

Simple and efficient R toolkit for ABC inference in population genetics

Martin Petr 1,2 Fernando Racimo 1,2

Approximate Bayesian Computation (ABC) is a powerful inference method for evolutionary biology. However, numerous technical challenges and significant amount of coding required often hinder its application to new modeling problems. We present a new R package, demografr, designed for building simple, efficient, and reproducible ABC pipelines entirely in R, using tree sequences for fast computation and the abc package as an inference engine. Additionally, demografr provides individual components for building other kinds of simulation-based inference workflows. Below, we demonstrate the features of demografr on a simple (but complete) ABC pipeline.

l demographic model

Demographic models for ABC inference are defined using the slendr R interface [bodkan.net/slendr]. Any function producing a slendr model can serve as model for a demografr inference, with arguments of the function representing model parameters. However, users are not limited by slendr: any SLiM or msprime script can be used as long as it creates a tskit tree sequence.

L. prior distributions

Priors are defined in an R list, using an intuitive domain-specific syntax of <parameter> ~ <sampling expression>. In addition to built-in R sampling functions like runif, users can define priors using custom-defined distributions.

5. observed summary statistics

Observed summary statistics are inputted as data frames or numeric vectors, and are collected in a named R list (one element of the list for each statistic).

4- tree-sequence summary statistics

The output of each simulation of the model (1.) will be a tree sequence. Therefore, simulated summary statistics can be computed very efficiently thanks to tskit [Kelleher et al., PLOS Comp Biol (2018)], via slendr's intuitive tskit R interface: allele frequency spectrum, f_3 , f_4 , nucleotide diversity, Identity-by-Descent sharing, and more are available. Summary functions are collected in a named R list in a format matching the list of the observed statistics (3.). File format conversion for computation in external software is therefore not needed.

Simulations are implicitly parallelized across CPU cores and/or remote servers via the R package future [Bengtsson, The R Journal (2021)]. During this process, user-defined summary statistic functions (4.) are internally applied to the simulated tree sequence produced by each simulation replicate.

O. ABC inference

perform abc() is a wrapper around the inference machinery of the abc package [Csilléry et al., Methods in Eco Evol (2012)], acting as a bridge feeding simulation results from demografr (5.) to the abc package in the correct format. Support for other methods and packages as inference "engines" is in progress.

/ analysis

The result of ABC inference (6.) is compatible with all the validation and diagnostic functions of the abc package, for which demografr provides convenient wrappers. Additionally, demografr provides its own functions for inspection and visualization of posterior distributions.

```
model <- function(Ne_A, Ne_B, Ne_C, Ne_D, T_AB, T_BC, T_CD, gf_BC) {</pre>
  popA <- population("popA", time = 1, N = Ne_A)
  popB <- population("popB", time = T_AB, N = Ne_B, parent = popA)</pre>
  popC <- population("popC", time = T_BC, N = Ne_C, parent = popB)</pre>
  popD <- population("popD", time = T_CD, N = Ne_D, parent = popC)</pre>
  gf <- gene_flow(from = popB, to = popC,</pre>
                    start = 9000, end = 9301, rate = gf_BC
  compile_model(
    populations = list(popA, popB, popC, popD), gene_flow = gf,
    generation_time = 1, simulation_length = 10000
  return(model)
 priors <- list(
   Ne_A \sim runif(100, 100000), Ne_B \sim runif(100, 100000),
   Ne_C \sim runif(100, 100000), Ne_D \sim runif(100, 100000),
   ... # split time and gene-flow priors omitted for brevity
 observed <- list(diversity = observed_diversity,
                   divergence = observed_divergence,
                               = observed_f4)
  compute_diversity <- function(ts) {</pre>
    samples <- extract_names(ts, split = "pop")</pre>
    ts_diversity(ts, sample_sets = samples)
  compute_divergence <- function(ts) { ... } # omitted for brevity</pre>
  compute_f4 <- function(ts) { ... } # omitted for brevity</pre>
  functions <- list(diversity = compute_diversity,</pre>
                     divergence = compute_divergence,
                                 = compute_f4)
   future::plan("multicore", workers = <number of CPUs>)
  data <- simulate_abc(</pre>
    model, priors, functions, observed, iterations = 100000,
     sequence_length = 10e6, recombination_rate = 1e-8
   abc <- perform_abc(data, engine = "abc",</pre>
                       tol = 0.05, method = "neuralnet")
    extract_summary(abc, "Ne_A"); plot_posterior(abc, "T")
                                            Posterior distribution of 'T' parameter(s)
                              Ne_A
                          -3.010207
  Min.:
  Weighted 2.5 % Perc.:
                         474.550342
  Weighted Median:
                        1782.484402
                                       S 6e-04
                                                                       T_AB
T_BC
T_CD
   Weighted Mean:
                        1840.627399
```

Visit bodkan.net/demografr for documentation and extensive tutorials!

Weighted Mode:

Max.:

Weighted 97.5 % Perc.: 3424.879882

1454.772471

4030.283888

3e-04