

Running MiSSE

Jeremy M. Beaulieu

Background

In Beaulieu and O’Meara (2016) we pointed out that the trait-independent HiSSE model is basically a model for traits and a separate model for shifts in diversification parameters, much like BAMM (though without priors, discontinuous inheritance of extinction probability, or other mathematical foibles). The hidden states can drive different diversification processes, and the traits just evolve under a regular Markovian trait model. At that point, there is no harm in just dropping the trait altogether and just focusing on diversification driven by unknown factors. This is what we meant by our HiSSE framework essentially forming a continuum from a purely trait-independent model (e.g., BAMM, or MEDUSA), to a completely trait-dependent model (e.g., BiSSE)(see discussion in Caetano et al., 2018). That is what this MiSSE function does – it sets up and executes a completely trait-free version of a HiSSE model. Thus, all that is required is a tree. The model allows up to 26 possible hidden states in diversification (denoted by A-Z). Transitions among hidden states are governed by a single global transition rate, q . A “shift” in diversification denotes a lineage tracking some unobserved, hidden state. An interesting byproduct of this assumption is that distantly related clades can actually share the same discrete set of diversification parameters.

Note that we refer to “hidden state” simply as a shorthand. We do not mean that there is a single, discrete character that is solely driving diversification differences. There is some heritable “thing” that affects rates, such as a combination of body size, oxygen concentration, trophic level, and, say, how many total species are competing for resources in an area. In other words, it could be that there is some single discrete trait that drives everything. However, it is more likely that a whole range of factors play a role, and we just slice them up into discrete categories, the same way we slice up mammals into carnivore / omnivore / herbivore or plants into woody / herbaceous when the reality is more continuous. This is true for HiSSE, but this concept is especially important to grasp for MiSSE.

Setting up a MiSSE model

The set up is similar to other functions in `hisse`, except there is no need to set up a transition model For the following example, we will use the cetacean phylogeny (e.g., whales and relatives) of Steeman et al. (2009).

```
suppressWarnings(library(hisse))

## Loading required package: ape
## Loading required package: deSolve
## Loading required package: GenSA
## Loading required package: subplex
## Loading required package: nloptr

phy <- read.tree("whales_Steemanetal2009.tre")
```

As with `hisse`, rather than optimizing λ_i and μ_i separately, `MiSSE` optimizes transformations of these variables. We let $\tau_i = \lambda_i + \mu_i$ define “net turnover”, and we let $\epsilon_i = \mu_i/\lambda_i$ define the “extinction fraction”. This reparameterization alleviates problems associated with over-fitting when λ_i and μ_i are highly correlated, but both matter in explaining the diversity pattern (see discussion of this issue in Beaulieu and O’Meara 2016). The number of free parameters in the model for both net turnover and extinction fraction are specified as index vectors provided to the function call. First, let us fit a single rate model:

```
turnover <- c(1)
eps <- c(1)
one.rate <- MiSSE(phy, f=1, turnover=turnover, eps=eps)
```

Pretty simple. Now to fit a model that contains two rate classes, we will simply expand out the turnover vector:

```
turnover <- c(1,2)
eps <- c(1,1)
two.rate <- MiSSE(phy, f=1, turnover=turnover, eps=eps)
```

Overall, MiSSE allows up to 26 possible hidden states in diversification (denoted by A-Z), and in this example since we fit two rate classes, we have two hidden states, A and B, impacting turnover rates. Here is the rest of the model set applied to the whale data set:

```
#rate classes A:C
turnover <- c(1,2,3)
eps <- c(1,1,1)
three.rate <- MiSSE(phy, f=1, turnover=turnover, eps=eps)
#rate classes A:D
turnover <- c(1,2,3,4)
eps <- c(1,1,1,1)
four.rate <- MiSSE(phy, f=1, turnover=turnover, eps=eps)
#rate classes A:E
turnover <- c(1,2,3,4,5)
eps <- c(1,1,1,1,1)
five.rate <- MiSSE(phy, f=1, turnover=turnover, eps=eps)
```

We can also let ϵ_i vary across the tree:

```
turnover <- c(1,2)
eps <- c(1,2)
two.rate.weps <- MiSSE(phy, f=1, turnover=turnover, eps=eps)

#rate classes A:C, but include eps as well:
turnover <- c(1,2,3)
eps <- c(1,2,3)
three.rate.weps <- MiSSE(phy, f=1, turnover=turnover, eps=eps)
```

However, in the case of cetaceans, allowing extinction fraction to vary does not provide any additional information to the model as all rate classes return nearly identical estimates (not shown).

I have already fit these models. Figure 1 shows the improvement in AIC as we increase the complexity of the model.

Plotting MiSSE reconstructions

Like with all other functions, we provide plotting functionality with `plot.misse.states()` for hidden state reconstructions of `class misse.states` output by our `MarginReconMiSSE()` function. And, as with other functions, a single `misse.states` object can be supplied and the plotting function will provide a heat map of the diversification rate parameter of choice, or a list of `misse.states` objects can be supplied and the function will “model-average” the results. For plotting rates, users can choose among turnover, net diversification (“net.div”), speciation, extinction, or extinction fraction (“extinction.fraction”). Below is an example of how to run the reconstruction function to obtain `misse.states` output from our two rate model for cetaceans. But, again, for simplicity, I have a file that contains the reconstructions and we can check that everything has loaded correctly and is of the proper `misse.states` class:

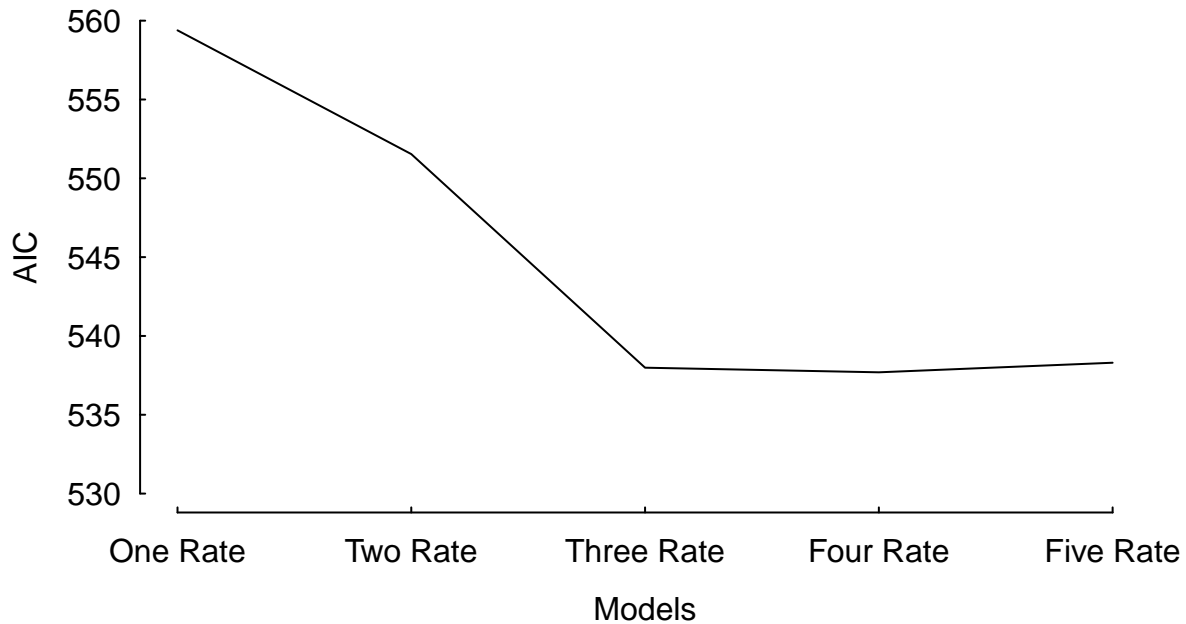


Figure 1: The fit of an incremental increase in the number of rate classes estimated under a MiSSE analysis of the cetacean phylogeny of Steeman et al. (2009). There is a clear reduction in AIC from one to three rate classes, which levels off at four rate classes and five rate classes returns an AIC that is about 1 unit higher than either the three and four rate class.

```
# two.rate.recon <- MarginReconMiSSE(phy=phy, f=1, hidden.states=2,
#pars=two.rate$solution, n.cores=3, aic=two.rate$AIC)
load("misse.vignette.Rsave") # Line above shows the command to create this result.
class(two.rate.recon)
```

```
## [1] "misse.states"
```

```
two.rate.recon
```

```
##
```

```
## Phylogenetic tree with 87 tips and 86 internal nodes.
```

```
##
```

```
## Tip labels:
```

```
## Balaena_mysticetus, Eubalaena_australis, Eubalaena_glacialis, Eubalaena_japonica, Caperea_marginata
```

```
## Node labels:
```

```
## 1, 1, 1, 1, 1, 1, ...
```

```
##
```

```
## Rooted; includes branch lengths.
```

Let's take a look at the reconstruction for the `two.rate` model reconstruction. I will simply supply the reconstruction object if `misse.states` class to the plotting function, `plot.misse.states()`, and plot net diversification (see Figure 2).

```
plot.misse.states(two.rate.recon, rate.param="net.div", show.tip.label=TRUE, type="phylogram",
fs=.25, legend="none")
```

```
## $rate.tree
```

```
## Object of class "contMap" containing:
```

```
##
```

```
## (1) A phylogenetic tree with 87 tips and 86 internal nodes.
```

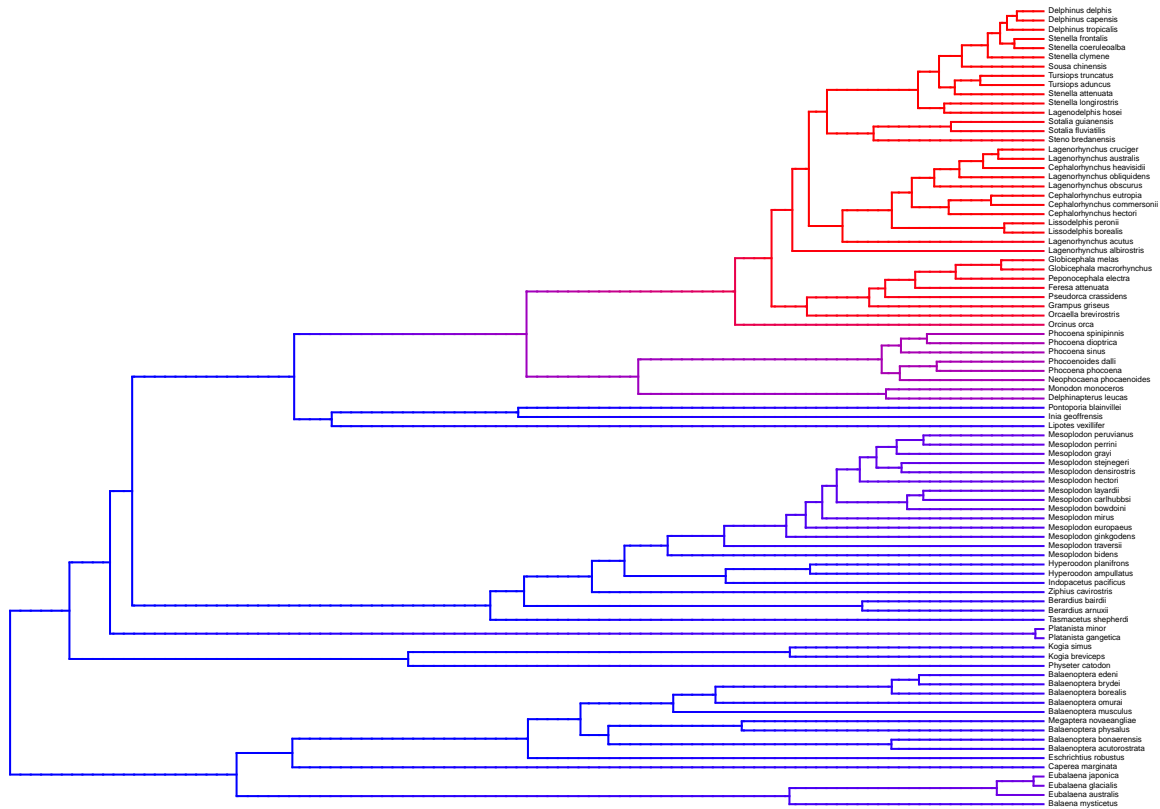


Figure 2: A two-rate class MiSSE analysis and reconstruction of the cetacean phylogeny of Steeman et al. (2009) shows a clear increase in the net diversification rate within the Delphinidae (dolphins) relative to all other cetaceans; there also seems to be a slightly elevated rates in the sister group of Delphinidae, the Monodontidae+Phocoenidae. Overall, this particular MiSSE model seems to correctly identify the source of ‘trait-independent’ diversification that can plague BiSSE analyses of simulated data sets on the cetacean tree (see Rabosky and Goldberg, 2015).

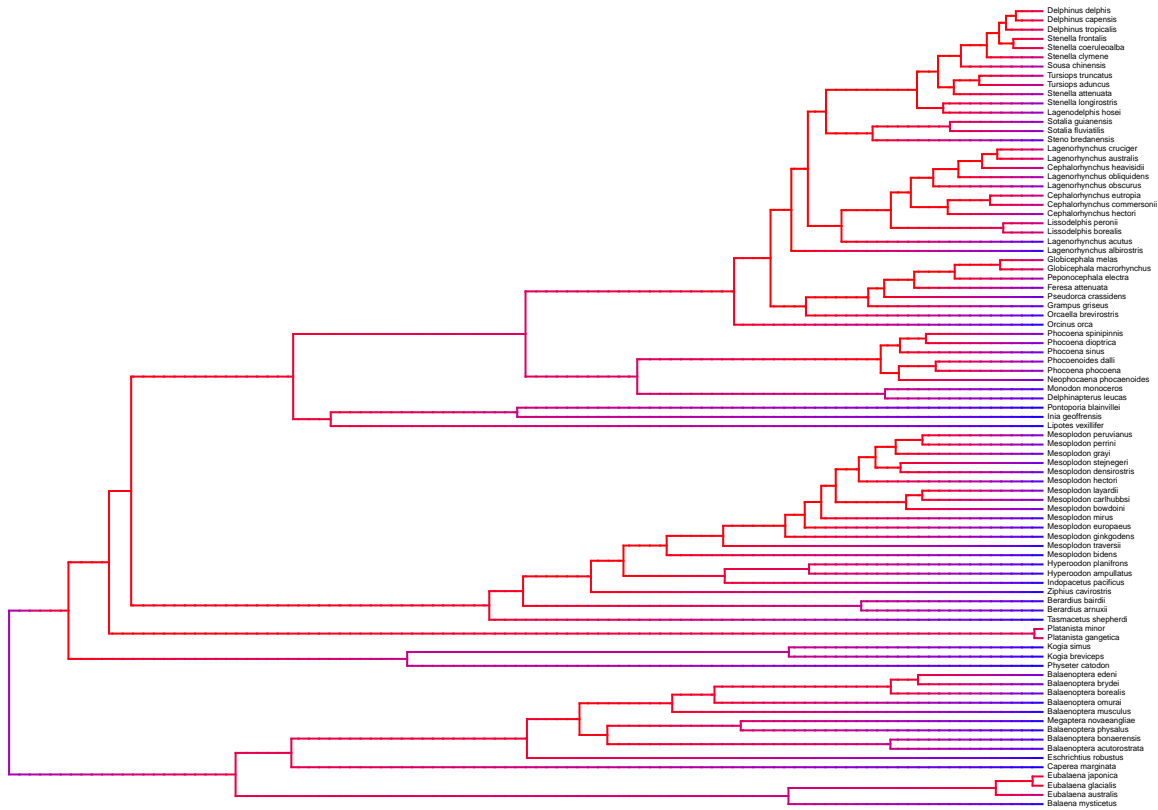


Figure 3: A model-averaged MiSSE analysis of the cetacean phylogeny of Steeman et al. (2009) shows an apparent slow down in the net diversification through time.

```
##
## (2) A mapped continuous trait on the range (0.074823, 0.205397).
```

Of course, in this example we have a set of models that includes models that contain upwards of five rate classes. Also, we know from AIC that there are three model – the three-, four-, and five rate class models – that are within 2 AIC units away from each. As we do in all other `hisse` functions, we allow for plots to be model-averaging of the rates at nodes and tips. To do this, we simply supply all reconstructions you want to average as elements in a list. There are many ways to generate a list, but here is one way:

```
misse.results.list = list()
misse.results.list[[1]] = one.rate.recon
misse.results.list[[2]] = two.rate.recon
misse.results.list[[3]] = three.rate.recon
misse.results.list[[4]] = four.rate.recon
misse.results.list[[5]] = five.rate.recon
```

And, as before, we simply supply this list to the plotting function, `plot.misse.states()`, and plot net diversification (see Figure 3).

```
plot.misse.states(misse.results.list, rate.param="net.div", show.tip.label=TRUE, type="phylogram",
                 fsize=.25, legend="none")
```

```
## $rate.tree
## Object of class "contMap" containing:
##
## (1) A phylogenetic tree with 87 tips and 86 internal nodes.
```

```
##
## (2) A mapped continuous trait on the range (0.047889, 0.229038).
```

Other considerations

Like with `hisse`, `GeoHiSSE`, and `MuHiSSE`, there are functions available for generating estimates of the uncertainty in the parameter estimates (i.e., `SupportRegionMiSSE()`), and to obtain model averages (i.e., `GetModelAveRates()`) for nodes, for tips, or for both to be used in post-hoc tests. Users are encouraged to read other vignettes and help pages provided for more information. For more conceptual discussions of these functions and ideas, readers are also encouraged to read Caetano et al. (2018).

There are two additional items that are worth mentioning. First, like with `MuHiSSE`, I would recommend users try multiple random starting points when optimizing any given model with `MiSSE`. In Nakov et al. (2018), we found that the default starting values often did not return the highest log likelihood. To alleviate this issue, we performed ≥ 50 maximum likelihood optimizations for each model, each initiated from a distinct starting point. All functions within `hisse` are provided with `starting.vals` option for these purposes.

Second, we note that `MiSSE` may seem slower than most other functions within `hisse`. This is somewhat intentional. Underneath the hood we have implemented a lot of checks to the integration for calculating probabilities along branches. This will mean that often times weird messages will spit out to the screen. For now, ignore them, the optimization “feels” these issues and takes necessary action. But this also means that users must pay particular attention to the complexity of the models they are fitting and critically think whether or not the parameters make sense. For example, in cetacean tree used above, I attempted to fit a model with three, four, and five hidden states, A, B, and C, but the reconstructions indicated a rather complicated, and highly uncertain, diversification history:

```
two.rate
```

```
##
## Fit
##           lnL           AIC           AICc           n.taxa
##      -271.7681       551.5363       552.0241       87.0000
## n.hidden.states
##           2.0000
##
## Model parameters:
##
## turnover0A      eps0A      turnover0B      eps0B      q0
## 7.489769e-02 2.061154e-09 2.052022e-01 2.061154e-09 3.743782e-03
```

```
three.rate
```

```
##
## Fit
##           lnL           AIC           AICc           n.taxa
##      -263.9931       537.9863       538.7270       87.0000
## n.hidden.states
##           3.0000
##
## Model parameters:
##
## turnover0A      eps0A      turnover0B      eps0B      turnover0C      eps0C
## 0.008069096 0.275802793 0.008083110 0.275802793 0.340152955 0.275802793
##           q0
## 0.131447160
```

Note that the likelihood was a significant improvement from the two rate model, but the transition rate, q , is roughly two orders of magnitude higher ($q = 0.131$) relative to the two-rate model estimate ($q = 0.004$). The same is true for the four and five rate class models. What does this mean? Well, we can convert this into the expected number of transitions by multiplying the rate by the sum of the branch lengths in the cetacean phylogeny:

```
expected.transitions.two <- 0.004 * sum(two.rate$phy$edge.length)
expected.transitions.two
```

```
## [1] 3.281109
```

```
expected.transitions.three <- 0.131 * sum(three.rate$phy$edge.length)
expected.transitions.three
```

```
## [1] 107.4563
```

At first glance, this might seem that something is off with the three-rate class model – e.g., 3 vs. 107 number of shifts among the different classes? However, examining the support region around the parameters can give an indication as to the overall reliability of these estimates:

```
load("misse.support.Rsave")
two.rate.support$ci[, "q0"]
```

```
##           0%           25%           50%           75%           100%
## 0.0009535839 0.0036835840 0.0065180114 0.0106515654 0.0206332391
```

```
three.rate.support$ci[, "q0"]
```

```
##           0%           25%           50%           75%           100%
## 0.06839226 0.11308939 0.13144716 0.15736921 0.22772718
```

In the case of the two rate class model, even though the MLE indicates very few transitions among the rate classes, there is a model within 2 log likelihood units that suggests as many as 17 expected transitions. In other words, the q estimate even under the two rate class model is fairly certain. We suspect that what is being picked up by these models is something that is both clade-specific (i.e., implied by the two rate class model) and time-dependent (i.e., implied by the three, four, and five rate class models), with the latter exerting the strongest influence on the overall fit.

References

- Beaulieu, J.M., and B.C. O’Meara. (2016). Detecting hidden diversification shifts in models of trait-dependent speciation and extinction. *Syst. Biol.* 65:583-601.
- Caetano, D.S., B.C. O’Meara, and J.M. Beaulieu. 2018. Hidden state models improve state-dependent diversification approaches, including biogeographic models. *Evolution*, 72:2308-2324.
- Nakov, T., Beaulieu, J.M., and Alverson, A.J. 2018. Freshwater diatoms diversify faster than marine in both planktonic and benthic habitats. *bioRxiv*, doi: <https://doi.org/10.1101/406165>.
- Rabosky, D.L., and E.E. Goldberg. (2015). Model inadequacy and mistaken inferences of trait-dependent speciation. *Syst. Biol.* 64:340-355.
- Steehan, M.E., B. M.B. Hebsgaard, E. Fordyce, S.Y.W. Ho, D.L. Rabosky, R. Nielsen, C. Rahbek, H. Glenner, M.V. Sorensen, and E. Willerslev. 2009. Radiation of extant cetaceans driven by restructuring of the oceans. *Syst. Biol.* 58:573-585.