

STAT 243 Final Project: ARS R Package

Jeffrey Kuo, Hsiang-Chuan Sha, Tessa Weiss

2022-12-14

Package Installation

Introduction

This project aims to create an R package called **ars** that simulates the adaptive rejection sampling algorithm described in the paper Adaptive Rejection Sampling for Gibbs Sampling by W.R. Gilks and P Wild. In this paper, Gilks and Wild propose a method for rejection sampling, which is valid for any univariate, log-concave probability density function. As part of this project, we created a main file called **ars.R** which contained the main **ars()** function used to calculate and produce the samples, with helper functions to complete each step in the algorithm defined in **ars_functions.R**.

Approach

The main function to calculate and produce the samples is called **ars()**. For purposes of organization and readability, we decided to use a functional programming approach where we defined a number of auxiliary functions to complete subtasks of the **ars()** function, and wrote all of these functions in a separate R script entitled NAME OF R SCRIPT HERE. The main **ars()** function takes in the following inputs from the user:

- **f** is the density function from which the user wants to sample from
 - requirement of the argument is that it must be a univariate, log-concave function
- **n** is the number of samples the user wants to generate from the function **f**
 - requirement of the argument is that it must be a positive integer value
- **bounds** denotes the left and right bounds of the domain of the function **f**, which is the set of points x for which $f(x) > 0$
 - requirement of the argument is that it be a vector of length 2, where the first element is the left (lower) bound and the second element is the right (upper) bound
- **x_init** is the initial point within the bounds used as one of the initial abscissae points for future sampling
 - requirement of the argument is that it is a single point defined within the bounds

We note that if the user does not provide the argument **n**, the default number of samples is set to 1000, and if the user does not provide the argument **bounds**, the default bounds are set to $(-\text{Inf}, \text{Inf})$. If the user does not provide the argument **x_init**, we define **x_init** by optimizing the function **f** with the following code:

```
if (is.na(x_init)) {  
  x_init <- optim(1, f, method = 'L-BFGS-B', control = list(fnscale=-1))$par  
}
```

The **control = list(fnscale=-1)** parameter allows us to find the value of the maximum of the function as opposed to the minimum, which is the default behavior of the **optim()** function. By finding the point x that maximizes the function **f** (i.e. the mode), we ensure that we can initialize the abscissae with points of positive probability.

The first task the function carries out is checking the validity of the user inputs. Thus, we implemented checks to ensure that **f** is a function, **n** and **x_init** are numeric, **bounds** is a vector of length two such that the first element is less than the second element (i.e. the lower bound is smaller than the upper bound), and that **x_init** is a point located within the bounds. After passing these initial argument checks, the **ars()** function moves on to defining two simple functions: **h()** and **hprime()**, where both functions take in a point **x**. **h()** returns the value $\log f(x)$ and **hprime()** returns the derivative of this value; that is, $\frac{d}{dx} \log f(x)$, where we used the external package **numDeriv** to compute the derivative. Next, the **ars()** function is broken down into the three steps Gilks and Wild suggest: (1) initialization, (2) sampling, and (3) updating.

Initialization Step

The basic idea of the initialization step is to initialize the abscissae points, which we define in the function as **Tk**. In order to carry out this step, we define an auxiliary function **initialize_abscissae()**, which uses the **x_init** point as a reference in order to output a vector of $k = 20$ initial abscissae points, which we then pass through the functions **h()** and **hprime()**. After defining these 20 initial abscissae points and computing their derivatives, we initialized the functions $u_k(x)$, $s_k(x)$, $l_k(x)$ and calculated the intersection of the tangent lines, z_j , for each x_j and x_{j+1} , where $j = 1, \dots, k-1$. As described in the paper by Gilks and Wild, we derive the rejection envelope on **Tk**, and calculate the value of every x_k on the rejection envelope by function **u_k(x)**. The function first identify the segment where the inputted value x lies, and calculate the $u_k(x)$ for the corresponding slope and intercept of the piecewise function. The value from squeezing function performed similarly by the function $l_k(x)$. Therefore we derived the value of both rejection sampling value and the squeezing value for each point within **Tk** from **u_k(k)**, and **l_k(x)** respectively.

For $s_k(x)$ we approximate density function by calculating the area of under the piecewise $\exp(u_j(x))$ for each segment (z_j, z_{j+1}) , and get the unnormalized probability q_j .

$$q_j = \begin{cases} \exp(u_j(z_j))(z_{j+1} - z_j) & \text{if } h'(x_j) = 0 \\ \exp(u_j(z_{j+1})) - \exp(u_j(z_j))/h'(x_j) & \text{if } h'(x_j) \neq 0 \end{cases}$$

The normalized probability is derived by dividing each q_j with the summation of all unnormalized probability $A = \sum_{j=1}^n q_j$. The normalized probability, where the true probability we used to sample x_* of each segment (z_j, z_{j+1}) therefore equals to

$$p_j = \frac{q_j}{A}$$

Sampling Step

We use the same method of sampling as Gilks and Wild derive. That is, we sample a x^* from our **sample_sk()** function. first, we sample a point w from a $U(0, 1)$ distribution which represents the probability of x equals to x^* .

$$w = P(x = x^*)$$

Next, for $h'(x_j) = 0$ we just scaling the original uniform distribution to (z_{j-1}, z_j) to get x^* . While for $h'(x_j) \neq 0$, we use the inverse CDF on segment which sampled from the set uniform to the CDF.

$$\frac{\int_{z_j}^{x^*} \exp(u_j(x)) dx}{\int_{z_j}^{z_{j+1}} \exp(u_j(x)) dx} = \frac{\exp(u_j(x^*)) - \exp(u_j(z_j))}{Ah'(x_j)}$$

Solve for x^* , we can derive that:

$$x^* = \begin{cases} z_j + t(z_{j+1} - z_j) & \text{if } h'(x_j) = 0 \\ \log(wAh'(x_j) + \exp(u_j(z_j)) - h(x_j))/h'(x_j) & \text{if } h'(x_j) \neq 0 \end{cases}$$

We then perform the following tests:

First is the squeezing test: if $w \leq e^{l_k(x^*) - u_k(x^*)}$, we add the value x^* to our sample vector. If this is not the case, we move on to the rejection test: if $w \leq e^{h(x^*) - u_k(x^*)}$, we add the value x^* to our sample vector. Otherwise, x^* is rejected. In the event a value of x^* leads us to perform the rejection test, we also complete the updating step, defined in the next section.

Updating Step

As just mentioned, if x^* fails the squeezing test, we perform the updating step, which works as follows: we append the value of x^* to our vector T_k such that we now have T_{k+1} . We then sort the values of T_{k+1} so they are in ascending order, and compute the $h'(x^*)$, placing it in its analogous location to its position in T_{k+1} . We then calculate the new intersection points of the tangent lines in of the points in T_{k+1} and return to the beginning of the sampling step. We complete this process until we have obtained a total of n samples, and the vector of n samples are returned from the `ars()` function to the user.

Notes on the Algorithm

While designing the code, one issue we came across was in comparing if two elements in a vector were equal. As a simple example, when trying to see if $x-y==0$ for two elements x and y which were the same, we would sometimes see a result of `FALSE`. We figured out that this was because of a numerical precision error, and our solution to this was to define a tolerance for a margin of error, which we chose to be the square root of Machine Epsilon, or `1e-8`, thus rewriting the above check as `abs(x-y)<=1e-8`, which would return `TRUE` as desired. Another issue we came across was in checking that the function `f` that the user provided was indeed a log-concave function by performing checks as our calculations proceeded rather than just checking if the entire function was log-concave at the beginning of our `ars()` function. The way we solved this was to use the definition of log-concavity noted in the Gilks and Wild paper. That is, a function $f(x)$ is log-concave if $h'(x)$ decreases monotonically with increasing x . Thus, we defined a function `check_log_concave()` which ensured this condition, using the same strategy to account for numerical precision as mentioned above. Below is the implementation:

```
check_log_concave <- function(x){
  # ensure that h'(x) is decreasing monotonically
  n <- length(x)
  assertthat::assert_that(sum(x[2:n] - x[1:n-1]) <= 1e-8) == n-1,
    msg = "Function is not log-concave")
}
```

We ran this function every time we defined new points in T_k (and thus calculated a new $h'(x)$), and this seemed to work in ensuring that the function `f` provided by the user was indeed log-concave.

Tests

DESCRIBE HERE HOW TO RUN THE TESTS

In order to run the tests for our function, we created an R script which utilized the R package `testthat`. We designed a number of tests in order to ensure that our function works as desired. The first set of tests was designed to ensure that the `ars()` function errors out when the user inputs a non-log-concave function `f`. Some known densities that are not log-concave include the pareto distribution, the t-distribution, and the lognormal distribution. Thus, we inputted these densities as our function `f` to make sure the `ars()` function returned an error. The second set of tests was designed to ensure that the `ars()` function errors out when the user enters invalid inputs. Some example tests we ran were inputting a value for `f` that was not a function and inputting a non-numeric argument for the number of samples `n`. The last set of tests was designed to check that our `ars()` function returned a vector of samples that resembled a vector of samples of the same length from its respective distribution. In order to do this, we implemented the Kolmogorov-Smirnov test.

As a basic summary, the Kolmogorov-Smirnov test is a nonparametric test based on the empirical CDF function that tests the equality of samples from a continuous distribution with samples from a reference distribution, given that the distribution is fully specified (i.e. rate, shape, etc. parameters are provided). The Kolmogorov-Smirnov test statistic is defined as

$$D = \sup_x |F_0(x) - F_{data}(x)|,$$

where $F_0(x)$ is the CDF of the hypothesized distribution (“true” distribution) and $F_{data}(x)$ is the empirical CDF of the observed data (data coming from `ars()`). The null and alternative hypotheses are

H_0 : the data comes from the specified distribution

H_a : at least one value does not match the specified distribution

In order to run this test in R, we used the function `ks.test()` and ensured that the obtained p-value was greater than $\alpha = 0.05$. We ran the test on multiple known distributions in R, where we generated a vector of n samples for a known distribution (for example the Normal distribution) using our `ars()` function, and compared it to a specified “true” distribution, which was a same-length vector of samples from the known distribution using built-in R functions, (`rnorm()`, for example).

We also wrote an R script to test that our auxiliary functions were providing reasonable results; namely `check_log_convave()`, `initialize_abscissae()`, `calc_z()`, `u()`, `l()`, `calc_probs()`, and `sample_sk()`. For the majority of these tests, we provided sample values for the functions to perform computations on, and then compared the result of the function with a manual computation carrying out the same task. Given that each of these functions work under a “black-box” methodology, the sample values provided were random and contained no real meaning, but were still effective in testing the basic tasks of each function.

Examples

To visualize the results of our `ars()` function, below is an implementation of the function utilizing some common densities (LIST DENSITIES HERE) with a histogram of the samples compared to the actual curve of the distribution displayed on top. CALL FUNCTION AND CREATE PLOTS.

References

Reference Gilks Paper, <https://www.statisticshowto.com/kolmogorov-smirnov-test/>, <https://www.itl.nist.gov/div898/handbook/eda/section3/eda35g.htm>