

Some other R packages for model-based multispecies analysis

Bert van der Veen

Department of Mathematical Sciences, NTNU

Questions so far?



Different packages

I will briefly go through the different model-based ordination packages

I will contrast each package to `gl1vm`

Each of these packages really warrants its own presentation

Examples with Bird counts (from CANOCO 5)

```

Y <- read.csv("../data/birdY.csv", header = TRUE, skip = 1, row.names = 1)
Y[is.na(Y)] <- 0;
Y <- Y[,order(colSums(ifelse(Y==0,0,1)),decreasing=TRUE)] #reorder by freq
Y <- as.matrix(Y)
row.names(Y)<-1:nrow(Y) # for VGAM
X <- read.csv("../data/birdX.csv", header = TRUE, skip = 1, row.names = 1)
X[,c(1:3,5:9)] <- scale(X[,c(1:3,5:9)])
X[, -c(1:3,5:9)] <- data.frame(lapply(X[, -c(1:3,5:9)], as.factor))
  
```

Bayesian Ordination and regression AnaLysis

Methods in Ecology and Evolution



Methods in Ecology and Evolution 2016, 7, 744–750

doi: 10.1111/2041-210X.12514

APPLICATION

BORAL – Bayesian Ordination and Regression Analysis of Multivariate Abundance Data in R

Francis K.C. Hui*

Mathematical Sciences Institute, The Australian National University, Canberra, ACT 0200, Australia

boral

- ▶ The first model-based ordination package for community ecology
- ▶ For unconstrained (or residual) ordination (and JSDM)
- ▶ Based on JAGS (Plummer, 2012)
- ▶ Writes the model to a file, loads it into JAGS, returns results
- ▶ Runs on a single MCMC chain

Features

- ▶ Covariates
- ▶ 4th corner model
- ▶ Row intercepts
- ▶ Structured LVs
- ▶ Natively includes gold-standard residuals

boral

This is boral version 2.0.2.

Please note that as of version 2.0, boral will no longer be regularly maintained and updated. However, if you spot any bugs/typos or have a specific feature requests, please contact the maintainer.

boral: code

```
model <- boral::boral(Y, X, formula.X = ~ Forest + Altit, lv.control=list(
```

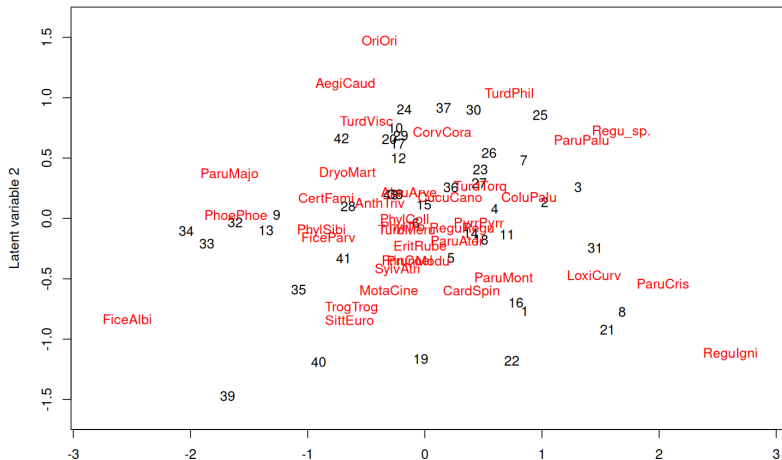
```
## Compiling model graph
##   Resolving undeclared variables
##   Allocating nodes
## Graph information:
##   Observed stochastic nodes: 1591
##   Unobserved stochastic nodes: 270
##   Total graph size: 9883
##
## Initializing model
```

```
boral::lvplot(model)
```

vignette: see paper

boral: plot

Biplot of latent variable posterior medians



boral



Article

Comparison of distance-based and model-based ordinations

David W. Roberts 

First published: 10 October 2019 | <https://doi.org/10.1002/ecy.2908>

Corresponding Editor: Helene H. Wagner.

boral: MCMC

- ▶ MCMC is (kind of) like optimisation, you need to check convergence
- ▶ MCMC needs “burn-in”, i.e., forget the initial state
- ▶ But samples of parameters are stored; so we can expect them
- ▶ MCMC needs to mix well (explore whole parameter space)
- ▶ The chain is stationary if we have reached a good state
- ▶ We can check this visually, or with statistics
- ▶ If it has not converged, it needs to be run longer (or your model is poorly formulated)

boral: calc.varpart

Partition variance per species over model terms

```
boral::calc.varpart(model)
```

```
## $varpart.X
##   EritRube   FrinCoel   ParuAter   PrunModu   SylvAtri   TurdMeru   PhylColl
## 0.22258298 0.70292858 0.05911173 0.55968993 0.43594145 0.24154892 0.16788261
##   CucuCano   PyrrPyrr   TurdTorg   ReguRegu   PhylTro   TrogTrog   PhylSibi
## 0.26558182 0.22193613 0.38170573 0.10209870 0.59288627 0.08752240 0.42714606
##   AnthTriv   TurdPhil   ParuMont   CardSpin   PhoePhoe   ColuPalu   FiceAlbi
## 0.74294135 0.08109488 0.15816026 0.16161584 0.34310262 0.09218656 0.05829110
##   ParuMajo   Regu_sp.   CorvCora   SittEuro   LoxiCurv   ParuCris   ReguIgni
## 0.09077412 0.13154426 0.34919838 0.58920719 0.19911165 0.29188367 0.13806011
##   AlauArve   FiceParv   MotaCine   OriOri   AegiCaud   CertFami   DryoMart
## 0.77674032 0.60020632 0.40306591 0.33714912 0.28856959 0.58945434 0.31248251
##   TurdVisc   ParuPalu
## 0.29422510 0.36944626
##
## $varpart.lv
##   EritRube   FrinCoel   ParuAter   PrunModu   SylvAtri   TurdMeru   PhylColl   CucuCano
## 0.7774170 0.2970714 0.9408883 0.4403101 0.5640586 0.7584511 0.8321174 0.7344182
##   PyrrPyrr   TurdTorg   ReguRegu   PhylTro   TrogTrog   PhylSibi   AnthTriv   TurdPhil
## 0.7780639 0.6182943 0.8979013 0.4071137 0.9124776 0.5728539 0.2570586 0.9189051
##   ParuMont   CardSpin   PhoePhoe   ColuPalu   FiceAlbi   ParuMajo   Regu_sp.   CorvCora
## 0.2412227 0.2822212 0.2522271 0.2372212 0.2172222 0.2222272 0.2222272 0.2222212
```

boral

Has a few other helpful functions:

- ▶ `get.enviro.cor` and `get.residualcor`
- ▶ `predict.boral` and `plot.boral`
- ▶ `coefspplot` and `ranefspplot`

boral: compared to gl1vm

boral	gl1vm
Bayesian	Frequentist
MCMC	Likelihood approximation
Slow	Fast
Correlated LVs	Correlated LVs
Single row effect	Multiple row effects
Stochastic Variable Selection	Adaptive shrinkage?

There is little reason to use boral at this point, except for the SVSS.



Hierarchical Modeling of Species Communities

APPLICATION

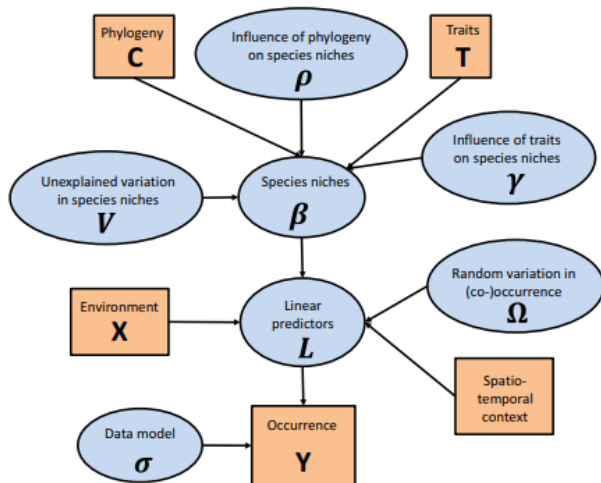
Methods in Ecology and Evolution

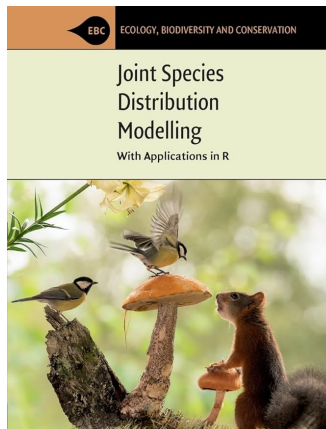


Joint species distribution modelling with the R-package HMSC

Gleb Tikhonov^{1,2} | Øystein H. Opedal^{2,3}  | Nerea Abrego⁴ | Aleksi Lehikoinen⁵ |
Melinda M. J. de Jonge⁶ | Jari Oksanen⁷ | Otso Ovaskainen^{2,3} 

HMSC





HMSC

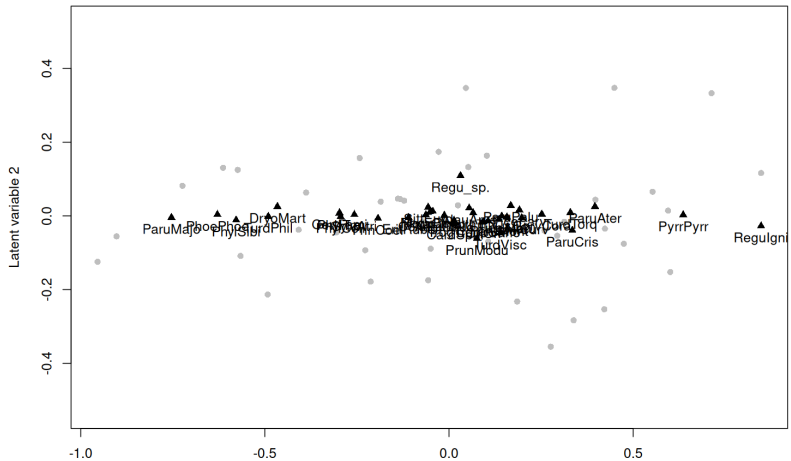
- ▶ Bayesian; fits with MCMC
- ▶ Custom Gibbs samplers
- ▶ Flexible package for multispecies hierarchical modeling
- ▶ Focuses on prediction and species associations
- ▶ Phylogenetic effects
- ▶ Efficiently implements spatial models with nearest neighbors
- ▶ 4th corner model
- ▶ Various extra random effects (intercepts and such)
- ▶ Effects can be specified at different sampling levels, including sets of LVs
- ▶ The “infinite factor model”
- ▶ Have a preprint on parallelisation
- ▶ Very little support for ordination
- ▶ Supports mixed response types

HMSC: code

need to set-up LVs

```
studyDesign = data.frame(sample=as.factor(1:nrow(Y)))
rL <- Hmsc::HmscRandomLevel(units = studyDesign$sample)
model <- Hmsc::Hmsc(Y, XFormula = ~Forest+Altit, XData= X,
distr = "lognormal poisson", studyDesign = studyDesign,
ranLevels = list(sample = rL))
# Run mcmc
run = Hmsc::sampleMcmc(model, samples = 1000, nChains = 3,
transient = 2500)
```

```
## Computing chain 1
## Chain 1, iteration 20 of 3500 (transient)
## Chain 1, iteration 40 of 3500 (transient)
## Chain 1, iteration 60 of 3500 (transient)
## Chain 1, iteration 80 of 3500 (transient)
## Chain 1, iteration 100 of 3500 (transient)
## Chain 1, iteration 120 of 3500 (transient)
```



HMSC

HMSC	gllvm
Bayesian	Frequentist
MCMC	Likelihood approximation
Slow (but parallel package)	Fast
normal, Bernoulli, Poisson, lognormal Poisson	Wide range of response types
Different effects at different sampling levels	One effect at a sampling level
Infinite factor model	Number of LVs fixed <u>a-priori</u>
Efficient spatial implementation	Spatial is a work in progress
Few tools for ordination	many tools for ordination

Ultimately, the focus of these two packages is very different. HMSC focuses on prediction and JSDMs, gllvm can do that, but its main focus is different (IMO).

ecoCopula

RESEARCH ARTICLE



Fast model-based ordination with copulas

Gordana C. Popovic¹ | Francis K. C. Hui² | David I. Warton¹

- ▶ Employs graphical models for determining species associations
- ▶ Requires a secondary model
- ▶ Is -very- fast for ordination (faster than NMDS!)
- ▶ Can estimate “direct associations” (not as quick)
- ▶ Supports mixed response types

ecoCopula: code

```
preModel <- ecoCopula::stackedsdm(Y, formula_X =~1, data = X)
model <- ecoCopula::cord(preModel)
plot(model, biplot=TRUE)
```

vignette

ecoCopula

ecoCopula	gllvm
Frequentist	Frequentist
Gaussian Copula	Likelihood approximation
Faster	Fast
A decent number of distributions	Wide range of response types
Direct species associations	Correlative
None	Many other random effects
Secondary model in parallel	Working on parallel computation
Native residuals	Native residuals
Biplot function	Biplot function
Marginal interpretation	Conditional interpretation

ecoCopula has a lot of potential due to its speed, but lacks in support, maintenance, and perhaps some maturity.

Vector Generalised Linear and Additive Models



Journal of Statistical Software

January 2010, Volume 32, Issue 10.

<http://www.jstatsoft.org/>

The VGAM Package for Categorical Data Analysis

Thomas W. Yee
University of Auckland

Vector Generalised Linera and Additive Models

- ▶ Package with a wide range of model types **VGLMs**
- ▶ Massive package with a lot of functionality
- ▶ An incredible range of response distributions
- ▶ Unconstrained and constrained ordination (fixed effects formulation)
- ▶ Quadratic and additive ordinations
- ▶ Supposed to fit quickly with IWLS
- ▶ In my experience, fitting is often difficult (errs often) and can be unstable
- ▶ Has some residuals
- ▶ Plotting functions are a bit different
- ▶ No random effects
- ▶ Now (recently) has doubly-constrained ordination!

Centers around `vglm()`, `vgam()`, `rrvglm()`, `cqo()`, `cao()`, `rcim()`

VGAM

The first (model-based) constrained ordination method

Ecological Monographs, 74(4), 2004, pp. 685–701
© 2004 by the Ecological Society of America

A NEW TECHNIQUE FOR MAXIMUM-LIKELIHOOD CANONICAL GAUSSIAN ORDINATION

THOMAS W. YEE¹

*Department of Statistics, University of Auckland, Private Bag 92019, Auckland, New Zealand, and
Department of Statistics and Applied Probability, 6 Science Drive 2, National University of Singapore, Singapore 117546*

VGAM: code

```

model1 <- VGAM::rcim(Y, Rank = 2, family = VGAM::poissonff)
VGAM::lvplot(model1)

```

```

# Could not get this to work :(
# model2 <- VGAM::rruglm(Y ~ model.matrix(~.,X[,1:4])[, -1], Rank = 2, fami

```

vignette: see reference card

VGAM: plot

VGAM

VGAM	gllvm
Frequentist	Frequentist
ML via IWLS	Approximate marginal likelihood
Fast	Fast
Incredible range of responses	Wide range of response types
Not robust fitting	Relatively robust
No random effects	Many other random effects
UQO, CQO, CAO	UQO, CQO
VGAMs	No smooths
Native residuals	Native residuals
Biplot function	Biplot function

VGAM has a lot of potentially useful tools, but I do not find it very usable.

glmmTMB

glmmTMB Balances Speed and Flexibility Among Packages for Zero-inflated Generalized Linear Mixed Modeling

by Mollie E. Brooks, Kasper Kristensen, Koen J. van Benthem, Arni Magnusson, Casper W. Berg,
Anders Nielsen, Hans J. Skaug, Martin Mächler, Benjamin M. Bolker




Journal of Statistical Software

March 2025, Volume 112, Issue 1.


doi: 10.18637/jss.v112.i01

Parsimoniously Fitting Large Multivariate Random Effects in glmmTMB

Maeve McGillicuddy 
UNSW Sydney

Gordana Popovic 
UNSW Sydney

Benjamin M. Bolker 

David I. Warton 

glmmTMB: functionality

- ▶ Kind of similar to gl1vm in that it uses approximate methods
- ▶ Laplace approximation with TMB (state-of-the-art)
- ▶ Great usability
- ▶ Can include many random-effects
- ▶ (Un)constrained random effects ordination
- ▶ Slower than gl1vm's VA (usually)
- ▶ Structured random effects (e.g., spatial), Phylogenetic intercepts
- ▶ Little support for ordinations
- ▶ Big on zero-inflated modelling

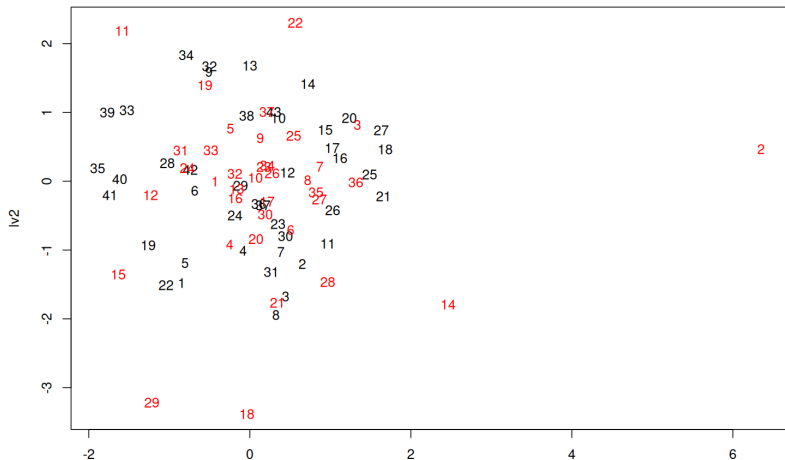
glmmTMB: code

Note: data needs to be in long format

```
# organize data into long format
tmp <- data.frame(Y)
tmp$id <- 1:nrow(tmp)
glmmDat <- reshape(tmp,
                    idvar = "id",
                    timevar = "col",
                    times = colnames(Y),
                    varying = list(colnames(Y)),
                    v.names = "y",
                    direction = "long")

model <- glmmTMB::glmmTMB(y ~ col + rr(col + 0|id, d = 2), data = glmmDat,
rrstuf <- extract_rr(model) # see here: https://github.com/glmmTMB/glmmTMB
plot(rbind(rrstuf$b,rrstuf$f1), type = "n");
text(rrstuf$b);text(rrstuf$f1, col = "red")
```

glmmTMB: plot



glmmTMB

glmmTMB	gllvm
Frequentist	Frequentist
Laplace	VA (default) or Laplace
Fast	Fast(er)
Wide range of response types	Wide range of response types
Many (structured) random effects	Many random effects
Can also fit with MCMC	No
Zero-inflated modeling	Work in progress
No residuals	Native residuals
No plotting function	Biplot function
Large community	Small community
Excellent developers	No comment :)

glmmTMB is especially useful if you want user friendliness and many other random effects.

Generalized Matrix Factorization

Generalized Matrix Factorization: efficient algorithms
for fitting generalized linear latent variable models
to large data arrays

Lukasz Kidziński

*Department of Bioengineering
Stanford University
Stanford, CA 94305, USA*

LUKASZ.KIDZINSKI@STANFORD.EDU

Francis K.C. Hui

*Research School of Finance, Actuarial Studies and Statistics
The Australian National University
Canberra, ACT 2601, Australia*

FRANCIS.HUI@ANU.EDU.AU

David I. Warton

*School of Mathematics and Statistics
and Evolution & Ecology Research Centre
The University of New South Wales
Sydney, NSW 2052, Australia*

DAVID.WARTON@UNSW.EDU.AU

Trevor Hastie

*Department of Statistics and Biomedical Data Science
Stanford University
Stanford, CA 94305, USA*

HASTIE@STANFORD.EDU

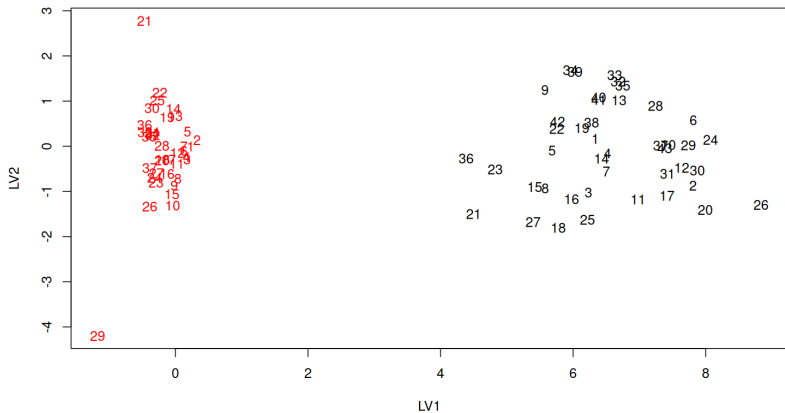
- ▶ Very quick; fits by penalized likelihood
- ▶ Unconstrained or residual ordination only
- ▶ No extra random-effects
- ▶ Can be unstable due to the approximation
- ▶ Stale package not on CRAN

gmf: code

```

# devtools::install_github("kidzik/gmf")
model <- gmf::gmf(Y, family = poisson(), p = 2)
plot(rbind(model$u,model$v), type = "n", xlab="LV1", ylab="LV2")
text(model$u)
text(model$v, col="red")
  
```

gmf: plot



gmf

gmf	gllvm
Frequentist	Frequentist
Penalized likelihood	VA or LA approximation
Fast(er) but unstable	Fast
A few response types	Wide range of response types
Fitting is fine	Relatively robust
No random effects	Many other random effects

A skeleton of a package, not very useful at this point.

RCM

A unified framework for unconstrained and constrained ordination of microbiome read count data

Stijn Hawinkel^{1*}, Frederiek-Maarten Kerckhof², Luc Bijmans^{3,4}, Olivier Thas^{1,4,5}

1 Department of Data Analysis and Mathematical Modelling, Ghent University, Ghent, Belgium, 2 Center for Microbial Ecology and Technology, Ghent University, Ghent, Belgium, 3 Quantitative Sciences, Janssen Pharmaceutical companies of Johnson and Johnson, Beerse, Belgium, 4 Center for Statistics, Hasselt University, Hasselt, Belgium, 5 National Institute for Applied Statistics Research Australia (NIASRA), University of Wollongong, Wollongong, Australia

* stijn.hawinkel@ugent.be

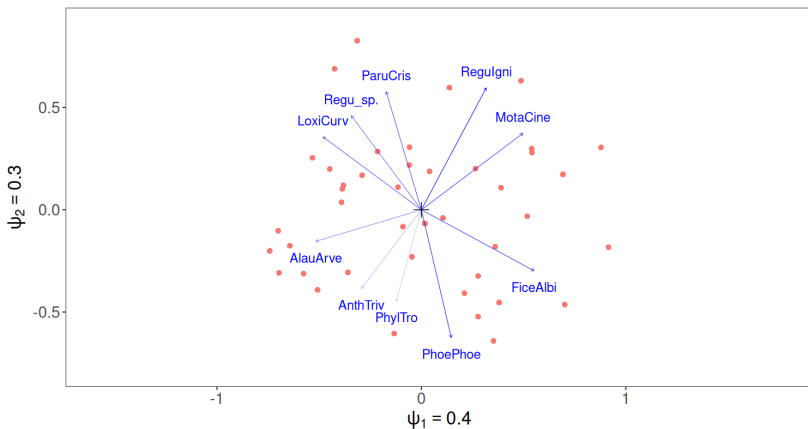
- ▶ Does both unconstrained and unconstrained ordination
- ▶ Even additive constrained ordination
- ▶ All based on fixed effects formulations, no random effects
- ▶ Only the negative binomial distribution
- ▶ Not a “true” statistical model (according to the authors)
- ▶ Permanova functionality
- ▶ Residual plots

RCM: code

```
# devtools::install_github("CenterForStatistics-UGent/RCM")
model <- RCM::RCM(Y, k = 2)
plot(model)
```

vignette

RCM: plot



RCM

RCM	gllvm
Frequentist	Frequentist
Maximum likelihood	Approximate marginal likelihood
Fast	Fast
Only NB	Wide range of response types
UO, CO, CQO, CAO	UO, CO, CQO
No random effects	Many other random effects

RCM seems good at what it does, but functionality is limited.

Community-level basis function models

RESEARCH ARTICLE

Methods in Ecology and Evolution



Spatiotemporal joint species distribution modelling: A basis function approach

Francis K. C. Hui¹ | David I. Warton² | Scott D. Foster³ | Christopher R. Haak⁴

- ▶ Most recent developments: GAM for multiple species
- ▶ Specifically targeted on spatio,temporal or spatio-temporal analysis
- ▶ This is something GLLVMs are not -terribly- good at yet (but very much an area of interest)
- ▶ Based on the idea of LVMs, but not with LVs
- ▶ Fitting using TMB
- ▶ I.e., JSMD-oriented, not ordination

CBFM

CBFM	gllvm
Frequentist	Frequentist
Penalized Quasi-likelihood	Approximate marginal likelihood
For large spatio-temporal problems	Not an option (yet?)
Wide range of response types	Wide range of response types
Post-hoc ordination	Is an ordination method
Can include extra "random effects" as smooths	Many other random effects
Parallelisation	Parallelisation
No traits or Phylogeny	Traits and Phylogeny

Sorry, no example yet. Bird data does not have coordinates, and CBFM only fits models with space it seems?

Summary

Package	cran ¹	UO ²	CO ³	CN ⁴	RE ⁵	CI ⁶	traits	Phylogeny	Space	framework ⁷
glvm	yes	yes	yes	yes	yes	yes	yes	yes	sort of	F
Boral	yes	yes	no	no	some	yes	yes	no	yes	B
HMSC	yes	yes	no	no	yes	yes	yes	yes	yes	B
ecoCopula	yes	yes	no	no	no	no	kind of	no	no	F
VGAM	yes	yes	yes	no	no	some	new?	no	no	F
glmmTMB	yes	yes	no	no	yes	yes	yes	sort of	Kind of	F
gmf	no	yes	no	no	no	no	no	no	no	F
RCM	no	yes	yes	no	no	no	no	no	no	F
CBFM	no	no	no	no	yes	yes	no	no	yes	F

¹cran: Package available on CRAN. ²UO: Unconstrained ordination. ³CO: Constrained. ⁴CN: Concurrent. ⁵RE: Random effects. ⁶CI: Confidence/Credible intervals. ⁷framework: The underlying framework of the model (F: Frequentist, B: Bayesian).

When to use what package?

- ▶ HMSC for extensive support for JSDMs
- ▶ glmmTMB for many (structured) random effects
- ▶ ecoCopula if you have a **huge** dataset and glvm is too slow
- ▶ CBFM for large spatial/temporal models

glvm for all your ordination needs

Summary

New software implementations are continuously being developed.
 Dimension reduction methods for ecology have entered a new era.

Most of existing implementations are based on the GLLVM
 framework (sjSDM, BayesComm, VGAM, RCM excluded)

Summary

New software implementations are continuously being developed.
 Dimension reduction methods for ecology have entered a new era.

Most of existing implementations are based on the GLLVM
 framework (sjSDM, BayesComm, VGAM, RCM excluded)

Summary

- ▶ It is important that we continue to explore new and better methods
- ▶ Especially the application of ordination methods is stuck in the past
- ▶ There is still a lot of work to be done on multivariate methods for community ecology
- ▶ There are more packages for model-based analysis that I have not mentioned
- ▶ E.g., jSDM, sjSDM, BayesComm