

Package ‘ML.MSBD’

July 9, 2018

Type Package

Title Maximum Likelihood Inference on Multi-State Trees

Version 1.1.0

Date 2018-06-20

Description Inference of a multi-states birth-death model from a phylogeny, comprising a number of states N , birth and death rates for each state and on which edges each state appears. Inference is done using a hybrid approach: states are progressively added in a greedy approach. For a fixed number of states N the best model is selected via maximum likelihood. Reference: J. Barido-Sottani and T. Stadler (2017) <doi:10.1101/215491>.

License GPL-3

Imports ape (>= 4.1)

Suggests knitr

RoxygenNote 6.0.1

VignetteBuilder knitr

NeedsCompilation no

Author Joelle Barido-Sottani [aut, cre]

Maintainer Joelle Barido-Sottani <joelle.barido-sottani@m4x.org>

Repository CRAN

Date/Publication 2018-07-09 09:30:17 UTC

R topics documented:

ML.MSBD-package	2
likelihood_MSBD	3
likelihood_MSBD_unresolved	5
ML_MSBD	7

Index	11
--------------	-----------

ML.MSBD-package

*Maximum Likelihood Inference on Multi-State Trees***Description**

Inference of a multi-states birth-death model from a phylogeny, comprising a number of states N , birth and death rates for each state and on which edges each state appears. Inference is done using a hybrid approach: states are progressively added in a greedy approach. For a fixed number of states N the best model is selected via maximum likelihood. Reference: J. Barido-Sottani and T. Stadler (2017) <doi:10.1101/215491>.

Details

```

Package:      ML.MSBD
Type:         Package
Title:        Maximum Likelihood Inference on Multi-State Trees
Version:      1.1.0
Date:         2018-06-20
Authors@R:    person("Joelle", "Barido-Sottani", email = "joelle.barido-sottani@m4x.org", role = c("aut", "cre"))
Description:  Inference of a multi-states birth-death model from a phylogeny, comprising a number of states N, birth and
License:      GPL-3
Imports:      ape (>= 4.1)
Suggests:    knitr
RoxygenNote: 6.0.1
VignetteBuilder: knitr
Author:       Joelle Barido-Sottani [aut, cre]
Maintainer:   Joelle Barido-Sottani <joelle.barido-sottani@m4x.org>

```

Index of help topics:

```

ML.MSBD-package      Maximum Likelihood Inference on Multi-State
                      Trees
ML_MSBD              Full Maximum Likelihood inference of birth and
                      death rates together with their changes along a
                      phylogeny under a multi-type birth-death model.
likelihood_MSBD      Likelihood calculation for randomly sampled
                      trees
likelihood_MSBD_unresolved
                      Likelihood calculation for unresolved trees

```

Author(s)

NA

Maintainer: NA

References

J. Barido-Sottani and T. Stadler. Accurate detection of HIV transmission clusters from phylogenetic trees using a multi-state birth-death model, *BioRxiv* 2017. (<https://www.biorxiv.org/content/early/2017/11/10/215491>)

See Also

[ape](#)

Examples

```
# Simulate a random phylogeny
set.seed(25)
tree <- ape::rtree(10)

# Calculate the log likelihood under a multi-states model with 2 states
# and full extant & extinct sampling
likelihood_MSBD(tree, shifts = matrix(c(2,1.8,2), nrow = 1),
  gamma = 0.05, lambdas = c(10, 6), mus = c(1, 0.5), sigma = 1)

# Infer the most likely multi-states birth-death model with full extant & extinct sampling
## Not run: ML_MSBD(tree, initial_values = c(0.1, 10, 1), sigma = 1, time_mode = "mid")
# Infer the most likely multi-states birth-death model with exponential decay
# and full extant & extinct sampling
## Not run: ML_MSBD(tree, initial_values = c(0.1, 10, 0.5, 1), sigma = 1,
  stepsize = 0.1, time_mode = "mid")
## End(Not run)
```

likelihood_MSBD

Likelihood calculation for randomly sampled trees

Description

Calculates the negative log likelihood of a multi-states model given a tree. This function is designed to work with constant extant and/or extinct sampling.

Usage

```
likelihood_MSBD(tree, shifts, gamma, lambdas, mus, lambda_rates = NULL,
  stepsize = NULL, uniform_weights = TRUE, p_lambda = 0, p_mu = 0,
  rho = 1, sigma = 0, rho_samplng = TRUE, add_time = 0,
  unresolved = FALSE)
```

Arguments

tree	Phylogenetic tree (in ape format) to calculate the likelihood on.
shifts	Matrix describing the positions (edges and times) of shifts. See 'Details'.
gamma	Rate of state change.

lambdas	Birth rates of all states.
mus	Death rates of all states.
lambda_rates	Rates of decay of birth rate for all states. To use exponential decay, stepsize should also be provided.
stepsize	Size of the step to use for time discretization with exponential decay, default NULL. To use exponential decay, lambda_rates should also be provided.
uniform_weights	Whether all states are weighted uniformly in shifts, default TRUE. If FALSE, the weights of states are calculated from the distributions p_lambda and p_mu. See 'Details'.
p_lambda	Prior probability distribution on lambdas, used if uniform_weights = FALSE.
p_mu	Prior probability distribution on mus, used if uniform_weights = FALSE.
rho	Sampling proportion on extant tips, default 1.
sigma	Sampling probability on extinct tips (tips are sampled upon extinction), default 0.
rho_sampling	Whether the most recent tips should be considered extant tips, sampled with sampling proportion rho. If FALSE, all tips will be considered extinct tips, sampled with sampling probability sigma. Should be TRUE for most macroevolution datasets and FALSE for most epidemiology datasets.
add_time	The time between the most recent tip and the end of the process (>=0). This is an internal variable used in calculations for unresolved trees.
unresolved	Whether this tree is the backbone of an unresolved tree. This is an internal variable used in calculations for unresolved trees.

Details

It is to be noted that all times are counted backwards, with the most recent tip positioned at 0.

The 'shifts' matrix is composed of 3 columns and a number of rows. Each row describes a shift: the first column is the index of the edge on which the shift happens, the second column is the time of the shift, and the third column is the index of the new state. For example the row vector (3,0.5,2) specifies a shift on edge number 3, at time 0.5, towards the state that has parameters lambdas[2], lambda_rates[2] and mus[2].

The weights w are used for calculating the transition rates q from each state i to j : $q_{i,j} = \gamma * w_{i,j}$. If `uniform_weights = TRUE`, $w_{i,j} = \frac{1}{N-1}$ for all i,j , where N is the total number of states. If `uniform_weights = FALSE`, $w_{i,j} = \frac{p_\lambda(\lambda_j)p_\mu(\mu_j)}{\sum_{k \neq i} p_\lambda(\lambda_k)p_\mu(\mu_k)}$ where the distributions p_λ and p_μ are provided by the inputs `p_lambda` and `p_mu`.

Value

The value of the negative log likelihood of the model given the tree.

Examples

```
# Simulate a random phylogeny
set.seed(25)
tree <- ape::rtree(10)

# Calculate the log likelihood under a constant birth-death model (i.e, no shifts)
# with full extant & extinct sampling
likelihood_MSBD(tree, shifts = c(), gamma = 0, lambdas = 10, mus = 1, sigma = 1)
# Calculate the log likelihood under a multi-states model with 2 states
# and full extant & extinct sampling
likelihood_MSBD(tree, shifts = matrix(c(2,1.8,2), nrow = 1),
  gamma = 0.05, lambdas = c(10, 6), mus = c(1, 0.5), sigma = 1)
# Calculate the log likelihood under a multi-states model with 2 states and exponential decay
# with full extant & extinct sampling
likelihood_MSBD(tree, shifts = matrix(c(2,1.8,2), nrow = 1),
  gamma = 0.05, lambdas = c(10, 6), mus = c(1, 0.5),
  sigma = 1, stepsize = 0.01, lambda_rates = c(0.1, 0.1))
```

likelihood_MSBD_unresolved

Likelihood calculation for unresolved trees

Description

Calculates the negative log likelihood of a multi-states model given a tree. This function is designed to work with unresolved trees, where tips represent collapsed clades. This sampling scheme is not recommended for epidemiology datasets. The MRCA times of collapsed clades and the number of collapsed lineages need to be provided for all tips. If neither is provided the function will default to random sampling. Extinct tips can be present outside of the unresolved parts, but not below the time(s) set for tcut.

Usage

```
likelihood_MSBD_unresolved(tree, shifts, gamma, lambdas, mus,
  lambda_rates = NULL, stepsize = NULL, uniform_weights = TRUE,
  p_lambda = 0, p_mu = 0, rho = 1, sigma = 0, rho_sampling = TRUE,
  lineage_counts = c(), tcut = NULL)
```

Arguments

tree	Phylogenetic tree (in ape format) to calculate the likelihood on.
shifts	Matrix describing the positions (edges and times) of shifts. See 'Details'.
gamma	Rate of state change.
lambdas	Birth rates of all states.
mus	Death rates of all states.

lambda_rates	Rates of decay of birth rate for all states. To use exponential decay, stepsize should also be provided.
stepsize	Size of the step to use for time discretization with exponential decay, default NULL. To use exponential decay, lambda_rates should also be provided.
uniform_weights	Whether all states are weighted uniformly in shifts, default TRUE. If FALSE, the weights of states are calculated from the distributions p_lambda and p_mu. See 'Details'.
p_lambda	Prior probability distribution on lambdas, used if uniform_weights = FALSE.
p_mu	Prior probability distribution on mus, used if uniform_weights = FALSE.
rho	Sampling proportion on extant tips, default 1.
sigma	Sampling probability on extinct tips (tips are sampled upon extinction), default 0.
rho_sampling	Whether the most recent tips should be considered extant tips, sampled with sampling proportion rho. If FALSE, all tips will be considered extinct tips, sampled with sampling probability sigma. Should be TRUE for most macroevolution datasets and FALSE for most epidemiology datasets.
lineage_counts	Number of lineages collapsed on each tip. Should be set to 1 for extinct tips.
tcut	Times of clade collapsing for each tip (i.e time of the MRCA of all collapsed lineages). Can be a single number or a vector of length the number of tips.

Details

It is to be noted that all times are counted backwards, with the most recent tip positioned at 0.

The 'shifts' matrix is composed of 3 columns and a number of rows. Each row describes a shift: the first column is the index of the edge on which the shift happens, the second column is the time of the shift, and the third column is the index of the new state. For example the row vector (3,0.5,2) specifies a shift on edge number 3, at time 0.5, towards the state that has parameters lambdas[2], lambda_rates[2] and mus[2].

The weights w are used for calculating the transition rates q from each state i to j : $q_{i,j} = \gamma * w_{i,j}$. If `uniform_weights = TRUE`, $w_{i,j} = \frac{1}{N-1}$ for all i,j , where N is the total number of states. If `uniform_weights = FALSE`, $w_{i,j} = \frac{p_\lambda(\lambda_j)p_\mu(\mu_j)}{\sum_{k \neq i} p_\lambda(\lambda_k)p_\mu(\mu_k)}$ where the distributions p_λ and p_μ are provided by the inputs `p_lambda` and `p_mu`.

Value

The value of the negative log likelihood of the model given the tree.

Examples

```
# Simulate a random phylogeny
set.seed(24)
tree <- ape::rcoal(10)
```

```
# Calculate the log likelihood under a constant birth-death model (i.e, no shifts)
# with unresolved tips
likelihood_MSBD_unresolved(tree, shifts = c(), gamma = 0, lambdas = 10, mus = 1,
                           lineage_counts = c(2,5,1,3,1,1,1,1,2,6), tcut = 0.05)
# Calculate the log likelihood under a multi-states model with 2 states and unresolved tips
likelihood_MSBD_unresolved(tree, shifts = matrix(c(2,0.7,2), nrow = 1),
                           gamma = 0.05, lambdas = c(10, 5), mus = c(1, 1),
                           lineage_counts = c(2,5,1,3,1,1,1,1,2,6), tcut = 0.05)
```

ML_MSBD

Full Maximum Likelihood inference of birth and death rates together with their changes along a phylogeny under a multi-type birth-death model.

Description

Infers a complete MSBD model from a phylogeny, including the most likely number of states, positions and times of state changes, and parameters associated with each state. Uses a greedy approach to add states and Maximum Likelihood inference for the other parameters.

Usage

```
ML_MSBD(tree, initial_values, uniform_weights = TRUE, p_lambda = 0,
         p_mu = 0, rho = 1, sigma = 0, rho_sampling = TRUE,
         lineage_counts = c(), tcut = 0, stepsize = NULL,
         no_extinction = FALSE, fixed_gamma = NULL, unique_lambda = FALSE,
         unique_mu = FALSE, optim_control = list(), attempt_remove = TRUE,
         max_nshifts = Inf, saved_state = NULL, save_path = NULL,
         time_mode = c("3pos", "tip", "mid", "root"), fast_optim = FALSE)
```

Arguments

<code>tree</code>	Phylogenetic tree (in ape format) to calculate the likelihood on.
<code>initial_values</code>	Initial values for the optimizer, to be provided as a vector in this order: gamma (optional), lambda, lambda decay rate (optional), mu (optional). See 'Details'.
<code>uniform_weights</code>	Whether all states are weighted uniformly in shifts, default TRUE. If FALSE, the weights of states are calculated from the distributions <code>p_lambda</code> and <code>p_mu</code> . See 'Details'.
<code>p_lambda</code>	Prior probability distribution on lambdas, used if <code>uniform_weights = FALSE</code> .
<code>p_mu</code>	Prior probability distribution on mus, used if <code>uniform_weights = FALSE</code> .
<code>rho</code>	Sampling proportion on extant tips, default 1.
<code>sigma</code>	Sampling probability on extinct tips (tips are sampled upon extinction), default 0.

<code>rho_sampling</code>	Whether the most recent tips should be considered extant tips, sampled with sampling proportion ρ . If FALSE, all tips will be considered extinct tips, sampled with sampling probability σ . Should be TRUE for most macroevolution datasets and FALSE for most epidemiology datasets.
<code>lineage_counts</code>	For trees with clade collapsing. Number of lineages collapsed on each tip. Should be set to 1 for extinct tips.
<code>tcut</code>	For trees with clade collapsing. Times of clade collapsing for each tip (i.e time of the MRCA of all collapsed lineages). Can be a single number or a vector of length the number of tips.
<code>stepsize</code>	Size of the step to use for time discretization with exponential decay, default NULL. To use exponential decay, an initial value for <code>lambda_rates</code> should also be provided.
<code>no_extinction</code>	Whether to use the Yule process ($\mu=0$) for all states, default FALSE. If TRUE no initial value for μ is needed.
<code>fixed_gamma</code>	Value to which γ should be fixed, default NULL. If provided no initial value for γ is needed.
<code>unique_lambda</code>	Whether to use the same value of λ for all states, default FALSE. If TRUE and exponential decay is active all states will also share the same value for <code>lambda_rate</code> .
<code>unique_mu</code>	Whether to use the same value of μ for all states, default FALSE.
<code>optim_control</code>	Control list for the optimizer, corresponds to control input in <code>optim</code> function, see <code>?optim</code> for details.
<code>attempt_remove</code>	Whether to attempt to remove shifts at the end of the inference, default TRUE. If FALSE, use a pure greedy algorithm.
<code>max_nshifts</code>	Maximum number of shifts to test for, default Inf.
<code>saved_state</code>	If provided, the inference will be restarted from this state.
<code>save_path</code>	If provided, the progress of the inference will be saved to this path after each optimization step.
<code>time_mode</code>	String controlling the time positions of inferred shifts. See 'Details'.
<code>fast_optim</code>	Whether to use the faster mode of optimization, default FALSE. If TRUE only rates associated with the state currently being added to the tree and its ancestor will be optimized at each step, otherwise all rates are optimized.

Details

It is to be noted that all times are counted backwards, with the most recent tip positioned at 0.

Five time modes are possible for the input `time_mode`. In `tip` mode, the shifts will be placed at 10% of the length of the edge. In `mid` mode, the shifts will be placed at 50% of the length of the edge. In `root` mode, the shifts will be placed at 90% of the length of the edge. In `3pos` mode, the three "tip", "mid" and "root" positions will be tested.

The weights w are used for calculating the transition rates q from each state i to j : $q_{i,j} = \gamma * w_{i,j}$. If `uniform_weights = TRUE`, $w_{i,j} = \frac{1}{N-1}$ for all i,j , where N is the total number of states. If `uniform_weights = FALSE`, $w_{i,j} = \frac{p_\lambda(\lambda_j)p_\mu(\mu_j)}{\sum_{k \neq i} p_\lambda(\lambda_k)p_\mu(\mu_k)}$ where the distributions p_λ and p_μ are provided by the inputs `p_lambda` and `p_mu`.

Initial values for the optimization need to be provided as a vector and contain the following elements (in order): an initial value for `gamma`, which is required unless `fixed_gamma` is provided, an initial value for `lambda` which is always required, an initial value for `lambda_decay_rate`, which is required if `stepsize` is provided, and an initial value for `mu`, which is required unless `no_extinction = TRUE`. An error will be raised if the number of initial values provided does not match the one expected from the rest of the settings, and the function will fail if the likelihood cannot be calculated at the initial values.

Value

Returns a list describing the most likely model found, with the following components:

<code>likelihood</code>	the negative log likelihood of the model
<code>shifts.edge</code>	the indexes of the edges where shifts happen, 0 indicates the root state
<code>shifts.time</code>	the time positions of shifts
<code>gamma</code>	the rate of state change
<code>lambdas</code>	the birth rates of all states
<code>lambda_rates</code>	if exponential decay was activated, the rates of decay of birth rate for all states
<code>mus</code>	the death rates of all states
<code>best_models</code>	a vector containing the negative log likelihood of the best model found for each number of states tested (<code>best_models[i]</code> corresponds to i states, i.e $i-1$ shifts)

All vectors are indexed in the same way, so that the state with parameters `lambdas[i]`, `lambda_rates[i]` and `mus[i]` starts on edge `shifts.edge[i]` at time `shifts.time[i]`.

Examples

```
# Simulate a random phylogeny
set.seed(25)
tree <- ape::rtree(10)

# Infer the most likely multi-states birth-death model
# with full extant & extinct sampling
## Not run: ML_MSBD(tree, initial_values = c(0.1, 10, 1), sigma = 1, time_mode = "mid")
# Infer the most likely multi-states birth-death model with exponential decay
# and full extant & extinct sampling
## Not run: ML_MSBD(tree, initial_values = c(0.1, 10, 0.5, 1), sigma = 1,
                    stepsize = 0.1, time_mode = "mid")
## End(Not run)

# Simulate a random phylogeny with extant samples
set.seed(24)
```

```
tree2 <- ape::rcoal(10)

# Infer the most likely multi-states Yule model with partial extant sampling
## Not run: ML_MSBD(tree2, initial_values = c(0.1, 10), no_extinction = TRUE,
                    rho = 0.5, time_mode = "mid")
## End(Not run)
# Infer the most likely multi-states birth-death model with full extant sampling
# and unresolved extant tips
## Not run: ML_MSBD(tree2, initial_values = c(0.1, 10, 1),
                    lineage_counts = c(2,5,1,3,1,1,1,1,2,6), tcut = 0.05, time_mode = "mid")
## End(Not run)
```

Index

*Topic **package**

ML.MSBD-package, [2](#)

ape, [3](#)

likelihood_MSBD, [3](#)

likelihood_MSBD_unresolved, [5](#)

ML.MSBD (ML.MSBD-package), [2](#)

ML.MSBD-package, [2](#)

ML_MSBD, [7](#)