

# Generalised Linear (Mixed) Models for data of multiple species

Bert van der Veen

Department of Mathematical Sciences, NTNU

## Species correlation

---

If we fit a GLM to data of multiple species, we assume **independence**

But, observations of the same species form groups. Co-occurring species have more similar observations than for other species

In GLMM language: **observations of species exhibit correlation**

- 1) Part of this can be explained by shared environmental responses
- 2) The other part remains

## The previous model

---

```
model4 <- gllvm::gllvm(y, X = X, formula = ~N03, num.lv = 0,  
                        family = "negative.binomial")
```

If we look at the correlations in the residuals of this model, we can see that.

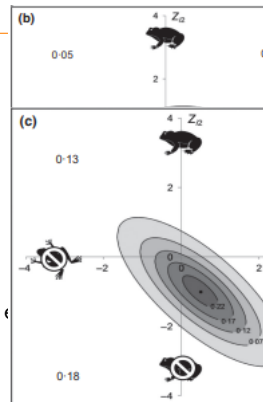
[illegible]



- The goal: to incorporate covariation of species

What induces c

- 1: Pollock e



## Interactions and co-occurrence

### ECOLOGY LETTERS

*Ecology Letters*, (2020) 23: 1050–1063

doi: 10.1111/ele.13525

IDEAS AND

PERSPECTIVES

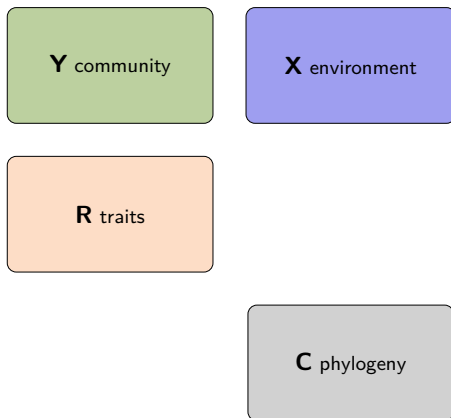
**Co-occurrence is not evidence of ecological interactions**

Interactions induce correlation, but we cannot infer interactions from (non-temporal) co-occurrence data.



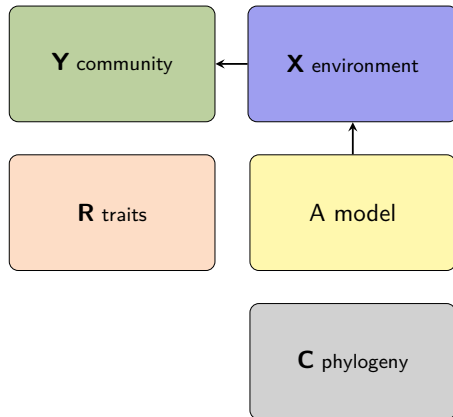
## Typical questions in the framework

---



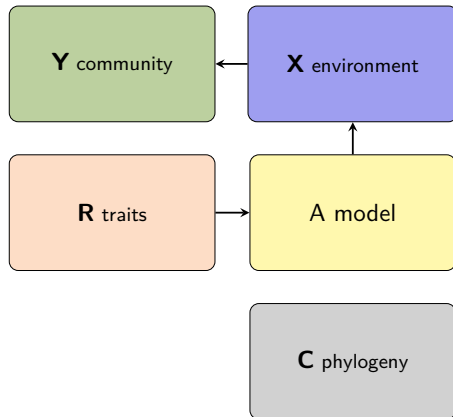
## Typical questions in the framework

---



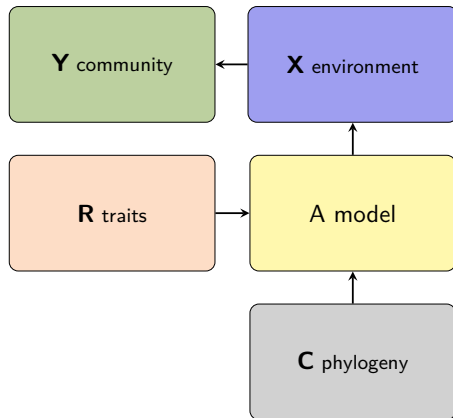
Q: How does the environment structure the community?  
**environmental filtering**

## Typical questions in the framework



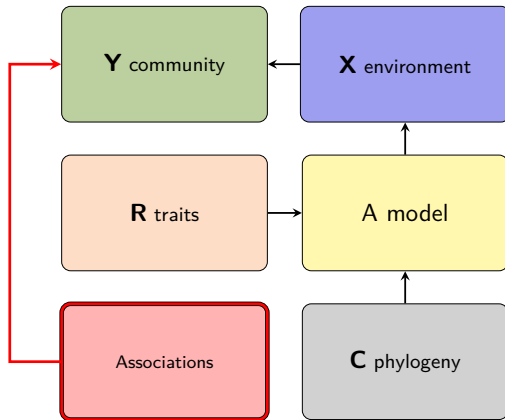
Q: How do traits affect species' responses to the environment?  
**environmental filtering**

## Typical questions in the framework



Q: Do species with shared evolutionary history co-occur?  
 (phylogenetic structuring) **environmental filtering**

## Typical questions in the framework



Q: Do species co-occur **after** the environment has been considered? **biotic filtering**

## Joint Species Distribution

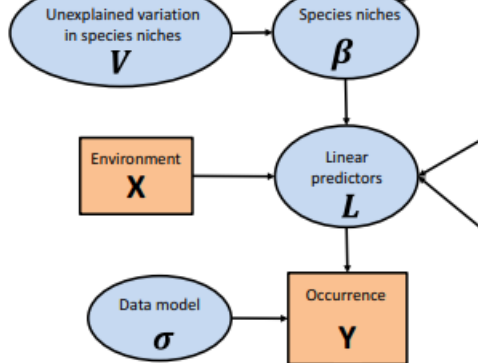


Figure 2: Figure from Ovaskainen et al. (2017)

# Joint Species Distribution Modeling

---

The idea of JSDMs is to incorporate *species associations*

- ▶ Species may co-occur due to biotic interactions
- ▶ Due to similar environmental preferences
- ▶ Or because they have a similar history

Either how, it results in correlations between responses

- ▶ we add  $\epsilon_i$  relative to the VGLM(M)
- ▶ This random effect takes care of the left-over (co)variation of species
- ▶ so we assume  $\epsilon_i \sim \mathcal{N}(\mathbf{0}, \Sigma)$
- ▶  $\Sigma$  is the matrix of **species associations**



## JSDM: the model

---

$$\eta_{ij} = \beta_{0j} + \dots + \epsilon_{ij} \quad (3)$$

- ▶  $\epsilon_i \sim \mathcal{N}(0, \Sigma)$
- ▶  $\Sigma$  is the matrix of *species associations*
- ▶ So we expect a positive values of species co-occur, and negative if they do not

## JSDM: it is a mixed-effects model

The JSDM is “just” a mixed-effects model. So we can fit it with available software:

In lme4:

```
glmer(abundance ~ species + x:species + (0+species|sites), data = data)
```

- ▶ There are  $p(p+1)/2$  correlations between species
- ▶ This model becomes (very) large very quickly
- ▶ Will usually not fit
- ▶ So we need to do something smart!

## Species associations

- ▶ Difficult to estimate: there are usually too many parameters
- ▶ Can only fit this way when there are (much) more sites than species
- ▶ The number of pairwise associations grows quadratically
  - ▶ 2 with 2 species, 6 for 4 species, 45 for 10 species, 4950 for 100

$$\Sigma = \begin{bmatrix} 1 & sp_{12} & \cdots & sp_{1j} \\ sp_{21} & 1 & \cdots & sp_{2j} \\ \vdots & & \ddots & \vdots \\ sp_{j1} & sp_{j2} & \cdots & 1 \end{bmatrix} \quad (4)$$

This very quickly becomes an issue for fitting models

## Ordination to the rescue

---

- ▶ GLLVMs were introduced as a technical solution to this problem
- ▶ We represent the covariance matrix with fewer **dimensions**:  

$$\Sigma \approx \mathbf{\Gamma}\mathbf{\Gamma}^\top$$

“The factor analytic solution” because factor analysis (Spearman, 1904) is the precursor of GLLVMs

## JSDM vs classical multivariate analysis

	Classic	JSDM
Focus	Low-dimensional space	Distributions
Goal	Inference	Prediction
Data type	Usually quantitative	Binary
Scale	Local	Regional
Covariates	Environmental	Bioclimatic
Presentation	Ordination diagram	Correlation plot/map
Audience	Community ecologists	Macro ecologists

## JSDM vs classical multivariate analysis

	Classic	JSDM
Focus	Low-dimensional space	Distributions
Goal	Inference	Prediction
Data type	Usually quantitative	Binary
Scale	Local	Regional
Covariates	Environmental	Bioclimatic
Presentation	Ordination diagram	Correlation plot/map
Audience	Community ecologists	Macro ecologists

That is not to say JSDMs cannot be used for non-binary data, for inference, or for local scales

## JSDM software implementations

---

There are many!

- ▶ Boral (Bayesian, slow and somewhat outdated)
- ▶ sJSDM (Bayesian, relatively slow, but faster than Boral)
- ▶ Hmsc (Bayesian, generally slow, loads of functionality)
- ▶ ecoCopla (Frequentist, very fast but limited functionality)
- ▶ CBFM (Frequentist, geared towards spatio-temporal analysis)
- ▶ sjSDM (Frequentist, very fast but limited functionality, requires python)
- ▶ glmmTMB (Frequentist, fast and very versatile, not purpose-coded)
- ▶ gllym (Frequentist, fast and very versatile, purpose-coded)



- ▶ Boral (Bayesian, slow and somewhat outdated)
- ▶ sJSDM (Bayesian, relatively slow, but faster than Boral)
- ▶ Hmsc (Bayesian, generally slow, loads of functionality)
- ▶ ecoCopla (Frequentist, very fast but limited functionality)
- ▶ CBFM (Frequentist, geared towards spatio-temporal analysis)
- ▶ sjSDM (Frequentist, very fast but limited functionality, requires python)
- ▶ glmmTMB (Frequentist, fast and very versatile, not purpose-coded)
- ▶ gllym (Frequentist, fast and very versatile, purpose-coded)



## Example with Eucalyptus data (Pollock et al. 2015)

```
Y <- read.csv("../data/eucalyptY.csv"),[-1]
X <- read.csv("../data/eucalyptX.csv"),[-1]
X <- data.frame(lapply(X, function(x){if(is.numeric(x)){scale(x)}else{as.factor(x)}}))
knitr::kable(head(X), format="latex", booktabs = T)
```

IDENT	elev	PerSlope	aspect	Northness	Sdclass	Rockiness	VallyBotFlat	topowe
MGS1	0.1972540	1.0236658	-1.1686010	0.3254470	-1.551741	1.0315338	-0.5939667	-0.407019
MGS5	0.1972540	0.7778142	-1.1686010	0.3254470	-1.551741	1.4558834	-0.5939667	-0.407019
MGS4	0.1757452	1.0236658	-0.8419304	-0.5684498	-1.551741	1.3144335	-0.5939667	-0.163936
MGS3	0.1757452	0.4090369	-0.8419304	-0.5684498	-1.426354	0.4657344	-0.5939667	-0.476998
MGS2	0.1327276	0.5319627	-0.3583180	-1.5163833	-1.426354	-0.2415148	-0.5939667	-1.003677
MGN1	0.6381843	1.3924432	0.7306972	0.0144213	-1.426354	1.5973333	-0.5939667	-0.889502

- ▶ 20 species
- ▶ 458 sites
- ▶ Soil covariates and a few bioclimatic

# Eucalyptus: fit a model

```
jsdm1 <- gllvm::gllvm(Y, X = X, formula = ~ Sandiness + cvTemp,
  family = "binomial", num.lv = 2, method = "EVA", starting.val = "zero")
```

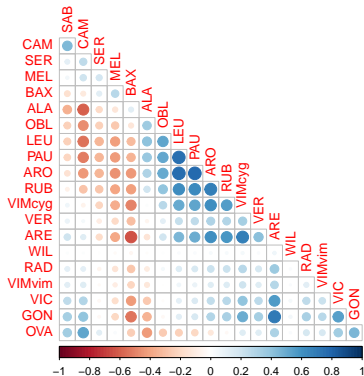
- ▶ method: controls the approximation (LA, VA (default), EVA)
- ▶ starting.val: how to generate initial values (multimodality)
- ▶ n.init: rerun multiple times

The gllvm R-package is fast, but might can take a minute.

Models should be refitted with 'n.init'.

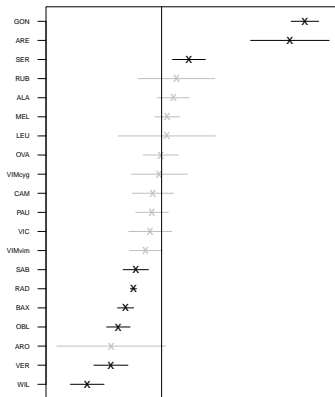
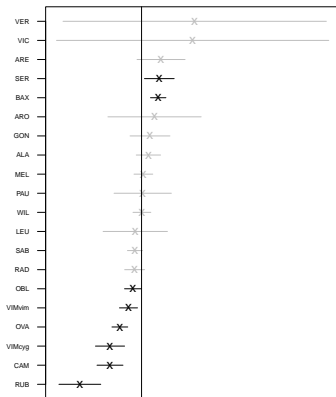
## Eucalyptus: residual associations

```
corrplot::corrplot(gllvm::getResidualCor(jsdm1), order = "AOE", type = "lower", diag = FALSE)
```



## Eucalyptus: environment effects

```
gllvm::coefplot(jsdm1)
```



## Eucalyptus: explained variation

---

```
jsdm2 <- gllvm::gllvm(Y, family = "binomial", num.lv = 2, method = "EVA",  
                     starting.val = "zero")  
gllvm::getResidualCov(jsdm1)$trace/gllvm::getResidualCov(jsdm2)$trace
```

```
## [1] 0.6172353
```

```
jsdm3 <- gllvm::gllvm(Y, X = X, formula = ~(0+Sandiness + cvTemp|1),
  family = "binomial", num.lv = 2,
  method = "EVA", starting.val = "zero")
```



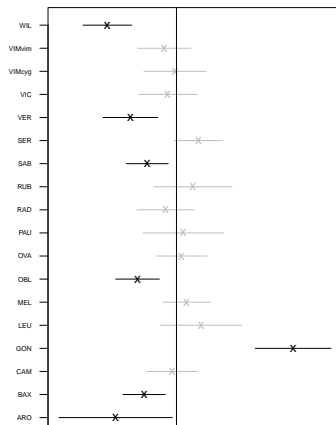
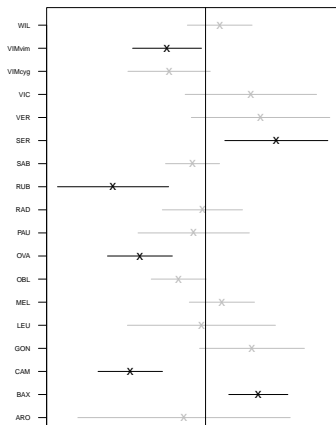
## Eucalyptus: with random effects

```
summary(jsdm3)
```

```
##
## Call:
## gllvm::gllvm(y = Y, X = X, formula = ~(0 + Sandiness + cvTemp |
##      1), family = "binomial", num.lv = 2, method = "EVA", starting.val =
##
## Family:  binomial
##
## AIC:   4486.426 AICc:   4487.341 BIC:   4942.273 LL:   -2179.2 df:   64
##
## Informed LVs:  0
## Constrained LVs:  0
## Unconstrained LVs:  2
##
## Formula:  ~(0 + Sandiness + cvTemp | 1)
```

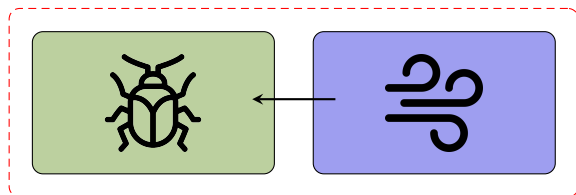
# Eucalyptus: with random effects

```
gllvm::randomCoeefplot(jsdm3)
```



- ▶ **Y**: community data
- ▶ **X**: environmental variables
- ▶ **TR**: species traits

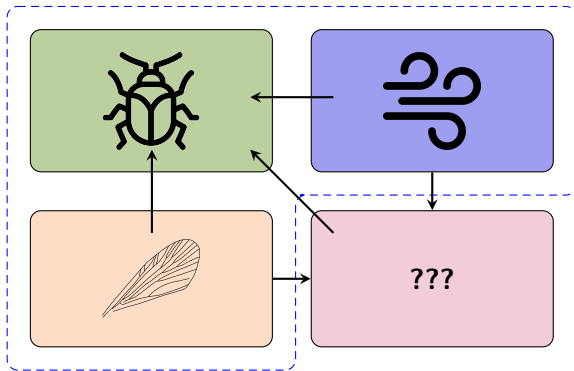
## Fourth corner analysis



Species-environment relationship: the methods

Site-trait relationship: VGLM(M)

## Fourth corner analysis



## Trait-based analysis

1. CWM + RDA *Doleddec et al. (1996)*
2. Double constrained ordination  
*Lebreton et al. (1988), ter Braak et al. (2018)*
3. Fourth corner (LV) Models *Brown et al. (2014), Ovaskainen et al. (2017), Niku et al. (2021)*

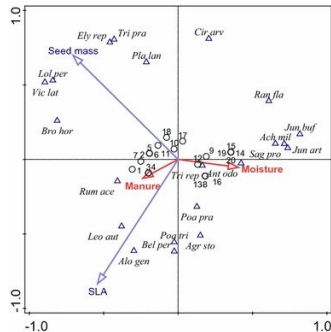


Figure 3: Quadriplot  
*ter Braak et al. (2018)*

## Fourth-corner LVMs

A kind of LVM (JSDM) that also includes traits

Received: 19 February 2020 | Revised: 6 April 2021 | Accepted: 9 April 2021

DOI: 10.1002/env.2683

SPECIAL ISSUE PAPER

WILEY

# Analyzing environmental-trait interactions in ecological communities with fourth-corner latent variable models

Jenni Niku<sup>1</sup> | Francis K. C. Hui<sup>2</sup> | Sara Taskinen<sup>1</sup> | David I. Warton<sup>3</sup>



## Fourth-corner LVMs

---

The model is very similar to before:

$$\eta_{ij} = \beta_{0j} + \mathbf{x}_i^\top \boldsymbol{\beta}_j + \epsilon_{ij} \quad (5)$$

but now, we are hierarchically modelling species' effects  $\beta_j$

$$\eta_{ij} = \beta_{0j} + \mathbf{x}_i^\top (\boldsymbol{\beta}_x + \mathbf{b}_j) + \mathbf{tr}_j^\top \mathbf{B}_{xtr} \mathbf{x}_i \quad (6)$$

- ▶  $\boldsymbol{\beta}_x$  species-common effects
- ▶  $\mathbf{b}_j$  species-specific effects
- ▶  $\mathbf{B}_{xtr}$  4th-corner coefficients



## Fourth-corner LVMs

---

- ▶ With a 4th corner LVM we can examine trait-environment relationships
- ▶ Figure out **why** species prefer particular conditions
  - ▶ In contrast to “just” which conditions they prefer
- ▶ While still incorporating **other reasons** for co-occurrence

## Example with Eucalyptus data

---

```
TR <- read.csv("../data/eucalyptTR.csv");row.names(TR) <- TR$TAXON
```

## Example with Eucalyptus data

```
jsdm4 <- gllvm::gllvm(Y, X, TR = TR,
  formula = ~ Sandiness + cvTemp + (Sandiness + cvTemp):
(MedianSLA + MaxHeight.m. + MedianSeedMass.mg.),
  randomX = ~Sandiness + cvTemp,
  family = "binomial", method = "EVA", starting.val = "zero")
```

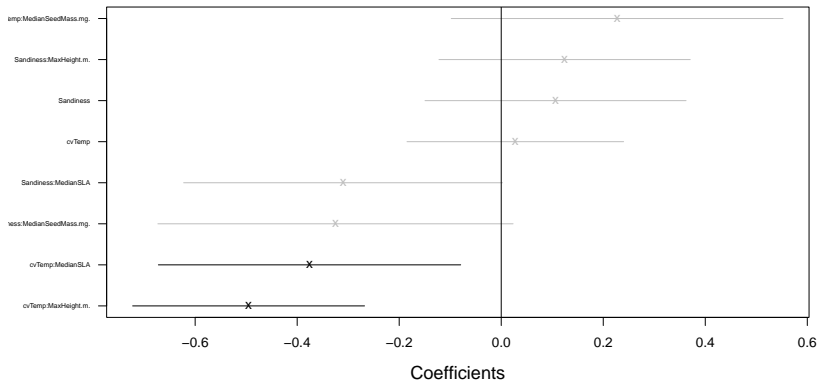
The formula formulation follows the equation: the first two terms are  $\beta_x$ , the next terms represent  $\mathbf{B}_{xtr}$ , and  $b_{kj}$  last





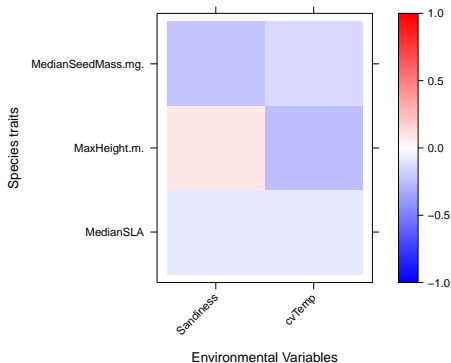
## Example with Eucalyptus data

```
gllvm::coefplot(jsdm4)
```



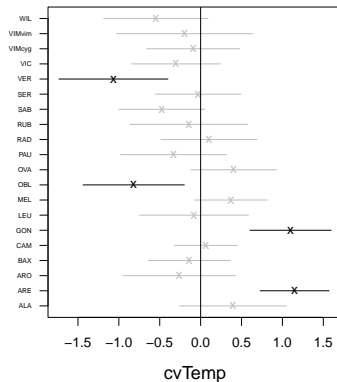
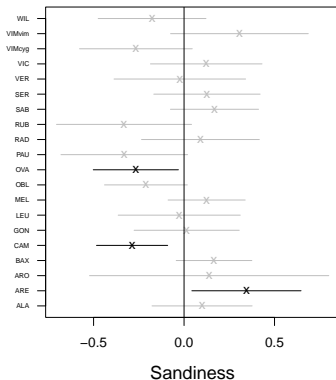
## Example with Eucalyptus data

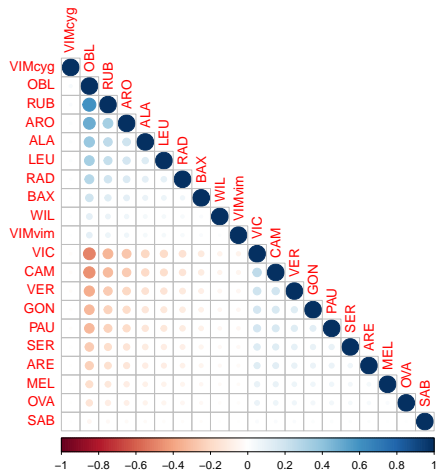
```
plot.4th <- lattice::levelplot(as.matrix(jsdm4$fourth.corner), xlab = "Environmental Variables",
                               ylab = "Species traits", col.regions = colort(100), cex.lab = 1.3,
                               at = seq(-1, 1, length = 100), scales = list(x = list(rot = 45)))
plot.4th
```



## Example with Eucalyptus data

```
gllvm::randomCoeffplot(jsdm4)
```





## Example with Eucalyptus data: hypothesis testing

```
anova(jsdm4, jsdm3)
```

```
## Model 1 : y ~ NULL
```

```
## Model 2 : ~ Sandiness + cvTemp + (Sandiness + cvTemp):(MedianSLA + Ma
```

##	Resid.Df	D	Df.diff	P.value
## 1	9096	0.00000	0	
## 2	9090	96.63163	6	0

We accept the alternative hypothesis: species responses to the environment are structured by traits

So far, we have been discussing **unstructured** species associations  
But what if we have information to provide?

(and we can predict for species without data)

$$\boldsymbol{\eta} = \mathbf{1}\beta_{0j}^\top + \mathbf{X}\mathbf{B} \quad (8)$$

- ▶  $\mathbf{B}$  are the random effects for covariates
- ▶ We assume  $\mathbf{B} \sim \mathcal{N}(\mathbf{0}, \Sigma_m \otimes \Sigma_r)$
- ▶  $\Sigma_r$  covariance matrix of random effects
- ▶  $\Sigma_m$  correlation matrix due to phylogeny

**We assume that all our random effects are structured by the Phylogeny**







## Phylogenetic signal

- ▶ 1: Fully phylogenetically structured responses
- ▶ 0: Normal (“iid”) random effects

When it is 0, it does not mean there is nothing going on.

Absence of phylogenetic signal:

- ▶ Scale mismatch
- ▶ Evolution moves very fast
- ▶ Too little information
- ▶ Traits are phylogenetically structured
- ▶ There are other (flexible) terms in the model
- ▶ Model misspecification

Presence of phylogenetic signal:

- ▶ Related species have similar “traits” (here: environmental

- ▶ unless species do not stably co-occur and/or evolution is still ongoing

# Example with fungi data (Abrego 2021)

Received: 1 November 2021 | Accepted: 20 December 2021

DOI: 10.1111/1365-2745.13839

## RESEARCH ARTICLE

Journal of Ecology



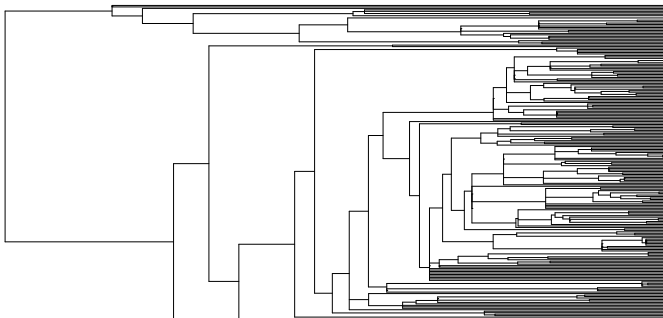
# Traits and phylogenies modulate the environmental responses of wood-inhabiting fungal communities across spatial scales

Nerea Abrego<sup>1,2</sup> | Claus Bässler<sup>3,4</sup> | Morten Christensen<sup>5</sup> | Jacob Heilmann-Clausen<sup>6</sup>



## Example with fungi data

---



## Example with fungi data

---

Phylogenetic models in `gllvm` use a **nearest neighbour approximation**

- ▶ We need to set the number of tips to consider on the tree
- ▶ The ordering of species matters!

```
covMat <- ape::vcv(tree)
e <- eigen(covMat)
distMat <- ape::cophenetic.phylo(tree)
ord <- gllvm::findOrder(covMat = covMat, distMat = distMat,
species <- colnames(covMat)[ord]
Y <- Y[, species]
covMat <- covMat[species, species]
distMat <- distMat[species, species]
```



## Ordering species

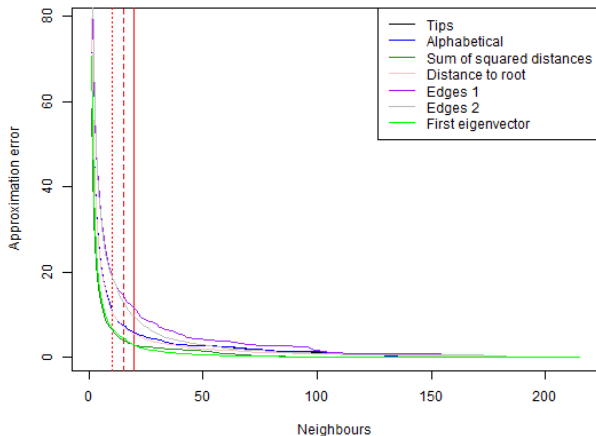


Figure 4: See vignette 7

## Example with fungi data

```
TMB::openmp(parallel::detectCores()-1, autpar = TRUE, DLL = "gllvm")
jsdm5 <- gllvm::gllvm(y = Y, X=X, family = "binomial", num.lv = 0, beta0com = TRUE,
  row.eff = ~(1 | REGION/RESERVE), studyDesign = X[,c("REGION", "RESERVE")],
  formula = ~(DBH.CM+AVERDP+I(AVERDP^2)+CONNECT10+TEMPR+PRECIP+log.AREA|1),
  colMat = list(covMat, dist = distMat), nn.colMat = 15, max.iter = 10e3, optim.method = "L-BFGS-B")
```

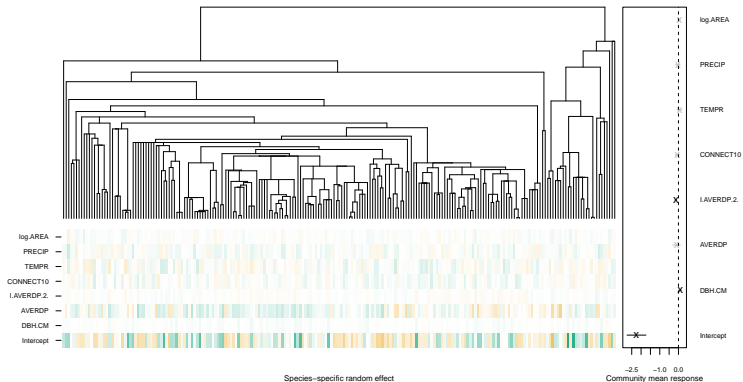
This takes a while to fit, but is really incredibly fast (it is just a complex model)

```
## Call:
## glvm::glvm(y = Y, X = X, formula = ~(DBH.CM + AVERDP + I(AVERDP^2) +
## CONNECT10 + TEMPR + PRECIP + log.AREA | 1), family = "binomial",
## num.lv = 0, studyDesign = X[, c("REGION", "RESERVE")], colMat = list(covMat,
## dist = distMat), row.eff = ~(1 | REGION/RESERVE), beta0com = TRUE,
## nn.colMat = 15, max.iter = 10000, optim.method = "L-BFGS-B")
##
## Family: binomial
##
## AIC: 103171.8 AICc: 103171.8 BIC: 103678.9 LL: -51539 df: 47
##
## Informed LVs: 0
## Constrained LVs: 0
## Unconstrained LVs: 0
##
## Formula: ~(DBH.CM + AVERDP + I(AVERDP^2) + CONNECT10 + TEMPR + PRECIP + log.AREA | 1)
## LV formula: ~ 0
## Row effect: ~(1 | REGION/RESERVE)
##
## Random effects:
## Name Signal Variance Std.Dev Corr
## Intercept 0.6037 1.0495 1.0244
## DBH.CM 0.6037 0.0051 0.0715 0.7642
## AVERDP 0.6037 0.1796 0.4238 0.0529 0.0458
## I.AVERDP.2 0.6037 0.0066 0.0815 -0.3550 -0.6454 -0.4397
## CONNECT10 0.6037 0.0401 0.2003 -0.3544 -0.0711 0.1790 -0.0091
## TEMPR 0.6037 0.0689 0.2625 0.0626 -0.2879 -0.0321 0.6438 0.3917
## PRECIP 0.6037 0.0440 0.2098 0.4461 0.5139 -0.4465 0.0148 -0.5701
## log.AREA 0.6037 0.0140 0.1184 -0.5196 -0.0889 0.0173 -0.3518 0.4538
##
##
##
##
```



## Example with fungi data

```
gllvm::phyloplot(jsdm5, tree)
```





## Summary

- ▶ JSDMs is a framework for analysing species co-occurrence data
- ▶ Focussed on prediction, but also suitable for inference
- ▶ We can also fit models with non-binary data (e.g., counts or biomass)
- ▶ The GLLVM framework is used here to implement JSDM efficiently
- ▶ We can incorporate random effects
- ▶ Phylogenetically structure species' effects
- ▶ Above all: we incorporate correlation of species