sixmodal example in Case 2



Obviously, from the unstable trace plot and high ACF value, it can be concluded that the four chains have not yet converged. In fact, this conclusion is doomed. Since the four chains are independently generated using algorithms, a similar conclusion can also be drawn according to case 1. But this time, because the initial point is the same, it will bring wrong conclusions to many convergence diagnostic methods.

2. ESS Results

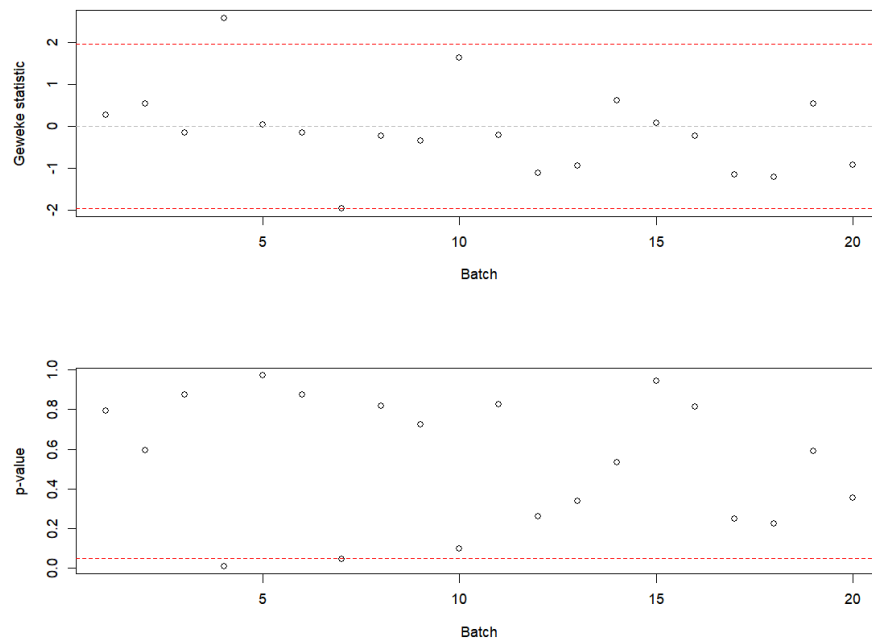
Table 8: ESS Result of the 4 chains in case 2 using Batch Means

X	Y
469.7163	2475.7983

The minimum ESS for these four two-dimensional chains is 7529, obviously neither X nor Y has reached this value. Therefore, we can draw the conclusion that the chain has not yet converged from the conclusion.

3. Geweke Results

Figures 20: Geweke result using multi-dimensional Geweke in case 2



Proportion of $p_value < 0.05 = 10\%$

Only 2 in 20 p-values below 0.05 considers that these four chains (X and Y of the other four chains are similar) have not yet converged, but in fact we know that these four chains have not yet converged, so the Geweke test gives a wrong conclusion.

4. Heidelberger and Welch Results

Table 9: Result using of the Heidelberger and Welch in the case 2

Chain	Stationarity test	Starting point	P-value	Halfwidth test	Mean	Halfwidth
Chain 1 \$ X	passed	1	0.722	failed	0.738	0.1385
Chain 1 \$ Y	passed	1	0.956	passed	-4.708	0.0241
Chain 2 \$ X	passed	1	0.771	failed	0.619	0.1326
Chain 2 \$ Y	passed	1	0.862	passed	-4.719	0.0242
Chain 3 \$ X	passed	1	0.582	failed	0.668	0.1306
Chain 3 \$ Y	passed	9001	0.175	passed	-4.682	0.0291
Chain 4 \$ X	passed	1	0.814	failed	0.581	0.1543
Chain 4 \$ Y	passed	1	0.867	passed	-4.738	0.0255

All chains pass the Heidelberger and Welch (1983) test for stationarity so Heidelberger and Welch (1983)'s diagnostics fail to detect the non-convergence of the chains to the target distribution.

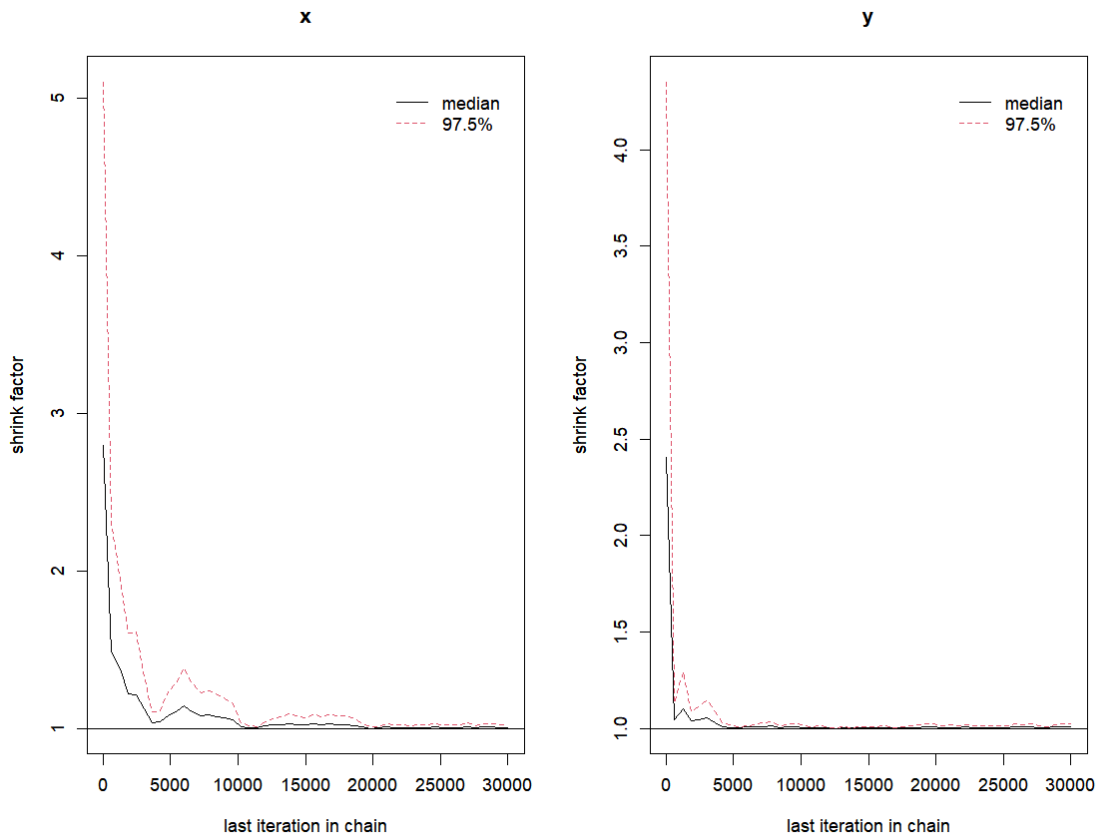
It is not surprising that all of the chains passed the bandwidth test in Y direction. Upon observation, it can be noticed that the majority of the P-values from the stability test of these four chains are larger than those in case1. This indicates that due to the same initial points, the test has a tendency to make judgments in a more erroneous direction.

5. Gelman-Rubin Results

Table 10: Result using of the Gelman-Rubin (4 chains) in case 2

Variable	Point est.	Upper C.I.
X	1.01	1.02
Y	1.01	1.02
Multivariate PSRF	1.01	

Figures 21: Iterative \hat{R} plot from four parallel chains in case 2



Unlike case1, the Gelman-Rubin test results for case2 provide completely incorrect conclusions. This is due to the fact that all four chains in the Y direction are consistently running within the same mode, resulting in a very small between-chain variance. This makes the Gelman-Rubin statistic unable to compute the correct result.

6. KL divergence Tool

Dixit and Roy (2017)'s bivariate KL divergence Tool1 are applied on the joint chain. The maximum symmetric KL divergence among the six pairs is 0.07. Tool1 subtly rejected the null hypothesis and concluded that the four chains have not yet converged. This is because Tool1 compares the differences in adaptive kernel densities between the chains, and in the Y direction, the chains did not separate into other modes. Additionally, the X direction was approximated as stationary. As a result, the conclusion obtained is not entirely reliable, but at least correct.

Table 7: Result of Tool 2 in case 2

Chain	Chain 5	Chain 6	Chain 7	Chain 8
T_2^*	0.8263013	0.8274083	0.8283718	0.830165

In contrast, Tool2 provided the correct conclusion. Similar to case1, the adaptive kernel densities of the chains in case 2 exhibit significant differences from the target distribution, indicating that they have not yet converged.

From a practical perspective, it is important to exercise caution when relying solely on empirical convergence diagnostic tools, especially in situations where the presence of multiple modes is suspected. Empirical diagnostics cannot precisely detect convergence. Therefore, we should consider trying updated sampling methods, such as thinning, which involves discarding all but every k th observation and is often used by MCMC practitioners to reduce high autocorrelations present in the Markov chain samples. Alternatively, we can explore new sampling methods like the multiset sampler [5] to ensure that the chain can traverse all modes. A potential future research problem is to theoretically verify the convergence (towards zero) of Dixit and Roy's (2017) statistics based on the KL divergence [4]. Another possible research question is to develop theoretically sound and computationally efficient MCMC convergence diagnostic methods for ultra-high-dimensional settings.

VI. Acknowledgements

First of all, I would like to express our gratitude to Professor Guanyang Wang for allowing me to conduct this research and introducing us to many relevant papers. Throughout the research process, Professor Wang provided me with continuous guidance and helped us deepen my understanding of MCMC convergence diagnostics.

Secondly, I want to thank Professor Jack Mardekian for his detailed lectures on Bayesian analysis and MCMC convergence diagnostics. We had multiple discussions, and Professor Mardekian provided numerous references to authors who have made significant contributions to this field.

Finally, I would like to express our gratitude to the authors of the R packages “coda” and “mcmcse” and the authors of the references cited in this research. Their unique perspectives have been crucial to this study. They have made the convergence diagnostics of MCMC possible and introduced innovative techniques, such as spectral variance estimators, to improve existing diagnostic methods.

VII. Reference

- [1] [SAS/STAT® 14.2 User’s Guide Introduction to Bayesian Analysis Procedures: Chapter 7](#)
- [2] [BATCH MEANS AND SPECTRAL VARIANCE ESTIMATORS IN MCMC](#)
(James M. Flegal and Galin L. Jones, 2010)
- [3] [Convergence diagnostics for Markov chain Monte Carlo](#)
(Vivekananda Roy, 2011)
- [4] [MCMC diagnostics for higher dimensions using Kullback Leibler divergence](#)
(Dixit, A. and Roy, V, 2017)
- [5] [The multiset sampler.](#)
(Leman, S. C., Chen, Y. and Lavine, M., 2009)

VIII. Appendix

<https://github.com/QiruPan/Convergence-diagnostics-for-MCMC>