

Statistical inference for SIS stochastic model: MLE and MH Tutorial in R

Dr Panayiota Touloupou

July 15, 2025

Abstract

Welcome to this R tutorial on fitting the Susceptible-Infected-Susceptible (SIS) model, considering it as a stochastic epidemic process based on the Binomial sampling scheme that have seen in our lectures. In this tutorial we will use two statistical methods: Maximum Likelihood Estimation (MLE) and Metropolis-Hastings (MH). This tutorial is designed for students new to R and will guide you through the process step-by-step.

1 Getting started with R (for students new to R)

R is a programming language and free software environment for statistical computing and graphics.

1.1 Install R and RStudio

We assume you don't yet have R on your computer. To get started, you will need to install two pieces of software:

- R, the actual programming language.
 - R has a homepage, <http://r-project.org/>, but the software itself is located for download on the CRAN, which you can find at <http://cran.r-project.org/>.
 - Chose your operating system (Macintosh, Windows, Linux or Unix) and select to install the most recent version.
- RStudio (you must have R installed to use RStudio), for working with R.
 - RStudio is simply an interface used to interact with R.
 - Install RStudio from <http://rstudio.org>.

1.2 RStudio interface (Brief Overview)

The RStudio screen is divided into 4 windows:

- Bottom left: **console** window (also called command window). This is the most important window since it is the active bit, where you can communicate directly with R. Here you can type simple commands after the “>” prompt and R will then execute your command. It is also where your non-graphical results appear.
- Top left: **editor** window (also called script window). Collections of commands (scripts) that can be edited and saved. Each script is basically just a text file where you can write longer sections of code. Also this area is a place where you can view datasets.
- Top right: environment and history window. In the environment window (also called workspace window) you can see what is in R’s memory, i.e. a list of data and variables, which contain values that R has been told to save from previous commands. The history window shows a history of commands that have been typed on the console.
- Bottom right: files/plots/packages/help window. Here you can do everything else - open files, look at and manipulate plots, install and load packages or use the help function and read help files.

1.3 Basic R syntax

Let’s try some basic commands. Type these into the **console** and press Enter (or into **editor** window and press the button Run).

```

1 # This is a comment. R ignores anything after a '#' on a line.
2
3 # Assigning values to variables
4 x <- 10
5 y = 5 # '=' also works for assignment, but '<-' is preferred in R
6 print(x + y)
7
8 # Vectors (ordered collections of values of the same type)
9 my_vector <- c(1, 2, 3, 4, 5) # 'c()' combines values into a vector
10 print(my_vector)
11 print(my_vector * 2)
12 print(sum(my_vector))
13
14 # Functions
15 my_function <- function(a, b) {
16   result <- a * b
17   return(result)
18 }
19 print(my_function(x, y))
20
21 # Data frames (like tables or spreadsheets)
22 my_data <- data.frame(
23   Name = c("Alice", "Bob", "Charlie"),
24   Age = c(25, 30, 22),
25   City = c("NY", "LA", "Chicago")
26 )
27 print(my_data)
28 print(my_data$Age) # Accessing a column by name
29 print(my_data[1, ]) # Accessing the first row

```

1.4 Installing and loading packages

R's power comes from its vast collection of packages. Packages are collections of functions, data, and compiled code in a well-defined format, designed to extend R's capabilities. Think of them as add-ons that provide specialized tools. With the standard R installation, most common packages are installed. At some point in the future if you want to use some functions that are not in the ordinary R installation, there is a good chance that there is a package that will fulfill your needs.

```
1 # Install packages (only need to do this once per R installation)
2 install.packages("ggplot2") # For nice plots
3 install.packages("coda")    # For MCMC diagnostics
4
5 # Load packages (need to do this every time you start a new R session
6 # and want to use them)
7 library(ggplot2)
8 library(coda)
```

1.5 Basic plotting

Visualizing data is crucial in statistics and modelling. R has built-in plotting functions, and packages like ggplot2 offer more advanced and aesthetically pleasing options.

```
1 # Simple plot using R's base plotting system
2 # 'plot()' is a versatile function, 'type = "l"' means draw lines.
3 # 'main', 'xlab', 'ylab' are for setting the title and axis labels.
4 plot(my_vector, type = "l", main = "My Simple Line Plot",
5      xlab = "Index", ylab = "Value")
```

2 The stochastic SIS model

The SIS model is an epidemic model where individuals can be in one of two states: Susceptible (S) or Infected (I). Infected individuals recover and immediately become susceptible again (no immunity).

Unlike deterministic models that describe average behavior with continuous equations, stochastic models incorporate randomness, reflecting that events (like infections or recoveries) happen probabilistically rather than with fixed rates. This is especially important in smaller populations where random fluctuations can significantly impact the epidemic trajectory.

2.1 Model events and probabilities

In this specific stochastic SIS model, we consider two types of events occurring at each discrete time step (e.g., daily, weekly):

1. **Number of new infections (X_t):** This is the number of susceptible individuals who become infected during the current time step.

- The number of new infections X_t is drawn from a Binomial distribution: $X_t \sim \text{Binomial}(S_t, p_{\text{inf}})$.
- The probability of infection for a single susceptible individual, p_{inf} , is given by: $1 - e^{-\frac{\beta I_t}{N}}$.
- A higher β or more I_t leads to a higher chance of infection.

2. **Number of new recoveries (Y_t):** This is the number of infected individuals who recover and return to the susceptible state during the current time step.

- The number of new recoveries Y_t is drawn from a Binomial distribution: $Y_t \sim \text{Binomial}(I_t, p_{\text{rec}})$.
- The probability of recovery for a single infected individual, p_{rec} , is given by: $1 - e^{-\gamma}$.
- A higher γ means individuals recover faster.

Where:

- S_t and I_t are the number of susceptible and infected individuals at time t .
- $N = S_t + I_t$ is the total population size (assumed constant throughout the simulation).
- β is the transmission rate (a positive value).
- γ is the recovery rate (a positive value).

2.2 Implementing the stochastic SIS model in R

Now, let's translate these probabilistic rules into an R function. This function will simulate the SIS model step-by-step over a given time period.

Note that for the SIS model, calculating the likelihood based only on the total number of susceptible or infected individuals at each time point is not sufficient. We need additional information. This is because individuals can recover and return to the susceptible state. If we only observe S_t or I_t , we cannot tell whether a change in S (or I) is due to a new infection or a recovery.

For example, if S decreases, it could be due to an infection; but if I also decreases while S increases, it could be a recovery. Without explicitly tracking the number of new infections (`new_inf`) and recoveries (`new_rec`) at each time step, the likelihood function cannot correctly attribute the changes to the underlying processes. Therefore, we need to modify our simulation function (compared to the one we used in lectures) to return `new_inf_history` and `new_rec_history` to provide this crucial information for the likelihood calculation.

```
1 # Define the Stochastic SIS model simulation function
2 # This function takes:
3 # N: Total population size
4 # beta: Transmission rate
```

```

5 # gamma: Recovery rate
6 # I0: Initial number of infected individuals
7 # T: Total number of time steps to simulate
8
9 simulate_SIS <- function(N, beta, gamma, I0, T) {
10   # Initialize vectors to store the number of Susceptible (S)
11   # and Infected (I) individuals
12   S <- numeric(T)
13   I <- numeric(T)
14
15   # Initialize vectors to store the *number of new infections* and
16   # *new recoveries* that occur at each time step. These will be
17   # crucial for calculating the likelihood function.
18   new_inf_history <- numeric(T)
19   new_rec_history <- numeric(T)
20
21   # Set initial conditions for S and I at time t=0
22   S[1] <- N - I0
23   I[1] <- I0
24   new_inf_history[1] <- 0 # No new events at time 0
25   new_rec_history[1] <- 0 # No new events at time 0
26
27   # Loop through each time step from t=1 to T-1
28   for (t in 1:(T-1)) {
29     # Calculate the probability of infection for a susceptible
30     # individual based on the current number of infected individuals
31     p_inf <- 1 - exp(-beta*I[t]/N)
32
33     # Calculate the probability of recovery for an infected individual
34     p_rec <- 1 - exp(-gamma)
35
36     # Simulate number of events in this time step using binomial
37     # distribution
38     # rbinom(n, size, prob) generates n random values from a binomial
39     # distribution size is the number of trials, prob is the
40     # probability of success on each trial
41     new_inf <- rbinom(1, S[t], p_inf)
42     new_rec <- rbinom(1, I[t], p_rec)
43
44     # Store the simulated number of events for this time step
45     new_inf_history[t+1] <- new_inf
46     new_rec_history[t+1] <- new_rec
47
48     # Update the number of Susceptible and Infected individuals for the
49     # next time step (t+1)
50     # S decreases by new infections and increases by new recoveries
51     S[t + 1] <- S[t] - new_inf + new_rec
52     # I increases by new infections and decreases by new recoveries
53     I[t + 1] <- I[t] + new_inf - new_rec
54   }
55   # Return a data frame
56   return(data.frame(time = 0:(T-1), S = S, I = I,
57     new_inf_obs = new_inf_history, new_rec_obs = new_rec_history))
58 }

```

2.3 Simulating the stochastic SIS model

Let's set some parameters and initial conditions and simulate the model over time.

```
1 # Set a seed for reproducibility of stochastic results.
2 # This ensures that if you run the code multiple times with the
3 # same seed, you will get the exact same sequence of random
4 # numbers, making your results reproducible.
5 set.seed(3)
6
7 # Simulate data using the 'simulate_SIS' function.
8 # N: Total population size
9 # beta: True transmission rate
10 # gamma: True recovery rate
11 # I0: Initial number of infected individuals
12 # T: Total number of time steps (e.g., days) for the simulation
13 SISdata <- simulate_SIS(N = 1000, beta = 0.3, gamma = 0.1, I0 = 10,
14   T = 150)
15
16 # View the first few rows of the simulation output.
17 head(SISdata)
18
19 # Plotting the simulation results
20 plot(SISdata$time, SISdata$I, type = "l", col = "blue",
21   ylab = "Number of individuals", xlab = "Time", ylim = c(0, 1000))
22 lines(SISdata$time, SISdata$S, col = "green")
23
24 # Add legend
25 legend("topright", legend = c("Infectious", "Susceptible"),
26   lty = c(1, 1), col = c("blue", "green"))
```

2.4 Student activity:

- **Observe variability:** Run the simulation code several times. How do the individual simulation paths differ, even with the same parameters? This is the essence of stochasticity!
- **Experiment with parameters:** Change `beta` and `gamma` values in the simulation code. How do they influence the average behavior (e.g., the peak number of infected individuals, the endemic equilibrium) and the variability of the epidemic?
- **Change initial conditions:** Vary `I0` (initial infected). Does it affect the long-term behavior or the initial spread? What happens if you start with very few infected individuals (e.g., `I0 = 1`)?
- **Change population size N:** Does changing `N` (while keeping `I0` proportional or constant) alter the magnitude of stochastic fluctuations? **Comment:** In smaller populations, stochastic effects are often more pronounced. Why do you think this is the case? Consider what happens when counts become very small.

3 Maximum Likelihood Estimation (MLE) for stochastic SIS

Maximum Likelihood Estimation (MLE) is a widely used statistical method for estimating the parameters of a model. The core idea is to find the parameter values that make the observed data most “likely” to have occurred under the assumed model, i.e. the values that maximize the likelihood function (which is the probability of observing the given data under the model).

To test our MLE approach, we first need some “observed” data. In a real-world scenario, this would be actual epidemic data (e.g., daily reported cases). For this tutorial, we will generate synthetic data from our stochastic SIS model using known “true” parameters. We then pretend we don’t know these true parameters and use MLE to try and recover them. The stochastic nature of the model itself will introduce the “noise” or variability that we would typically see in real data.

3.1 Generating Synthetic Data

To test our MLE approach, we’ll first generate some “observed” data from our stochastic SIS model. The stochastic nature itself will provide the “noise”.

```
1 # True parameters (these are what we want to estimate)
2 true_beta <- 0.3
3 true_gamma <- 0.1
4
5 # Initial conditions and time for data generation
6 N_data <- 1000
7 IO_data <- 5
8 SO_data <- N_data - IO_data
9 T_data <- 150
10
11 # Simulate the true SIS stochastic model to generate observed data
12 set.seed(123) # For reproducibility
13 observed_data <- simulate_SIS(N = N_data, beta = true_beta,
14   gamma = true_gamma, IO = IO_data, T = T_data)
```

3.2 Defining the log likelihood function for SIS

The likelihood of observing a specific number of new infections (X_t) at time t is given by a Binomial distribution with parameters S_t and p_{inf} (probability of infection). Similarly, for new recoveries (Y_t), it’s a Binomial distribution with parameters I_t and p_{rec} (probability of recovery).

The total log-likelihood is the sum of the log-likelihood of all observed events across all time steps.

```
1 # Define the log-likelihood Function for the SIS model
2 # This function takes:
3 # params: A vector of parameters (e.g., c(beta, gamma))
```

```

4 # data: The observed data.
5
6 loglik_SIS <- function(params, data) {
7   beta = params[1]
8   gamma = params[2]
9
10  # Ensure parameters are valid (e.g., positive).
11  # If any parameter is non-positive, return -Inf for the
12  # log-likelihood.
13  # This tells the optimizer that these parameters are invalid.
14  if (beta <= 0 || gamma <= 0) {
15    return(-Inf)
16  }
17
18  loglik <- 0 # Initialize the total log-likelihood to zero
19  # Calculate the total population size from the initial states
20  N <- data$S[1] + data$I[1]
21
22  # Loop through each time step in the observed data
23  for (t in 1:(nrow(data) - 1)) {
24
25    # Get the number of Susceptible and Infected individuals
26    # at time t and t+1
27    St <- data$S[t]
28    It <- data$I[t]
29
30    # Get the observed number of new infections and new recoveries
31    # that occurred between time t and time t+1. These come directly
32    # from our observed data.
33    new_inf <- data$new_inf_obs[t + 1] # Events observed at t+1,
34                                         # from state at t
35    new_rec <- data$new_rec_obs[t + 1] # Events observed at t+1, from
36                                         # state at t
37
38    # Calculate probabilities
39    p_inf <- 1 - exp(-beta*It/N)
40    p_rec <- 1 - exp(-gamma)
41
42    # Avoid probabilities of exactly 0 or 1, which can cause log(0) or
43    # log(1-1) issues in the dbinom function.
44    # We clip them to be slightly away from the boundaries.
45    p_inf <- min(max(p_inf, 1e-10), 1 - 1e-10)
46    p_rec <- min(max(p_rec, 1e-10), 1 - 1e-10)
47
48    # Add log-likelihood contribution for this step
49    # dbinom(x, size, prob, log = TRUE) gives the log of the
50    # binomial probability mass function.
51    loglik <- loglik +
52      dbinom(new_inf, St, p_inf, log = TRUE) +
53      dbinom(new_rec, It, p_rec, log = TRUE)
54  }
55  return(loglik) # Return the total log-likelihood
56 }

```


3.3 Performing MLE using optim

R's `optim` function is a powerful general-purpose optimization tool. It attempts to find the set of parameters that minimize a given function. In our case, it will maximize the log-likelihood by using `control = list(fnscale = -1)`, thereby finding the parameters that maximize the likelihood of our observed data. Write `?optim` in console to find more details.

```
1 ?optim
2
3 # Set Initial guesses for parameters
4 # These are the starting points for the optimizer's search.
5 initial_params <- c(beta = 0.2, gamma = 0.05)
6
7 # Run the optimization using 'optim()'.
8 # 'par': The initial guesses for the parameters.
9 # 'fn': The function to be optimized (our log-likelihood function).
10 # 'data': Our simulated 'SISdata' that contains the observed S, I,
11 # new_inf_obs, new_rec_obs.
12 # 'control = list(fnscale = -1)': This is crucial! It tells 'optim'
13 # to maximize 'fn' instead of minimizing it.
14 mle_result <- optim(par = initial_params,
15                     fn = loglik_SIS, data = observed_data,
16                     control = list(fnscale = -1))
17 print("MLE Results (initial run):")
18 print(mle_result)
19 # Look for 'convergence = 0' which indicates successful convergence.
20 # 'par' will give you the estimated parameter values.
21 # 'value' is the maximum log-likelihood found.
22
23 # Extract estimated parameters
24 estimated_beta_mle <- mle_result$par[1]
25 estimated_gamma_mle <- mle_result$par[2]
26
27 cat("\nEstimated Beta (MLE):", estimated_beta_mle, "\n")
28 cat("Estimated Gamma (MLE):", estimated_gamma_mle, "\n")
29 cat("True Beta:", true_beta, "\n")
30 cat("True Gamma:", true_gamma, "\n")
```

3.4 Repeated MLE estimation and variability

In the previous section, you performed MLE once. However, because our data is generated from a stochastic model, each time you simulate data, you'll get a slightly different trajectory. This means that if you run the MLE process on different simulated datasets, you'll get slightly different parameter estimates.

Let's explore this variability by repeating the data simulation and MLE estimation process multiple times.

```
1 # Define the number of times to repeat the MLE process
2 num_repetitions <- 25
3
4 # Create empty vectors to store the estimated beta and gamma from
```

```

5 # each repetition
6 estimated_betas <- numeric(num_repetitions)
7 estimated_gammas <- numeric(num_repetitions)
8
9 # Set the fixed parameters for data generation for this exercise
10 fixed_N <- 1000
11 fixed_I0 <- 10
12 fixed_T <- 150
13 fixed_true_beta <- 0.3
14 fixed_true_gamma <- 0.1
15
16 # Loop to repeat the simulation and MLE estimation
17 for (i in 1:num_repetitions) {
18   # Set a unique seed for each repetition to get different stochastic
19   # data realizations
20   set.seed(100 + i)
21
22   # Simulate new SIS data for this repetition
23   current_SISdata <- simulate_SIS(N = fixed_N, beta = fixed_true_beta,
24     gamma = fixed_true_gamma, I0 = fixed_I0, T = fixed_T)
25
26   # Run MLE for the current simulated data
27   # Use the same initial guesses and bounds as before
28   current_mle_result <- optim(par = initial_params,
29     fn = loglik_SIS, data = current_SISdata,
30     control = list(fnscale = -1))
31
32   # Store the estimated parameters
33   estimated_betas[i] <- current_mle_result$par[1]
34   estimated_gammas[i] <- current_mle_result$par[2]
35 }
36
37 # Now, visualize the distribution of your 25 estimates using boxplots.
38 # Boxplots are great for showing the median, quartiles, and outliers
39 # of a distribution.
40 # Set up a 1x2 plotting layout
41 par(mfrow = c(1, 2), mar = c(4, 4, 2, 1))
42
43 # Boxplot for estimated Beta values
44 boxplot(estimated_betas, main = "Distribution of Estimated Beta",
45   ylab = "Estimated Beta", col = "lightblue", border = "black")
46 # Add a line for the true beta value
47 abline(h = fixed_true_beta, col = "red", lwd = 2, lty = 2)
48
49 # Boxplot for estimated Gamma values
50 boxplot(estimated_gammas, main = "Distribution of Estimated Gamma",
51   ylab = "Estimated Gamma", col = "lightgreen", border = "black")
52 # Add a line for the true gamma value
53 abline(h = fixed_true_gamma, col = "red", lwd = 2, lty = 2)
54
55 # Reset plotting layout
56 par(mfrow = c(1, 1))
57

```

3.5 Student Activities for MLE:

- Execute the “Repeated MLE Estimation and Boxplots” code block.

- Observe the boxplots:
 - Where is the median (the line inside the box) of your boxplots located relative to the true parameter values (red dashed line)?
 - How wide are the boxes? This indicates the interquartile range (middle 50% of your estimates).
 - Are there any “outliers” (points beyond the whiskers)?
 - **Comment:** What does the spread of the boxplots tell you about the precision of your MLE estimates for this stochastic model? How does this relate to the concept of statistical uncertainty?
- Vary `num_repetitions`: Try increasing `num_repetitions`: (e.g., to 100 or 500). Does the distribution of estimates change? Does it become more concentrated around the true value?
- Vary `N` (Population Size): Change `N` (e.g., to 100, 1000, 10000) when generating the data, then re-run the “Repeated MLE Estimation and Boxplots” section.

Comment: How does increasing the population size (`N`) affect the variability (spread) of your boxplots? Does more data generally lead to more precise estimates for stochastic models?
- Vary `I0` (Initial Infected): Change `I0` (e.g., to 1, 50, 100) when generating the data, then re-run the “Repeated MLE Estimation and Boxplots” section.

Comment: How does the initial number of infected individuals affect the estimation? Is it harder to estimate parameters if `I0` is very small (e.g., leading to early extinction of the epidemic in some stochastic runs)?

4 Metropolis-Hastings (MH) for stochastic SIS

While MLE provides a single “best” point estimate for parameters, it doesn’t directly tell us about the uncertainty around those estimates. Metropolis-Hastings (MH) is a Markov Chain Monte Carlo (MCMC) algorithm that allows us to sample from complex probability distributions, particularly the posterior distribution in Bayesian inference. Instead of a single estimate, MH gives us a distribution of plausible parameter values, which is incredibly valuable for understanding uncertainty.

4.1 Introduction to Bayesian Inference (Recap)

- **Prior distribution** $p(\theta)$: Our beliefs about the parameters ($\theta = (\beta, \gamma)$) before seeing the data.
- **Likelihood function** $p(Data|\theta)$: The probability of observing the data given the parameters (same as in MLE, but now based on the stochastic model).
- **Posterior distribution** $p(\theta|Data)$: Our updated beliefs about the parameters after seeing the data. It’s proportional to Prior \times Likelihood.

The MH algorithm constructs a Markov chain whose stationary distribution is the target posterior distribution. By running the chain for a long time, the samples generated will approximate samples from the posterior.

4.2 Implementing the MH Algorithm

The core idea of MH is to propose a new set of parameters, calculate an acceptance ratio, and either accept or reject the new parameters.

```
1 # Log-Prior Function (assuming exponential priors)
2 log_prior <- function(beta, gamma, lambda_beta = 1, lambda_gamma = 1) {
3   dexp(beta, rate = lambda_beta, log = TRUE) +
4   dexp(gamma, rate = lambda_gamma, log = TRUE)
5 }
6
7
8 # Metropolis-Hastings Algorithm
9 MH_sampler_SIS <- function(data, n_iter, beta_init, gamma_init,
10                             lambda_beta, lambda_gamma, proposal_sd) {
11
12   # Initialize current parameters
13   beta <- beta_init
14   gamma <- gamma_init
15   samples <- matrix(NA, n_iter, 2)
16
17   # Counter for accepted proposals
18   accepted_count <- 0
19
20   for (i in 1:n_iter) {
21     beta_prop <- rnorm(1, beta, proposal_sd)
22     gamma_prop <- rnorm(1, gamma, proposal_sd)
23
24     if (beta_prop > 0 && gamma_prop > 0) {
25       loglik_curr <- loglik_SIS(c(beta, gamma), data)
26       loglik_prop <- loglik_SIS(c(beta_prop, gamma_prop), data)
27
28       logprior_curr <- log_prior(beta, gamma, lambda_beta,
29                                   lambda_gamma)
30       logprior_prop <- log_prior(beta_prop, gamma_prop, lambda_beta,
31                                   lambda_gamma)
32
33       log_accept_ratio <- (loglik_prop + logprior_prop) -
34         (loglik_curr + logprior_curr)
35
36       if (log(runif(1)) < log_accept_ratio) {
37         # accept
38         beta <- beta_prop
39         gamma <- gamma_prop
40         accepted_count <- accepted_count + 1
41       }
42     }
43
44     # Store the current (accepted or re-used) parameters in the chain
45     samples[i, ] <- c(beta, gamma)
46   }
47
48   cat("Acceptance Rate:", accepted_count / n_iter, "\n")
49
50   colnames(samples) <- c("beta", "gamma")
51   return(as.data.frame(samples))
52 }
```

4.3 Running the MH Sampler

Now, let's run the Metropolis-Hastings algorithm to generate samples from the posterior distribution of our β and γ parameters.

```
1 # True parameters (these are what we want to estimate)
2 true_beta <- 0.3
3 true_gamma <- 0.1
4
5 # Initial conditions and time for data generation
6 N_data <- 1000
7 IO_data <- 5
8 SO_data <- N_data - IO_data
9 T_data <- 150
10
11 # Simulate the true SIS stochastic model to generate observed data
12 set.seed(123) # For reproducibility
13 observed_data <- simulate_SIS(N = N_data, beta = true_beta,
14                               gamma = true_gamma, IO = IO_data, T = T_data)
15
16
17 # Run the MH sampler
18 mh_chain <- MH_sampler_SIS(data = observed_data,
19                             n_iter = 10000, beta_init = 0.2, gamma_init = 0.2,
20                             lambda_beta = 1, lambda_gamma = 1, proposal_sd = 0.01)
```

4.4 Analyzing MCMC Output

It's crucial to analyze the MCMC chain to ensure it has converged and is sampling effectively.

```
1 # Discard burn-in period (e.g., first 10% of iterations)
2 burn_in <- n_iter * 0.1
3 mh_chain_post_burnin <- mh_chain[-(1:burn_in), ]
4
5 #install.packages("coda")
6 library("coda")
7
8 # Convert to 'mcmc' object for coda package functions
9 mcmc_object <- as.mcmc(mh_chain_post_burnin)
10
11 # 1. Trace Plots: Show the values of parameters over iterations.
12 # Should look like "fuzzy caterpillars".
13 plot(mcmc_object)
14
15 # 2. Summary Statistics: Mean, median, credible intervals
16 # (similar to confidence intervals).
17 summary(mcmc_object)
18
19 # 3. Autocorrelation Plots: Show correlation between
20 # samples at different lags. Should drop quickly.
21 autocorr.plot(mcmc_object)
22
23 # 4. Effective Sample Size (ESS): How many independent samples you
```

```

24 # effectively have. Higher is better.
25 effectiveSize(mcmc_object)
26
27 # Compare MH estimates (mean of posterior) with MLE and true values
28 cat("\nMH Posterior Means:\n")
29 print(colMeans(mh_chain_post_burnin))
30 cat("\nMLE Estimates:\n")
31 print(mle_result$par)
32 cat("\nTrue Parameters:\n")
33 print(c(beta = true_beta, gamma = true_gamma))

```

4.5 Student Activities for MH:

1. Multiple Chains:

- It's good practice to run multiple independent MCMC chains from different starting points to check for convergence.
- Run the `MH_sampler_SIS` function multiple times (e.g., 2-3 times), each with a different `set.seed()` and different initial parameters, `beta_init` and `gamma_init`.
- Combine the chains (e.g., using `rbind` or `mcmc.list` from `coda`) and plot them together using `plot(as.mcmc(combined_chain))`.
- **Comment:** Do the chains converge to the same region of the parameter space? What does this tell you about the reliability of your MCMC results and whether the algorithm has found the true posterior?

2. Different proposal distributions:

- Experiment with `proposal_sd` values (e.g., `c(0.001, 0.0005)` for smaller steps, or `c(0.05, 0.02)` for larger steps).
- Observe the “Acceptance Rate” printed by the function.
- **Comment:** How does the `proposal_sd` affect the acceptance rate? What happens to the trace plots and autocorrelation plots if the acceptance rate is very low or very high? (Aim for an acceptance rate between 20-50%).

3. Comparing MLE and MH:

- Compare the point estimates from MLE (`mle_result$par`) with the mean/median of the posterior distributions from MH (`colMeans(mh_chain_post_burnin)` or `summary(mcmc_object)`).
- **Comment:** Are they similar? What are the advantages of having a full posterior distribution (from MH) compared to a single point estimate (from MLE) when dealing with stochastic models and uncertainty? (Hint: Think about how each method quantifies or represents uncertainty).

Experiment with the code and the suggested activities to deepen your understanding!!!