

Introduction to SIDER

Kevin Healy

2017-02-16

This package estimates Trophic Discrimination Factors (TDF) based on the imputation function within the MCMCglmm package and includes functionality to include the error associated with building phylogenetic trees using the MulTree package.

Installation

To install SIDER its dependency MulTree package must first be installed directly from GitHub using the following:

```
# check if you have devtools and install if you do not
if(!require(devtools)) install.packages("devtools")

# check if you have mulTree and install if you do not using devtools
if(!require(mulTree)) devtools::install_github("TGuillerme/mulTree", ref = "master")
```

Following this you can then install SIDER directly from GitHub using the following:

```
# Check if you have SIDER and install from github if you do not
if(!require(SIDER)) devtools::install_github("healyke/SIDER", ref = "master")
```

Load the package

```
# You could load mulTree if you want, but you do not need to as you can
# call its functions using e.g. mulTree::tree.bind as below.
#library(mulTree)
library(SIDER)
```

Read in the data

SIDER has a data file which contains discrimination factors for a range of species. It also contains what tissue was used to measure the isotopic value and the basic ecology of the species including whether it is a herbivore etc and whether it is marine or terrestrial. This data is used to run a model and impute a value for the new species. Let's look at the data.

```
#read in the data
SIDER.data <- read.csv(file = system.file("extdata",
                                           "SIDER_data.csv",
                                           package = "SIDER"),
                      header = TRUE,
                      stringsAsFactors = FALSE)

# view the first 10 rows of the dataframe
head(SIDER.data)
```

```
##           species      habitat taxonomic.class tissue diet.type
## 1 Rattus_norvegicus terrestrial      mammalia  liver herbivore
## 2 Rattus_norvegicus terrestrial      mammalia  liver  carnivore
```

```
## 3 Rattus_norvegicus terrestrial      mammalia liver omnivore
## 4 Rattus_norvegicus terrestrial      mammalia liver omnivore
## 5 Rattus_norvegicus terrestrial      mammalia blood herbivore
## 6 Rattus_norvegicus terrestrial      mammalia blood carnivore
##   source.iso.13C source.iso.15N delta13C delta15N
## 1      -25.3      4.8      3.1      1.2
## 2      -16.2     12.3      3.1      0.6
## 3      -24.5      7.1      1.5      3.2
## 4      -16.7      7.4      2.2      3.0
## 5      -25.3      4.8      1.5      3.6
## 6      -16.2     12.3      0.7      3.7
```

```
#read in the phylogenetic information
#first the mammal trees
mammal_trees <- ape::read.tree(system.file("extdata",
                                           "3firstFritzTrees.tre",
                                           package = "SIDER"))

#then the bird trees
bird_trees   <- ape::read.tree(system.file("extdata",
                                           "3firstJetzTrees.tre",
                                           package = "SIDER"))

#combine them together using the tree.bind function from the mulTree package
combined_trees <- mulTree::tree.bind(x = mammal_trees,
                                     y = bird_trees,
                                     sample = 2,
                                     root.age = 250)
```

As may we want to include the error associated with building phylogenies into our analysis we take a sample of the possible trees (See Healy et al 2014). In this case we combine the mammal and bird phylogenies at a rooted age of 250 mya and take a sample of two possible trees.

Testing the new data: recipeSider

In order to estimate a trophic enrichment factor for a new species we need to check that the species is already present in our phylogeny and check what data is available for the new species.

recipeSider checks for the presence of the following data: tissue type (“blood”, “claws”, “collagen”, “feather”, “hair”, “kidney”, “liver”, “milk”, “muscle”), habitat (“terrestrial”, “marine”), and diet.type (“carnivore”, “herbivore”, “omnivore”, “pellet”).

```
#####function that checks the data for some species we want to estimate TEF for
new.data.test <- recipeSider(species = "Meles_meles",
                             habitat = "terrestrial",
                             taxonomic.class = "mammalia",
                             tissue = "blood",
                             diet.type = "omnivore",
                             tree = combined_trees)
```

If the species is not in the phylogeny already, such as say the Komodo dragon (*Varanus komodoensis*), or we are missing values as donated by NA, say tissue type, we get an error message to indicate what is missing from our data. N.B, the following code will throw a stop error if evaluated (it is silenced in this vignette).

```
#####function that checks the dat for some species we want to estimate TEF for
new.data.dummy <- recipeSider(species = "Varanus_komodoensis",
```

```
habitat = "terrestrial",
taxonomic.class = "mammalia",
tissue = "NA",
diet.type = "omnivore",
tree = combined_trees)
```

The recipeSider function also formats the data for the new species data so that it can be combined with the data already available within the package using the prepareSider function.

Formatting the new data: prepareSider

We now need to format the data by combining both the isotopic data already available within the package and the data from the new species and matching these species to the included phylogeny. We also include what isotope we want to estimate a trophic discrimination value for (either “carbon” or “nitrogen”).

```
tdf_data_c <- prepareSider(new.data.test,
                           SIDER.data,
                           combined_trees,
                           "carbon")
```

We now have a mulTree class object, which is required by the imputation analysis. It contains the matched phylogenies, in this case two phylogenies

```
tdf_data_c$phy
```

```
## 2 phylogenetic trees
```

and a dataset containing TDF and related data with the new species for which you want to estimate a trophic enrichment factor at the top with a NA for either delta13C or delta15N depending on isotope.

```
head(tdf_data_c$data)
```

```
##           sp.col      habitat taxonomic.class tissue diet.type
## 1      Meles_meles terrestrial      mammalia  blood  omnivore
## 120 Rattus_norvegicus terrestrial      mammalia  liver  herbivore
## 2      Rattus_norvegicus terrestrial      mammalia  liver  carnivore
## 3      Rattus_norvegicus terrestrial      mammalia  liver  omnivore
## 4      Rattus_norvegicus terrestrial      mammalia  liver  omnivore
## 5      Rattus_norvegicus terrestrial      mammalia  blood  herbivore
## source.iso.13C delta13C          animal
## 1              NA      NA      Meles_meles
## 120          -25.3      3.1 Rattus_norvegicus
## 2              -16.2      3.1 Rattus_norvegicus
## 3              -24.5      1.5 Rattus_norvegicus
## 4              -16.7      2.2 Rattus_norvegicus
## 5              -25.3      1.5 Rattus_norvegicus
```

Running the analysis: prepareSider

With the data formatted as a mulTree object we can decide on a model which will impute the new species estimate. In this case we will run the full model to estimate delta13C with the fixed factors of diet type and habitat type.

```
formula.c <- delta13C ~ diet.type + habitat
```

and random terms that includes the “animal” term which is required to include phylogeny into the analysis, sp.col to allow multiple species entries into the analysis and account for variation within a species, and tissue type.

```
random.terms <- ( ~ animal + species + tissue)
```

As we rely on Bayesian imputation to estimate the missing value we also need to specify a prior, in this case we use a non-informative prior as recommended in the MCMCglmm guidelines.

```
prior <- list(R = list(V = 1/4, nu=0.002),
             G = list(G1=list(V = 1/4, nu=0.002),
                     G2=list(V = 1/4, nu=0.002),
                     G3=list(V = 1/4, nu=0.002)))
```

along with the number of iterations to run the chain (nitt), the burn-in (burnin), the sampling thinning (thin), the number of chains to run (no.chains). (See MCMCglmm guidelines.)

```
# These produce results that pass convergence and avoid autocorrelation
# in the chains, but take a while to run.
# nitt <- c(1200000)
# burnin <- c(200000)
# thin <- c(500)
# no.chains <- c(2)

# NB settings only for testing and to build the vignettes quickly
nitt <- c(10)
burnin <- c(1)
thin <- c(1)
parameters <- c(nitt, thin, burnin)
no.chains <- c(2)
```

We need to check that our MCMC chains are converging so we use the Gelman and Rubin diagnostic to check the convergence and also check that the estimated parameters have an effective sample size >1000.

```
convergence = c(1.1)
ESS = c(1000)
```

As the function exports the model output to avoid memory issues within R when running over multiple phylogenies make sure you have set the working directory to somewhere appropriate. In this vignette example, we will use the default temporary directory determined by `tempdir()`, but ordinarily, this would be a folder of your own choosing located somewhere sensible on your local machine. We **strongly** advise that you do not use this `tempdir` in your own analyses! In this example, I store the original working directory to `origwd` so that we can return the R session to the original working directory after using the temporary one. Then we can finally we can run the analysis using the `imputeSider` function. This model will normally take approximately 5 minutes, however for brevity we will run a much shorter chain.

```
origwd <- getwd() # get the current, or original working directory.
setwd(tempdir()) #
TDF_est.c <- imputeSider(mulTree.data = tdf_data_c,
                       formula = formula.c,
                       random.terms = random.terms,
                       prior = prior,
                       output = "test_c_run",
                       parameters = parameters,
                       chains = no.chains,
```

```

convergence = convergence,
ESS = ESS)

## Loading required package: MCMCglmm
## Loading required package: Matrix
## Loading required package: coda
## Loading required package: ape
##
## 2017-02-16 - 08:51:52: MCMCglmm performed on tree 1
## Convergence diagnosis:
## Effective sample size is > 1000: FALSE
## 9; 9; 67.64715; 9; 9; 9; 9; 9; 9; 9; 9; 9; 9; 64.21112; 1.673816; 9; 9; 9
## C1.Sol.(Intercept), C1.Sol.diet.typeherbivore, C1.Sol.diet.typeomnivore, C1.Sol.diet.typepellet, C1.S
## All levels converged < 1.1: FALSE
## 1.476484; 1.192773; 2.010944; 2.162236; 3.09478; 1.884597; 9.821599; 4.90655
## Individual models saved as: test_c_run-tree1_chain*.rda
## Convergence diagnosis saved as: test_c_run-tree1_conv.rda
##
## 2017-02-16 - 08:51:52: MCMCglmm performed on tree 2
## Convergence diagnosis:
## Effective sample size is > 1000: FALSE
## 9; 9; 9; 2.972582; 9; 9; 9; 24.82875; 9; 9; 9; 9; 9; 9; 9; 9; 9; 9; 9
## C1.Sol.(Intercept), C1.Sol.diet.typeherbivore, C1.Sol.diet.typeomnivore, C1.Sol.diet.typepellet, C1.S
## All levels converged < 1.1: FALSE
## 8.798069; 1.625256; 1.645504; 1.692321; 34.84772; 5.752168; 6.981101; 3.812354
## Individual models saved as: test_c_run-tree2_chain*.rda
## Convergence diagnosis saved as: test_c_run-tree2_conv.rda
##
## 2017-02-16 - 08:51:52: MCMCglmm successfully performed on 2 trees.
## Total execution time: 0.464828 secs.
## Use read.mulTree() to read the data as 'mulTree' data.
## Use summary.mulTree() and plot.mulTree() for plotting or summarizing the 'mulTree' data.

imputeSider now runs the selected amount of chains (no.chains) for each of the sampled phylogeny and
exports the resulting MCMC chains to the working directory. Hence in this case we have run two chains for
each of the two sampled trees so there should be four files in your working directory ending with something
like run-tree1_chain2.rda.

```

```
list.files()
```

```
## [1] "Introduction-to-SIDER.html" "Introduction-to-SIDER.Rmd"
## [3] "Untitled.Rmd"
```

These are the full MCMCglmm model outputs for each of the chains run and can be imported back into the R for full inspection if required using the `read.mulTree` function. Notice also that the two other files ending with something similar to `run-tree2_conv`. These give a description of the convergence diagnostics of the chains for each tree. The results of these diagnostics for each tree are also printed in R after running `tefMcmcgmm` returning whether the Effective sample size exceeded ESS for all estimated parameters (with the number for each parameter given as a list), and whether all chains converged for each parameter.

All these models can be separately imported into R using `read.mulTree()`. However, since we are only interested in the estimated TDF for our species `tefMcmcgmm` only imports the imputed posterior distribution of the estimated TDF for our species.

The contents of `TDF_est.c$tdf_global` is the posterior estimate of the imputed TDF aggregated across all

chains in the model run (TDF_est.c\$tdf_estimates is a list containing each posterior chains separately). As a mcmc class object, the many functions of the coda package provide options for summarising and plotting the distribution.

```
# Explore the names of the list created by SIDER::imputeSider
names(TDF_est.c)
```

```
## [1] "tdf_estimates" "tdf_global"
```

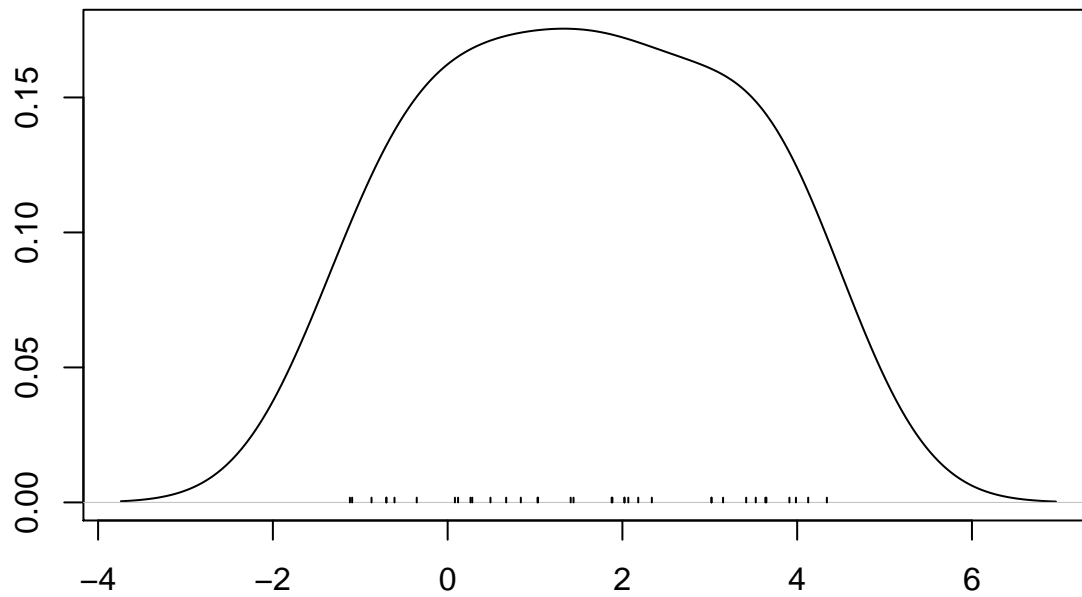
```
# Calculate summary statistics of the posterior.
# Specifically, the mean and standard deviation would be
# taken from here and used in a mixing model analysis using
# MixSIAR, MixSIR or SIAR for example.
summary(TDF_est.c$tdf_global)
```

```
##
## Iterations = 1:36
## Thinning interval = 1
## Number of chains = 1
## Sample size per chain = 36
##
## 1. Empirical mean and standard deviation for each variable,
##    plus standard error of the mean:
##
##           Mean           SD      Naive SE Time-series SE
##           1.5661          1.6880          0.2813          0.6398
##
## 2. Quantiles for each variable:
##
##    2.5%    25%    50%    75%   97.5%
## -1.096  0.226  1.659  3.053  4.153
```

```
# Credible intervals and the mode of the posterior are obtained
# using the hdrdcde package
hdrdcde::hdr(TDF_est.c$tdf_global, prob = c(50, 95, 99))
```

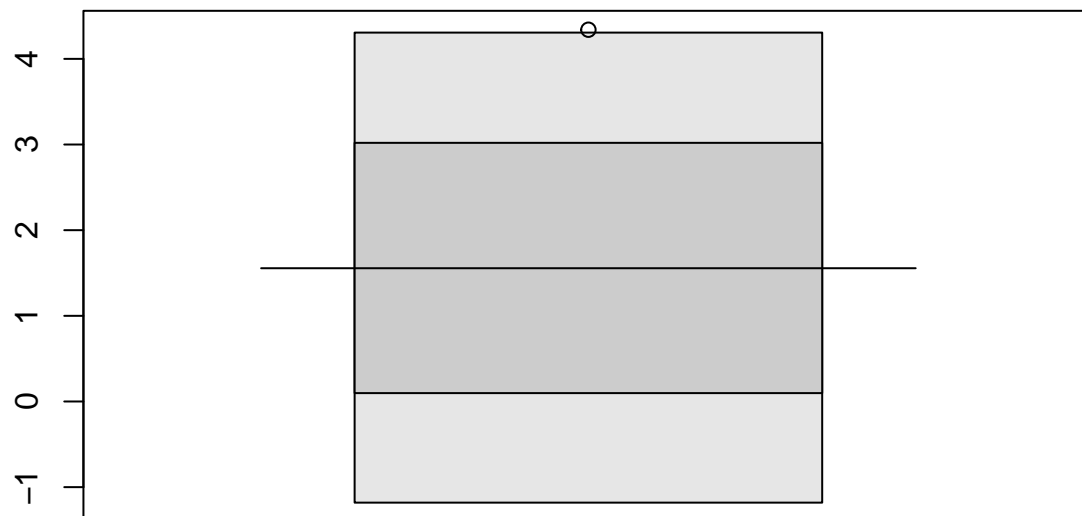
```
## $hdr
##           [,1]      [,2]
## 99% -1.1800431  4.307061
## 95% -1.0993383  4.226152
## 50%  0.1047338  3.019592
##
## $mode
## [1] 1.545137
##
## $alpha
##           1%           5%           50%
## 0.08683256 0.08861431 0.11141480
```

```
# You can also create density plots of the posterior
coda::densplot(TDF_est.c$tdf_global)
```



N = 36 Bandwidth = 0.8738

```
# Or SIAR / SIBER style density boxplots using the hdrnde package
hdrnde::hdr.boxplot(TDF_est.c$tdf_global)
```



```
# return the working directory to the original location (away from the tempdir())
# NB only for use in this vignette example. You probably do not need to
# do this for your own analyses.
setwd(origwd)
```