

Package ‘rubias’

January 29, 2018

Type Package

Title Bayesian Inference from the Conditional Genetic Stock Identification Model

Version 0.1.0

Maintainer Eric C. Anderson <eric.anderson@noaa.gov>

Description Implements Bayesian inference for the conditional genetic stock identification model. It allows inference of mixed fisheries and also simulation of mixtures to predict accuracy. A full description of the underlying methods is in review at Canadian Journal of Fisheries and Aquatic Sciences.

License CC0

LazyData TRUE

Depends R (>= 3.3.0)

Imports dplyr, ggplot2, gtools, magrittr, Rcpp (>= 0.12.5), readr, rlang, stringr, tibble, tidyr

LinkingTo Rcpp

RoxygenNote 6.0.1

Suggests knitr, rmarkdown, tidyverse

VignetteBuilder knitr

NeedsCompilation yes

Author Eric C. Anderson [aut, cre],
Ben Moran [aut]

Repository CRAN

Date/Publication 2018-01-29 10:12:38 UTC

R topics documented:

alewife	2
assess_pb_bias_correction	2
assess_reference_loo	4
assess_reference_mc	5

blueback	7
chinook	7
chinook_collection_levels	8
chinook_mix	8
chinook_repunit_levels	9
infer_mixture	9
perfect_chinook	10
perfect_chinook_mix	11
read_gsi_sim	11
rubias	12
self_assign	12
sim_spec_examples	13
small_chinook_mix	13
small_chinook_ref	14
write_gsi_sim_mixture	14
write_gsi_sim_reference	15
Index	17

alewife	<i>Microsat data from alewife herring reference populations</i>
---------	---

Description

Standard two-column genetic data with lots of other columns preceding it. Can be fed directly into rubias because it has at least the columns sample_type, collection, repunit and indiv.

Format

A tibble.

Source

<http://datadryad.org/resource/doi:10.5061/dryad.80f4f>

assess_pb_bias_correction	<i>Test the effects of the parametric bootstrap bias correction on a reference dataset through cross-validation</i>
---------------------------	---

Description

This is a rewrite of bias_comparison(). Eric didn't want the plotting to be wrapped up in a function, and wanted to return a more informative data frame.

Usage

```
assess_pb_bias_correction(reference, gen_start_col, seed = 5, nreps = 50,
  mixsize = 100, alle_freq_prior = list(const_scaled = 1))
```

Arguments

reference	a two-column format genetic dataset, with a "repunit" column specifying each individual's reporting unit of origin, a "collection" column specifying the collection (population or time of sampling) and "indiv" providing a unique name
gen_start_col	the first column containing genetic data in reference. All columns should be genetic format following this column, and gene copies from the same locus should be adjacent
seed	the random seed for simulations
nreps	The number of reps to do.
mixsize	The size of each simulated mixture sample.
alle_freq_prior	a one-element named list specifying the prior to be used when generating Dirichlet parameters for genotype likelihood calculations. Valid methods include "const", "scaled_const", and "empirical". See ?list_diploid_params for method details.

Details

Takes a reference two-column genetic dataset, pulls a series of random "mixture" datasets with varying reporting unit proportions from this reference, and compares the results of GSI through standard MCMC vs. parametric-bootstrap MCMC bias correction

The amount of bias in reporting unit proportion calculations increases with the rate of misassignment between reporting units (decreases with genetic differentiation), and increases as the number of collections within reporting units becomes more uneven.

Output from the standard Bayesian MCMC method demonstrates the level of bias to be expected for the input data set, and parametric bootstrapping is an empirical method for the removal of any existing bias.

Value

bias_comparison returns a list; the first element is a list of the relevant rho values generated on each iteration of the random "mixture" creation. This includes the true rho value, the standard result rho_mcmc, and the parametric bootstrapped rho_pb.

The second element is a dataframe listing summary statistics for each reporting unit and estimation method. mse, the mean squared error, summarizes the deviation of the rho estimates from their true value, including both bias and other variance. mean_prop_bias is the average ratio of residual to true value, which gives greater weight to deviations at smaller values. mean_bias is simply the average residual; unlike mse, this demonstrates the direction of the bias.

Examples

```
## Not run:
## This takes too long to run in R CMD CHECK
ale_bias <- assess_pb_bias_correction(alewife, 17)

## End(Not run)
```

assess_reference_loo *Simulate mixtures and estimate reporting group and collection proportions.*

Description

From a reference dataset, this creates a genotype-logL matrix based on simulation-by-individual with randomly drawn population proportions, then uses this in two different estimates of population mixture proportions: maximum likelihood via EM-algorithm and posterior mean from MCMC.

Usage

```
assess_reference_loo(reference, gen_start_col, reps = 50, mixsize = 100,
  seed = 5, alpha_repunit = 1.5, alpha_collection = 1.5,
  resampling_unit = "individual", alle_freq_prior = list(const_scaled = 1),
  printSummary = FALSE, return_indiv_posteriors = FALSE)
```

Arguments

reference	a two-column format genetic dataset, with "repunit", "collection", and "indiv" columns, as well as a "sample_type" column that has some "reference" entries
gen_start_col	the first column of genetic data in reference
reps	number of reps of mixture simulation and MCMC to do
mixsize	the number of individuals in each simulated mixture
seed	a random seed for simulations
alpha_repunit	If a vector, this is the dirichlet parameter for simulating the proportions of reporting units. Gets recycled to the number of reporting units. Default is 1.5. Otherwise, this could be a two-column data frame. The first column must be named "repunit" and the second one must be one of "dirichlet", "ppn", or "cnt", according to whether you wish to specify dirichlet parameters, or proportions, or exact counts, respectively, for each population. If you want to make multiple simulations, pass in a list of data frames or of individual dirichlet parameters. For examples, see sim_spec_examples .
alpha_collection	The dirichlet parameter for simulating proportions of collections within reporting units. Default = 1.5. If this is a data frame then the first column must be "collection" and the second must be one of "dirichlet", "ppn", "cnt", "sub_dirichlet",

"sub_ppn". If you want to provide multiple different scenarios. You can pass them in as a list. If alpha_repunit or alpha_collection is a list with length greater than 1, the shorter will be recycled. For examples, see [sim_spec_examples](#).

resampling_unit

what unit should be resampled. Currently the choices are "individuals" (the default) and "gene_copies". Using "individuals" preserves missing data patterns available in the reference data set. We also have "gene_copies_with_missing" capability, but it is not yet linked into this function.

alle_freq_prior

a one-element named list specifying the prior to be used when generating Dirichlet parameters for genotype likelihood calculations. Valid methods include "const", "scaled_const", and "empirical". See `?list_diploid_params` for method details.

printSummary if TRUE a summary of the reference samples will be printed to stdout.

return_indiv_posteriors

if TRUE, output is a list of 2. The first entry, `mixing_proportions`, contains the true (simulated) and estimated mixture proportions for each scenario, iteration, and collection. The second, `indiv_posteriors`, contains the posterior probability of assignment to each collection for each scenario, iteration, and individual. If FALSE, output is a single data frame, `mixing_proportions`

Examples

```
# very small number of reps so it is quick enough for example
ale_dev <- assess_reference_loo(alewife, 17, reps = 5)
```

assess_reference_mc	<i>Partition a reference dataset and estimate reporting group and collection proportions</i>
---------------------	--

Description

From a reference dataset, this draws (without replacement) a simulated mixture dataset with randomly drawn population proportions, then uses this in two different estimates of population mixture proportions: maximum likelihood via EM-algorithm and posterior mean from MCMC.

Usage

```
assess_reference_mc(reference, gen_start_col, reps = 50, mixsize = 100,
  seed = 5, alpha_repunit = 1.5, alpha_collection = 1.5,
  min_remaining = 5, alle_freq_prior = list(const_scaled = 1))
```

Arguments

reference	a two-column format genetic dataset, with "repunit", "collection", and "indiv" columns, as well as a "sample_type" column that has some "reference" entries.
gen_start_col	the first column of genetic data in reference
reps	number of reps to do
mixsize	the number of individuals in each simulated mixture.
seed	a random seed for simulations
alpha_repunit	The dirichlet parameter for simulating the proportions of reporting units. Default = 1.5
alpha_collection	The dirichlet parameter for simulating proportions of collections within reporting units. Default = 1.5
min_remaining	the minimum number of individuals which should be conserved in each reference collection during sampling without replacement to form the simulated mixture
alle_freq_prior	a one-element named list specifying the prior to be used when generating Dirichlet parameters for genotype likelihood calculations. Valid methods include "const", "scaled_const", and "empirical". See ?list_diploid_params for method details.

Details

This method is referred to as "Monte Carlo cross-validation". The input parameters for `assess_reference_mc` are more restrictive than those of `assess_reference_loo`. Rather than allowing a `data.frame` to specify Dirichlet parameters, proportions, or counts for specific reporting units and collections, `assess_reference_mc` only allows vector input (default = 1.5) for `alpha_repunit` and `alpha_collection`. These inputs specify the uniform Dirichlet parameters for all reporting units and collections, respectively.

For mixture proportion generation, the rho values are first drawn using a stick-breaking model of the Dirichlet distribution, but with proportions capped by `min_remaining`. Stick-breaking is then used to subdivide each reporting unit into collections. In addition to the constraint that mixture sampling without replacement cannot deplete the number of individuals in each collection below `min_remaining`, a similar constraint is placed upon the number of individuals left in reporting units, determined as `min_remaining * (# collections in reporting unit)`.

Note that this implies that the data are only truly Dirichlet distributed when no rejections based on `min_remaining` occur. This is a reasonable certainty with sufficient reference collection sizes relative to the desired mixture size.

Examples

```
# only 5 reps, so it doesn't take too long. Typically you would
# do many more
ale_dev <- assess_reference_mc(alewife, 17, 5)
```

blueback

Microsat data from blueback herring reference populations

Description

Standard two-column genetic data with lots of other columns preceding it. Can be fed directly into rubias because it has at least the columns sample_type, collection, repunit and indiv.

Format

A tibble.

Source

<http://datadryad.org/resource/doi:10.5061/dryad.80f4f>

chinook

SNP data from chinook reference populations

Description

Chinook salmon baseline data similar to that which can be downloaded from <http://datadryad.org/resource/doi:10.5061/dryad.574sv/1>. This data set includes 91 SNPs and 7301 fish and is what the Dryad data became after we converted from TaqMan to SNPtype assays (being forced to toss some loci) and tossed out a bunch of lousy historical samples from Trinity River.

Format

A tbl_df-ed (from dplyr) data frame with 7,301 rows and 185 variables. The first three columns are

repunit (chr) the reporting unit that the individual is in

pop (chr) the population from which the individual was sampled

ID (chr) Unique identifier of the individual fish

The remaining columns are two columns for each locus. These columns are named like, "Locus.1" and "Locus.2" for the first and second gene copies at that locus. For example, "Ots_104569-86.1" and "Ots_104569-86.2". The locus columns are ints and missing data is denoted by NA.

Source

<http://datadryad.org/resource/doi:10.5061/dryad.574sv/1>

chinook_collection_levels

a vector that gives a desired sort order of the chinook collections

Description

This is just an example of what one would use as levels in order to get the [chinook](#) collections in a desired sort order after analysis. The issue here is collection in the input data frame to most functions must be a character vector, not a factor. But, after analysis you can always make them a factor again and use a vector like this one to specify the levels.

Source

Made it up!

chinook_mix

SNP data from Chinook salmon taken in May/August 2015 from California fisheries

Description

This has data from 91 SNP markers (a subset of the 95 markers in the [chinook](#) baseline data set).

Format

A `tbl_df`-ed (from `dplyr`) data frame with 2256 rows and 193 variables. The first four columns are meta data. The remaining columns are two columns for each locus. These columns are named like, "Locus.1" and "Locus.2" for the first and second gene copies at that locus. For example, "Ots_104569-86.1" and "Ots_104569-86.2". The locus columns are ints and missing data is denoted by NA.

Source

Southwest Fisheries Science Center, Santa Cruz, CA

 chinook_repunit_levels

a vector that gives a desired sort order of the chinook repunits

Description

This is just an example of what one would use as levels in order to get the `chinook` repunits in a desired sort order after analysis. The issue here is that repunit in the input data frame to most functions must be a character vector, not a factor. But, after analysis you can always make them a factor again and use a vector like this one to specify the levels.

Source

Made it up!

 infer_mixture

Estimate mixing proportions and origin probabilities from one or several mixtures

Description

Takes a mixture and reference dataframe of two-column genetic data, and a desired method of estimation for the population mixture proportions (MCMC, PB) Returns the output of the chosen estimation method

Usage

```
infer_mixture(reference, mixture, gen_start_col, method = "MCMC",
  alle_freq_prior = list(const_scaled = 1), reps = 2000, burn_in = 100,
  pb_iter = 100, sample_int_Pi = 1)
```

Arguments

reference	a dataframe of two-column genetic format data, proceeded by "repunit", "collection", and "indiv" columns. Does not need "sample_type" column, and will be overwritten if provided
mixture	a dataframe of two-column genetic format data. Must have the same structure as reference dataframe, but "collection" and "repunit" columns are ignored. Does not need "sample_type" column, and will be overwritten if provided
gen_start_col	the first column of genetic data in both data frames
method	a choice between "MCMC", "PB" methods for estimating mixture proportions
alle_freq_prior	a one-element named list specifying the prior to be used when generating Dirichlet parameters for genotype likelihood calculations. Valid methods include "const", "scaled_const", and "empirical". See <code>?list_diploid_params</code> for method details.

reps	the number of iterations to be performed in MCMC
burn_in	how many reps to discard in the beginning of MCMC when doing the mean calculation. They will still be returned in the traces if desired.
pb_iter	how many bootstrapped data sets to do for bootstrap correction using method PB. Default is 100.
sample_int_Pi	how many iterations between storing the mixing proportions trace. Default is 1. Can't be 0. Can't be so large that fewer than 10 samples are taken from the burn in and the sweeps.

Details

"MCMC" estimates mixing proportions and individual posterior probabilities of assignment through Markov-chain Monte Carlo, while "PB" does the same with a parametric bootstrapping correction. All methods use a uniform $1/(\# \text{ collections or RUs})$ prior for the mixing proportions.

Value

Tidy data frames in a list with the following components: `mixing_proportions`: the estimated mixing proportions of the different collections. `indiv_posteriors`: the posterior probs of fish being from each of the collections. `mix_prop_traces`: the traces of the mixing proportions. Useful for computing credible intervals. `bootstrapped_proportions`: If using method "BH" this returns the bootstrap corrected reporting unit proportions.

Examples

```
mcmc <- infer_mixture(reference = small_chinook_ref,
  mixture = small_chinook_mix,
  gen_start_col = 5,
  method = "MCMC",
  reps = 200)
```

perfect_chinook *perfect-assignment genetic data for chinook.*

Description

This is just like the `chinook` data, but only has 7 loci and all loci are fixed in fortuitous patterns so that every single collection is easily resolved. This is primarily useful for testing purposes.

Source

Made it up!

perfect_chinook_mix	<i>perfect-assignment mixture genetic data for chinook.</i>
---------------------	---

Description

This is similar to the [chinook_mix](#) data, but only has 7 loci and all loci are fixed in fortuitous patterns so that every single collection is easily resolved. This is primarily useful for testing purposes. The name of the individual has its collection inside the colons.

Source

Made it up!

read_gsi_sim	<i>read a gsi_sim formatted input file into a tibble that rubias can use</i>
--------------	--

Description

Note that this relies on a system call to awk. It probably won't work on Windows.

Usage

```
read_gsi_sim(path, sample_type, repunits = NULL)
```

Arguments

path	path to the gsi_sim file
sample_type	should be "reference" or "mixture" depending on what kind of file it is
repunits	the gsi_sim reporting units file. Has no effect if sample_type is "mixture". If sample_type is "reference" and this is left as NULL, then all collections will be put in the the "default_repu".

rubias	<i>rubias: Bayesian inference from the conditional genetic stock identification model</i>
--------	---

Description

Read the "rubias-overview" vignette for information on data input formats and how to use the package

the rubias main high-level functions

The following functions are wrappers, designed for user-friendly input and useful output:

[infer_mixture](#) is used to perform genetic stock identification. Options include standard MCMC and the parametric bootstrap bias correction.

[self_assign](#) does simple self-assignment of individuals in a reference data set to the various collections in the reference data set.

[assess_reference_loo](#) does leave-one-out based simulations to predict how accurately GSI can be done.

[assess_reference_mc](#) uses Monte-Carlo cross-validation based simulations to predict how accurately GSI can be done.

[assess_pb_bias_correction](#) attempts to demonstrate how much (or little) improvement can be expected from the parametric bootstrap correction given a particular reference data set.

genetic data format

See the vignette.

example data

[alewife](#), [blueback](#), and [chinook](#) are genetic data sets that are useful for playing around with rubias and testing it out.

self_assign	<i>Do leave-one-out self-assignment of individuals in a reference baseline</i>
-------------	--

Description

Returns a tidy data frame

Usage

```
self_assign(reference, gen_start_col, preCompiledParams = NULL,
  alle_freq_prior = list(const_scaled = 1))
```

Arguments

reference	a two-column format genetic dataset, with "repunit", "collection", and "indiv" columns, as well as a "sample_type" column that has some "reference" entries
gen_start_col	the first column of genetic data in reference
preCompiledParams	Users should never use this option. It is here only so that this function can be called on a precompiled set of parameters with infer_mixture. Don't use this, unless you are one of the package developers...
alle_freq_prior	a one-element named list specifying the prior to be used when generating Dirichlet parameters for genotype likelihood calculations. Valid methods include "const", "scaled_const", and "empirical". See ?list_diploid_params for method details.

Value

a tibble ...

Examples

```
ale_sa <- self_assign(alewife, 17)
```

sim_spec_examples	<i>List of example ways of specifying repunit and collection quantities in simulations</i>
-------------------	--

Description

This is just a list of tibbles that can be passed to the alpha_repunit or the alpha_collection parameters in, for example, [assess_reference_loo](#).

Source

Made it up!

small_chinook_mix	<i>Small sample of SNP data from Chinook salmon taken in May/August 2015 from California fisheries</i>
-------------------	--

Description

This is simply a sample of 100 fish from [chinook](#).

Source

Southwest Fisheries Science Center, Santa Cruz, CA

small_chinook_ref	<i>SNP data from selected chinook reference populations</i>
-------------------	---

Description

A small number of populations from the Chinook salmon baseline data similar to that which can be downloaded from <http://datadryad.org/resource/doi:10.5061/dryad.574sv/1>. This data set includes 91 SNPs and 909 fish.

Format

A tbl_df-ed (from dplyr) data frame with 909 rows and 185 variables. The first three columns are

repunit (chr) the reporting unit that the individual is in

pop (chr) the population from which the individual was sampled

ID (chr) Unique identifier of the individual fish

The remaining columns are two columns for each locus. These columns are named like, "Locus.1" and "Locus.2" for the first and second gene copies at that locus. For example, "Ots_104569-86.1" and "Ots_104569-86.2". The locus columns are ints and missing data is denoted by NA.

Source

<http://datadryad.org/resource/doi:10.5061/dryad.574sv/1>

write_gsi_sim_mixture	<i>Write a mixture data frame to gsi_sim format baseline and repunits file</i>
-----------------------	--

Description

Note, this is only intended to work with integer-valued alleles, at the moment. It was just written for testing and verifying that things are working correctly.

Usage

```
write_gsi_sim_mixture(mix, gen_start_col, mixprefix)
```

Arguments

mix	mixture data frame
gen_start_col	column in which the genetic data start
mixprefix	path to write the mixture file to. The mixture collection name + .txt will be appended to this. This path can include directories if they exist. An example would be <code>"/my_gsi_data/mixture"</code> . This is a required argument.

Examples

```
# this writes to file prefix "mixfile" in a temporary directory
dd <- tempdir()
prefix <- file.path(dd, "mixfile")

# print that
prefix

# note that in practice you will probably want to specify
# your own directory...

# run the function
write_gsi_sim_mixture(chinook_mix, 5, prefix)

# see where those files live:
dir(dd, pattern = "mixfile*", full.names = TRUE)
```

```
write_gsi_sim_reference
```

Write a reference data frame to gsi_sim format baseline and repunits file

Description

Note, this is only intended to work with integer-valued alleles, at the moment. It was just written for testing and verifying that things are working correctly.

Usage

```
write_gsi_sim_reference(ref, gen_start_col, baseout = "baseline.txt",
  repout = "repunits.txt")
```

Arguments

ref	reference data frame
gen_start_col	column in which the genetic data start
baseout	path to write the baseline file to. Required.
repout	path to write the repunits file to. Required.

Examples

```
# create a temp directory to put example outputs
dd <- tempdir()
basefile <- file.path(dd, "baseline.txt")
repunitsfile <- file.path(dd, "repunits.txt")

# print those
basefile
```

```
repunitsfile

# note that in practice you will probably want to specify
# your own filepaths...

# run the function
write_gsi_sim_reference(alewife, 17, basefile, repunitsfile)
```


Index

alewife, [2](#), [12](#)
assess_pb_bias_correction, [2](#), [12](#)
assess_reference_loo, [4](#), [12](#), [13](#)
assess_reference_mc, [5](#), [12](#)

blueback, [7](#), [12](#)

chinook, [7](#), [8–10](#), [12](#), [13](#)
chinook_collection_levels, [8](#)
chinook_mix, [8](#), [11](#)
chinook_repunit_levels, [9](#)

infer_mixture, [9](#), [12](#)

perfect_chinook, [10](#)
perfect_chinook_mix, [11](#)

read_gsi_sim, [11](#)
rubias, [12](#)
rubias-package (rubias), [12](#)

self_assign, [12](#), [12](#)
sim_spec_examples, [4](#), [5](#), [13](#)
small_chinook_mix, [13](#)
small_chinook_ref, [14](#)

write_gsi_sim_mixture, [14](#)
write_gsi_sim_reference, [15](#)