# The rise and rise of joint species distribution models (JSDMs) in ecology

Francis K.C. Hui Australian National University

Multivariate abundance data

Gen 1: MGLMM

Gen 2: Latent variables/factor analysis

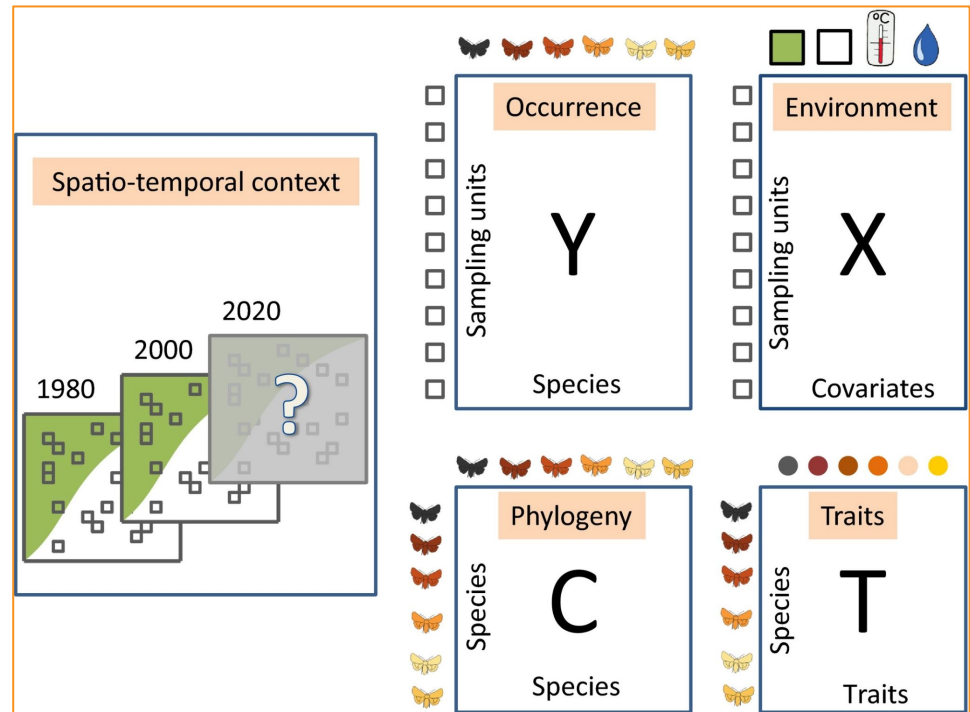Gen 2+: LVMs with all the extras

Closing thoughts

# Disclaimer

- This is an opinionated review/perspective talk, so you will see a decent chunk of my and my collaborators' works
  - Apologies for this!
  - Thank you to all who have/continue to inspire me to work on JSDMs
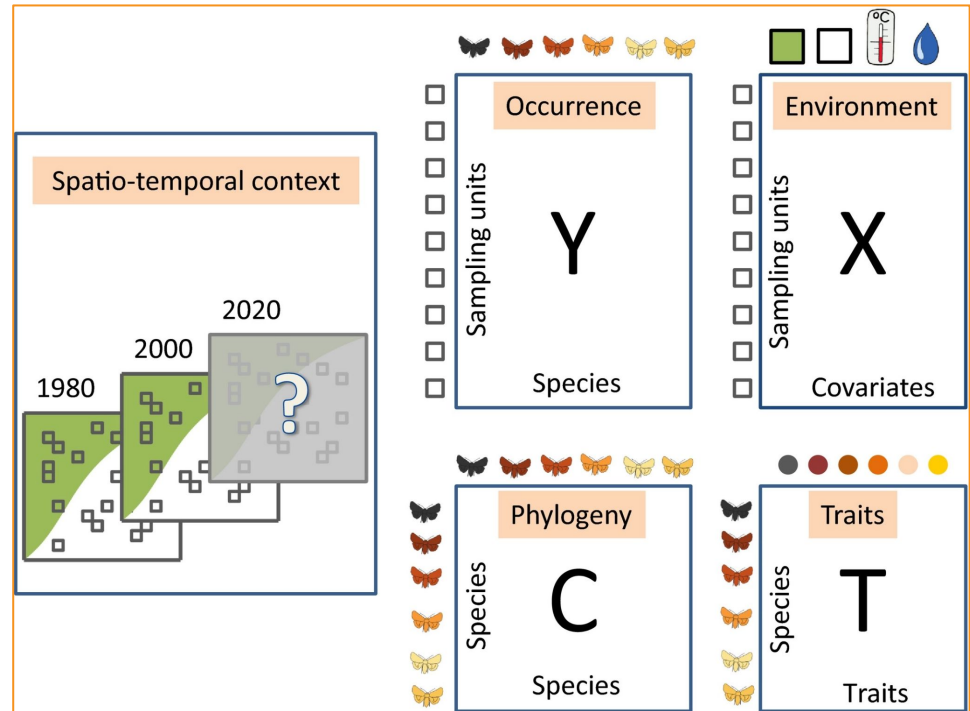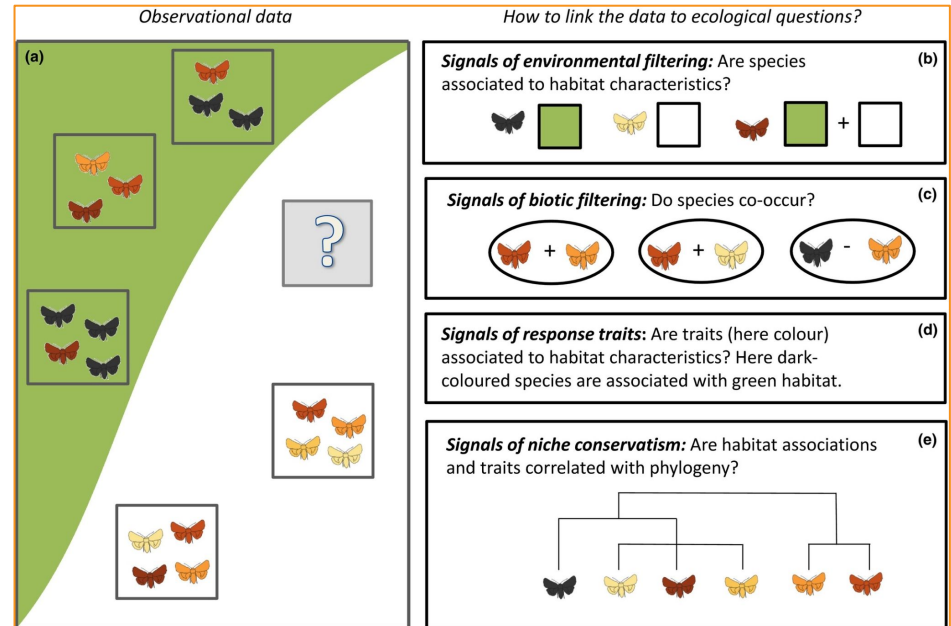
# Multivariate abundance data



3

# Multivariate abundance data

- Some common features:
  - Multiple correlated responses (high-dimensional)
  - Non-continuous responses with evident mean-variance relationship
  - Non-linear Y-X relationships



ECOLOGY LETTERS

Volume 20, Issue 5
May 2017
Pages 561-576

Idea and Perspective · Open Access

How to make more out of community data? A conceptual framework and its implementation as models and software

Otso Ovaskainen, Gleb Tikhonov, Anna Norberg, F. Guillaume Blanchet, Leo Duan, David Dunson, Tomas Roslin, Nerea Abrego

The effect of environmental and



4

# Question/s of interests

- Depends on the data you have:
  - (a) is a multivariate prediction problem
  - (b) -> how is Y and X related?
  - (c) -> how are the columns of Y related?
  - (d) + (e) -> how do T & C mediate/drive the Y–X relationship?



ECOLOGY LETTERS

Volume 20, Issue 5
May 2017
Pages 561-576

Idea and Perspective | 🔓 Open Access | ⓒ ⓘ

How to make more out of community data? A conceptual framework and its implementation as models and software

Otso Ovaskainen ✉ Gleb Tikhonov, Anna Norberg, F. Guillaume Blanchet, Leo Duan, David Dunson, Tomas Roslin, Nerea Abrego

Figures  References  Related  Information

The effect of environmental and

*Observational data*  |  *How to link the data to ecological questions?*

(a)

**Signals of environmental filtering:** Are species associated to habitat characteristics? (b)

**Signals of biotic filtering:** Do species co-occur? (c)

**Signals of response traits:** Are traits (here colour) associated to habitat characteristics? Here dark-coloured species are associated with green habitat. (d)

**Signals of niche conservatism:** Are habitat associations and traits correlated with phylogeny? (e)
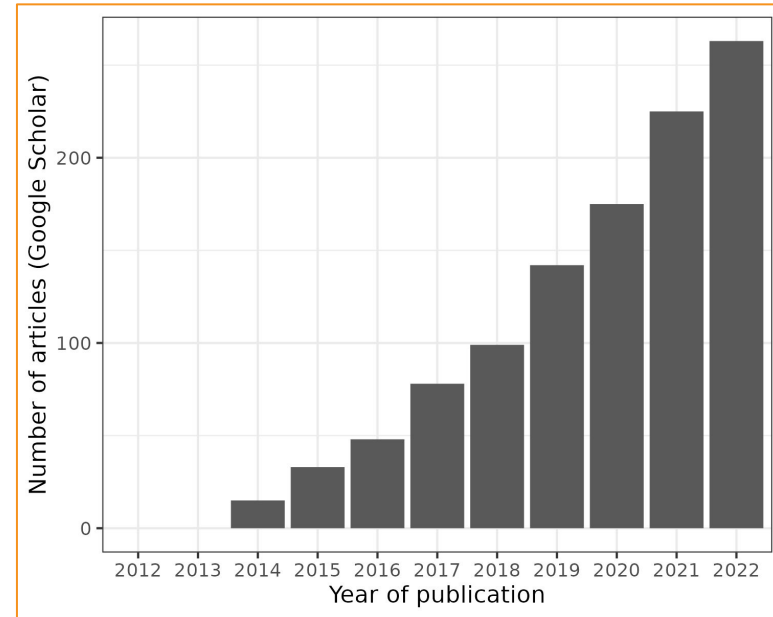
5

# Enter the joint species distribution model

- Loosely speaking, a joint species distribution model (JSDM) refers to a statistical method that <mark>simultaneously models all species</mark>
  - Accounts for the fact that species may be correlated with each other (after adjusting for measured predictor)
  - A single, potentially high-dimensional log-likelihood function
  - The sources of this (residual) correlation could be many…

# Enter the joint species distribution model

- Loosely speaking, a joint species distribution model (JSDM) refers to a statistical method that ==simultaneously models all species==
  - Accounts for the fact that species may be correlated with each other (after adjusting for measured predictor)
  - ==A single, potentially high-dimensional log-likelihood function==
  - The sources of this (residual) correlation could be many...

- JSDMs are basically a counterpart to stacked species distribution models (SSDMs), which model each species separately
  - Log-likelihood function comprises the sum of independent species contributions e.g., fit a GLM/GAM/GLMM/ML etc...to each species

# Enter the joint species distribution model

- A Google Scholar search of four key JSDM phrases (as of 23 November 2022)
  - `Joint species distribution models`
  - `Model-based ordination`
  - `Joint dynamic species distribution models`
  - `Hierarchical modeling of species communities`

- This is probably an underestimate of JSDM's rise...

# Enter the joint species distribution model

- A Google Scholar search of four key JSDM phrases (as of 8 October 2022)
  - Joint spec
  - Model-base
  - Joint dyna models
  - Hierarchic communitie

- This is probably JSDM's rise...

# Gen 1: MGLMMs

- Multivariate generalized linear mixed model (MGLMM)
  - Model residual between–species correlations using a multivariate random intercept

Figures References Related Information

Recommended

A comparison of joint species distribution models for presence–absence data

Figures References Related Information

Recommended

BOULDER COUNTY OPEN SPACE BUTTERFLY DIVERSITY AND ABUNDANCE

# Gen 1: MGLMMs

- Multivariate generalized linear mixed model (MGLMM)
  - Model residual between–species correlations using a multivariate random intercept
  - Exponential family is being used "loosely" here to cover many response distributions

Consider a set of species $j = 1, \ldots, m$ recorded at a set of observational units $i = 1, \ldots, N$, along with measured covariates $\boldsymbol{x}_i$. Then a vanilla JSDM is defined as

$$g(\mu_{ij}) = \eta_{ij} = \boldsymbol{x}_i^\top \boldsymbol{\beta}_j + e_{ij}$$
$$[\boldsymbol{e}_i] = \mathcal{N}(\boldsymbol{0}, \boldsymbol{\Sigma})$$
$$[y_{ij} | \boldsymbol{e}_i] = \mathsf{Exp\text{-}Fam}(\mu_{ij}, \boldsymbol{\phi}_j)$$
$$\ell(\boldsymbol{\Psi}) = \sum_{i=1}^N \log \left( \int \prod_{j=1}^m f(y_{ij} | \mu_{ij}, \boldsymbol{\phi}_j) f(\boldsymbol{e}_i) \, d\boldsymbol{e}_i \right)$$
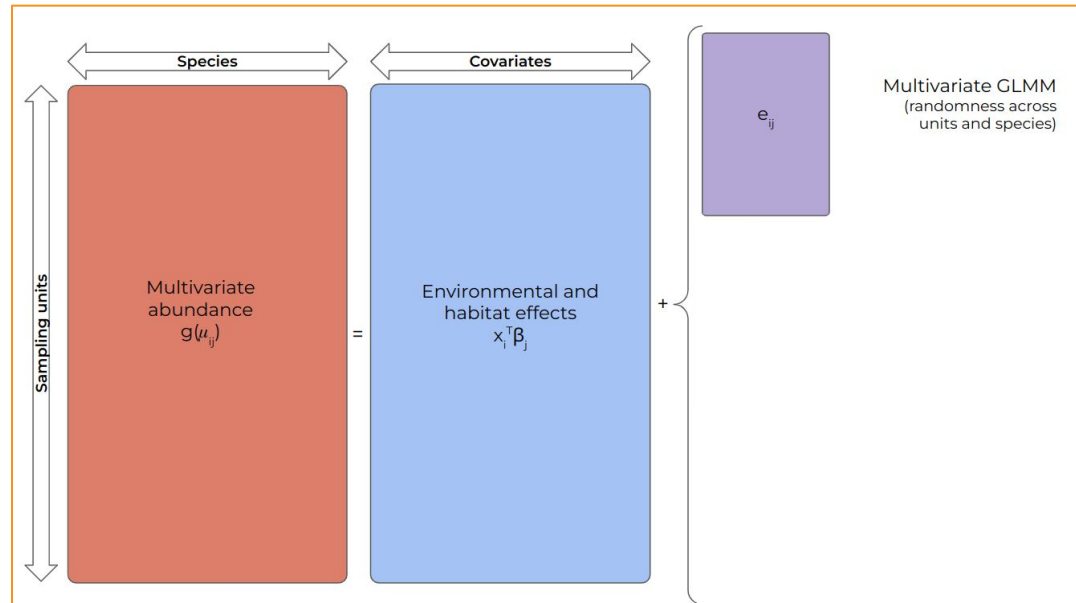
# Gen 1: MGLMMs

- Multivariate generalized linear mixed model (MGLMM)
  - Model residual between-species correlations using a multivariate random intercept
  - Exponential family is being used "loosely" here to cover many response distributions

# Gen 1: MGLMMs

- Multivariate generalized linear mixed model (MGLMM)
  - Model residual between-species correlations using a multivariate random intercept
  - Exponential family is being used "loosely" here to cover many response distributions
  - Pretty flexible (at least for correlations/symmetric associations)
  - Number of parameters scale as $m^2$, so great if m is not large (compared to N)
  - Lots of random effects, scaling as Nm

Consider a set of species $j = 1, \ldots, m$ recorded at a set of observational units $i = 1, \ldots, N$, along with measured covariates $\boldsymbol{x}_i$. Then a vanilla JSDM is defined as

$$g(\mu_{ij}) = \eta_{ij} = \boldsymbol{x}_i^\top \boldsymbol{\beta}_j + e_{ij}$$
$$[\boldsymbol{e}_i] = \mathcal{N}(\boldsymbol{0}, \boldsymbol{\Sigma})$$
$$[y_{ij}|\boldsymbol{e}_i] = \mathsf{Exp\text{-}Fam}(\mu_{ij}, \boldsymbol{\phi}_j)$$
$$\ell(\boldsymbol{\Psi}) = \sum_{i=1}^N \log \left( \int \prod_{j=1}^m f(y_{ij}|\mu_{ij}, \boldsymbol{\phi}_j) f(\boldsymbol{e}_i) \, d\boldsymbol{e}_i \right)$$

13

# Gen 1: MGLMMs

- Multivariate generalized linear mixed model (MGLMM)
  - Model residual between–species correlations using a multivariate random intercept
  - Exponential family is being used "loosely" here to cover many response distributions
  - Pretty flexible (at least for correlations/symmetric associations)
  - Number of parameters scale as $m^2$, so great if m is not large (compared to N)
  - Lots of random effects, scaling as Nm

- Largely overtaken by Gen 2 JSDMs, but advances continue to be made...
  - Translating ideas from sparse graphical model/network/ML literature

# Gen 2: Latent variable/factor analytic models

- Generalized linear latent variable models (LVMs)
  - Model residual between–species correlations using rank–reduction

## Trends in Ecology & Evolution

Volume 30, Issue 12, December 2015, Pages 766-779

Review

### So Many Variables: Joint Modeling in Community Ecology

David I. Warton [1], F. Guillaume Blanchet [2], Robert B. O'Hara [3], Otso Ovaskainen [4, 5], Sara Taskinen [6], Steven C. Walker [2], Francis K.C. Hui [7]

Show more ∨

+ Add to Mendeley    🔸 Cite

https://doi.org/10.1016/j.tree.2015.09.007          Get rights and content

---

Published: 24 August 2017

### Generalized Linear Latent Variable Models for Multivariate Count and Biomass Data in Ecology

Jenni Niku ✉, David I. Warton, Francis K. C. Hui & Sara Taskinen

*Journal of Agricultural, Biological and Environmental Statistics* **22**, 498–522 (2017) | Cite this article

**2453** Accesses | **28** Citations | **1** Altmetric | Metrics

---

## Methods in Ecology and Evolution

BRITISH ECOLOGICAL SOCIETY

Research Article | 🔓 Free Access

### Generating realistic assemblages with a joint species distribution model

David J. Harris ✉

First published: 05 January 2015 | https://doi.org/10.1111/2041-210X.12332 | Citations: 80

Volume 6, Issue 4
April 2015
Pages 465-473

Figures  References  Related  Information

Recommended

Defining and evaluating

15

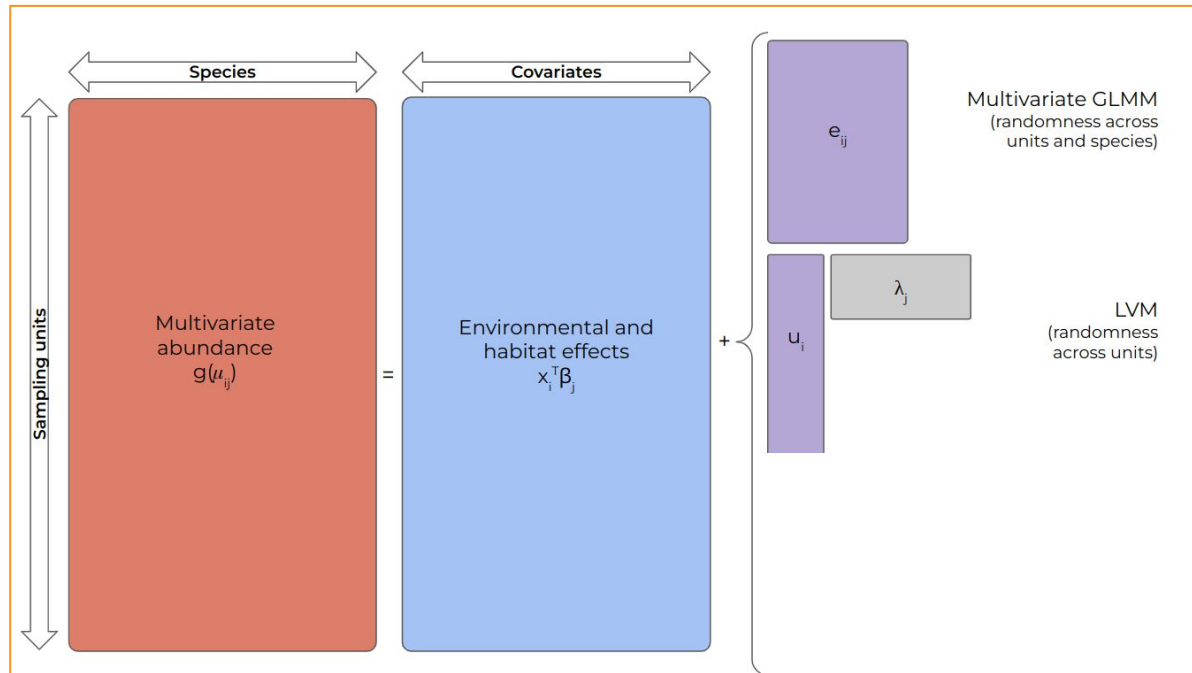# Gen 2: Latent variable/factor analytic models

- Generalized linear latent variable models (LVMs)
  - Model residual between–species correlations using rank–reduction

Consider a set of species $j = 1, \ldots, m$ recorded at a set of observational units $i = 1, \ldots, N$, along with covariates $\boldsymbol{x}_i$. Then a (basic) LVM is defined as

$$g(\mu_{ij}) = \eta_{ij} = \boldsymbol{x}_i^\top \boldsymbol{\beta}_j + \boldsymbol{u}_i^\top \boldsymbol{\lambda}_j$$

$$[\boldsymbol{u}_i] = \mathcal{N}(\boldsymbol{0}, \boldsymbol{I}_d); \ d \ll m$$

$$[y_{ij}|\boldsymbol{u}_i] = \text{Exp-Fam}(\mu_{ij}, \boldsymbol{\phi}_j); \ \text{Cov}(\eta_{ij}, \eta_{ij'}) = \boldsymbol{\lambda}_j^\top \boldsymbol{\lambda}_{j'}$$

$$\ell(\boldsymbol{\Psi}) = \sum_{i=1}^{N} \log \left( \int \prod_{j=1}^{m} f(y_{ij}|\mu_{ij}, \boldsymbol{\phi}_j) f(\boldsymbol{u}_i) \, d\boldsymbol{u}_i \right)$$

# Gen 2: Latent variable/factor analytic models

- Generalized linear latent variable models (LVMs)
  - Model residual between-species correlations using rank-reduction

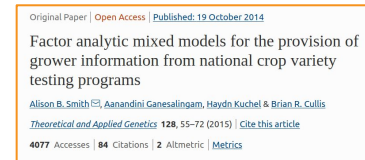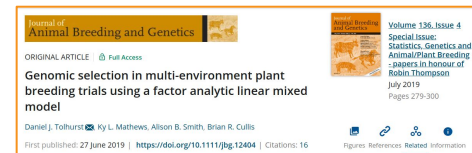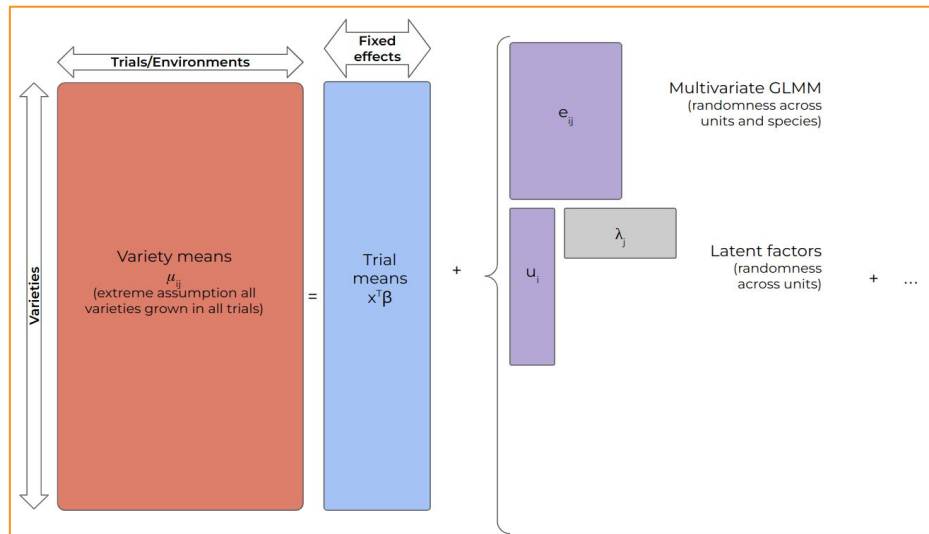# Gen 2: Latent variable/factor analytic models

- Generalized linear latent variable models (LVMs)
  - Model residual between–species correlations using rank–reduction
  - Less flexible than MGLMMs, but probably good enough in most scenarios?
  - Number of parameters scales as m, so can handle (a lot) more species
  - Less random effects than MGLMMs, scaling as Nd

Consider a set of species $j = 1, \ldots, m$ recorded at a set of observational units $i = 1, \ldots, N$, along with covariates $\boldsymbol{x}_i$. Then a (basic) LVM is defined as

$$g(\mu_{ij}) = \eta_{ij} = \boldsymbol{x}_i^\top \boldsymbol{\beta}_j + \boldsymbol{u}_i^\top \boldsymbol{\lambda}_j$$

$$[\boldsymbol{u}_i] = \mathcal{N}(\boldsymbol{0}, \boldsymbol{I}_d); \ d \ll m$$

$$[y_{ij}|\boldsymbol{u}_i] = \text{Exp-Fam}(\mu_{ij}, \boldsymbol{\phi}_j); \ \text{Cov}(\eta_{ij}, \eta_{ij'}) = \boldsymbol{\lambda}_j^\top \boldsymbol{\lambda}_{j'}$$

$$\ell(\boldsymbol{\Psi}) = \sum_{i=1}^{N} \log \left( \int \prod_{j=1}^{m} f(y_{ij}|\mu_{ij}, \boldsymbol{\phi}_j) f(\boldsymbol{u}_i) \, d\boldsymbol{u}_i \right)$$

# Gen 2: Latent variable/factor analytic models

- Generalized linear latent variable models (LVMs)
  - Model residual between-species correlations using rank-reduction

- LVMs are not new news! Examples include psychometrics, agriculture (MET)

# Gen 2: Latent variable/factor analytic models

- LVMs are not new news, but they took off in ecology!

# Gen 2: Latent variable/factor analytic models

- LVMs are not new news, but they took off in ecology!
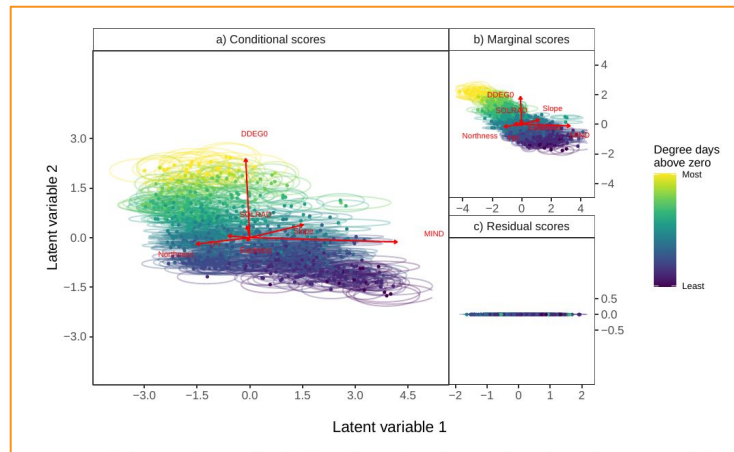  - Model–based unconstrained/partial/concurrent ordination, when d is small



Fig. 3. Model-based unconstrained ordination of the interannual mean catch of all fishing coves for all main taxa according to a functional taxonomy combining class and trophic group. We used a Tweedie distribution function with a log link function to model biomass responses and included random row effects.



$$g(\mu_{ij}) = \eta_{ij} = \boldsymbol{x}_i^\top \boldsymbol{\beta}_j + \boldsymbol{z}_i^\top \boldsymbol{\lambda}_j - \frac{1}{2}\boldsymbol{z}_j^\top \boldsymbol{D}_j \boldsymbol{z}_i$$

$$\boldsymbol{z}_i = \boldsymbol{C}^\top \boldsymbol{x}_{\mathsf{lv},i} + \boldsymbol{u}_i$$

$$[\boldsymbol{u}_i] = \mathcal{N}(\boldsymbol{0}, \boldsymbol{I}_d); \ d \ll m$$

$$[y_{ij}|\boldsymbol{u}_i] = \mathsf{Exp\text{-}Fam}(\mu_{ij}, \boldsymbol{\phi}_j)$$

# Gen 2: Latent variable/factor analytic models

- LVMs are not new news, but they took off in ecology!
  - Model–based unconstrained/partial/concurrent ordination, when d is small
  - Latent variables interpreted as unobserved environmental predictors
    - Neat interpretation but practically not very useful

# Gen 2: Latent variable/factor analytic models

- LVMs are not new news, but they took off in ecology!
  - Model–based unconstrained/partial/concurrent ordination, when d is small
  - Latent variables interpreted as unobserved environmental predictors
    - Neat interpretation but practically not very useful
  - Rank–reduction concept used in other community ecology contexts
    - Vector autoregressive models; community–level drivers/regulators

Species $j = 1, \ldots, m$ at time $t = 1, \ldots, T$

$$\log(\mu_{tj}) = \boldsymbol{x}_t^\top \boldsymbol{\beta}_j + (\boldsymbol{Q}\boldsymbol{c}_j + \boldsymbol{d}_j)^\top \log(\boldsymbol{\mu}_{t-1}) + \boldsymbol{u}_t^\top \boldsymbol{\lambda}_j + \delta_j$$

$$\dim(\boldsymbol{Q}) = m \times q; \; \dim(\boldsymbol{c}_j) = q \times 1; \; q \ll m$$

$$\boldsymbol{d}_j = (0, 0, \ldots, 0, d_j, 0, \ldots, 0)$$

$$[\boldsymbol{u}_t] = \mathcal{N}(\boldsymbol{0}, \boldsymbol{I}_d); \; d \ll m$$
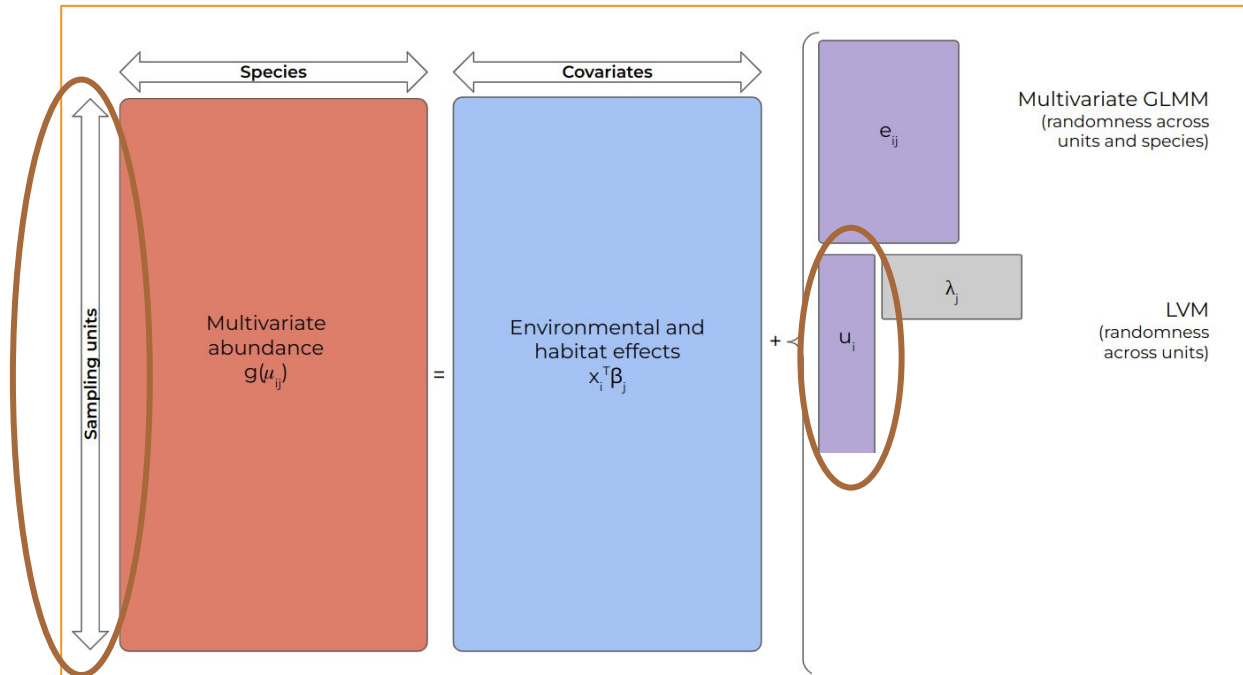
$$[\delta_j] = \mathcal{N}(0, \sigma^2)$$

# Gen 2+: LVMs with all the extras

- Current JSDM paradigm
  - Make LVMs more flexible and/or computationally more scalable

# Gen 2+: LVMs with all the extras

- <u>Example 1</u>: Spatio-temporal LVMs

# Gen 2+: LVMs with all the extras

- ● <u>Example 1</u>: Spatio–temporal LVMs
  - ○ Many flavours e.g., tensor–product or additive LVs, dynamic loadings
  - ○ Faster approximations/algorithms e.g., LVs + SPDE/NNGP/GPP

Consider a set of species $j = 1, \ldots, m$ recorded at a set of units $i = 1, \ldots, N$, each unit having a space–time coordinate $(\boldsymbol{s}_i, t_i)$. Then a (basic) spatio–temporal LVM is defined as

$$g\{\mu_j(\boldsymbol{s}_i, t_i)\} = \eta_j(\boldsymbol{s}_i, t_i) = \boldsymbol{x}(\boldsymbol{s}_i, t_i)^\top \boldsymbol{\beta}_j + \boldsymbol{u}(\boldsymbol{s}_i, t_i)^\top \boldsymbol{\lambda}_j$$

$$[\boldsymbol{u}_{.,k}] = [\{u_k(\boldsymbol{s}_1, t_1), \ldots, u_k(\boldsymbol{s}_N, t_N)\}] = \mathcal{N}(0, \boldsymbol{\Sigma}_k^{sp} \otimes \boldsymbol{\Sigma}_k^{time}); \ \boldsymbol{\Sigma}_k^{sp} \Rightarrow \mathsf{Matern}(\boldsymbol{\theta}_k^{sp}), \boldsymbol{\Sigma}_k^{time} \Rightarrow \mathsf{Matern}(\boldsymbol{\theta}_k^{time}); k = 1, \ldots, d$$

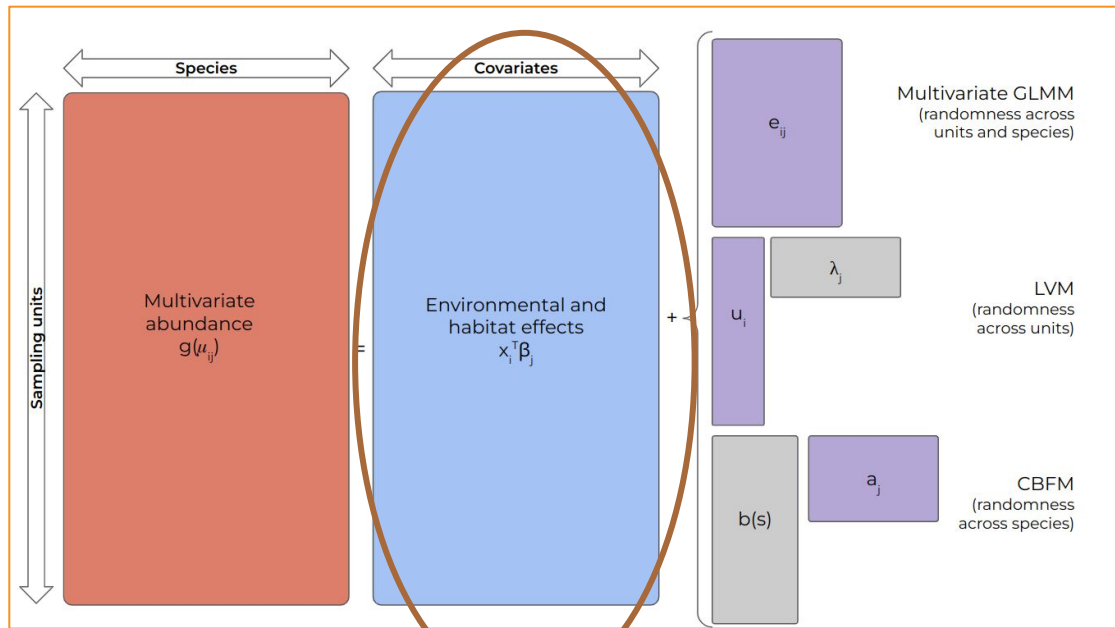$$[y_{ij} | \boldsymbol{u}_i] = \mathsf{Exp\text{-}Fam}(\mu_{ij}, \boldsymbol{\phi}_j)$$

$$\ell(\boldsymbol{\Psi}) = \log \left( \int \prod_{i=1}^N \prod_{j=1}^m f(y_j(\boldsymbol{s}_i, t_i) | \mu_j(\boldsymbol{s}_i, t_i), \boldsymbol{\phi}_j) \prod_{k=1}^d f(\boldsymbol{u}_{.,k}) \prod_{k=1}^d d\boldsymbol{u}_{.,k} \right)$$

Note $\mathrm{Cov}\{\eta_j(\boldsymbol{s}, t), \eta_{j'}(\boldsymbol{s}', t')\} = \sum_{k=1}^d \lambda_{jk} \Sigma_{k,ss'}^{sp} \Sigma_{k,tt'}^{time} \lambda_{j'k}$.

# Gen 2+: LVMs with all the extras

- <u>Example 2:</u> Borrow strength across species

# Gen 2+: LVMs with all the extras

- **Example 2:** Borrow strength across species
  - Traits mediate species mean responses to environment ("fourth–corner" models)
  - Phylogeny drives (dis)similarity in response to environment (phylogenetic LVMs)



Consider a set of species $j = 1, \ldots, m$ recorded at a set of observational units $i = 1, \ldots, N$, along with a set of $p$ covariates $\boldsymbol{x}_i$, an $m \times t$ trait matrix $\boldsymbol{T}$, and phylogenetic correlation matrix $\boldsymbol{C}$. Then a (basic) hierarchical LVM is defined as

$$g(\mu_{ij}) = \eta_{ij} = \boldsymbol{x}_i^\top \boldsymbol{\beta}_j + \boldsymbol{u}_i^\top \boldsymbol{\lambda}_j$$

$$[(\boldsymbol{\beta}_1, \ldots, \boldsymbol{\beta}_m)] = \mathcal{N}\{\text{vec}(\boldsymbol{K}\boldsymbol{T}^\top), \boldsymbol{V} \otimes (\rho\boldsymbol{C} + (1-\rho)\boldsymbol{I}_m)\}$$

$$\dim(\boldsymbol{K}) = p \times t; \ t < p; \ \dim(\boldsymbol{V}) = p \times p$$

$$[\boldsymbol{u}_i] = \mathcal{N}(\boldsymbol{0}, \boldsymbol{I}_d); \ d \ll m$$

$$[y_{ij}|\boldsymbol{u}_i] = \text{Exp-Fam}(\mu_{ij}, \boldsymbol{\phi}_j)$$

$$\ell(\boldsymbol{\Psi}) = \log\left( \int \prod_{i=1}^{N}\prod_{j=1}^{m} f(y_{ij}|\mu_{ij}, \boldsymbol{\phi}_j) \prod_{i=1}^{N} f(\boldsymbol{u}_i) \prod_{j=1}^{m} f(\boldsymbol{\beta}_j) \prod_{i=1}^{N} d\boldsymbol{u}_i \prod_{j=1}^{m} d\boldsymbol{\beta}_j \right)$$

Note $\text{Cov}(\eta_{ij}, \eta_{i'j'}) = \rho C_{j,j'} \boldsymbol{x}_i^\top \boldsymbol{V} \boldsymbol{x}_{i'} + \boldsymbol{\lambda}_j^\top \boldsymbol{\lambda}_{j'}$ for $j \neq j'$.
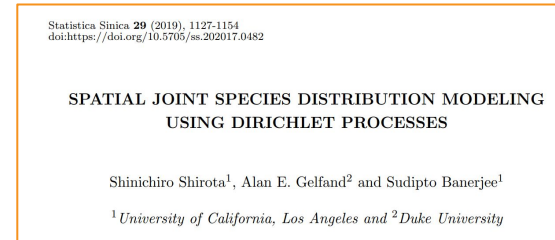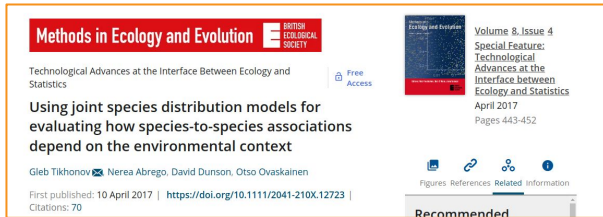
28

# Gen 2+: LVMs with all the extras

- <u>Example 3</u>: Borrow strength across species (in the loadings)

# Gen 2+: LVMs with all the extras

- **Example 3:** Borrow strength across species (in the loadings)
  - Clustering process on the loadings matrix (archetypal species associations)
  - Regress loadings against measured covariates (environment dependent associations)

Consider a set of species $j = 1, \ldots, m$ recorded at a set of observational units $i = 1, \ldots, N$, along with covariates $\boldsymbol{x}_i$. Then a (basic) loading-clustered LVM is defined as

$$g(\mu_{ij}) = \eta_{ij} = \boldsymbol{x}_i^\top \boldsymbol{\beta}_j + \boldsymbol{u}_i^\top \boldsymbol{Z}^\top \boldsymbol{q}(\boldsymbol{k}_j)$$
$$\dim(\boldsymbol{Z}) = r \times d; \; r \gg d; \; \dim\{\boldsymbol{q}(\boldsymbol{k}_j)\} = N \times 1$$
$$[\boldsymbol{k}_j] = \mathcal{DP}(\boldsymbol{\alpha}, \{1, 2, \ldots, r\});$$
$$[\boldsymbol{z}_{.l}] = \mathcal{N}(\boldsymbol{0}, \boldsymbol{W}); \; l = 1, \ldots, r$$
$$[\boldsymbol{u}_i] = \mathcal{N}(\boldsymbol{0}, \boldsymbol{I}_d); \; d \ll m$$
$$[y_{ij}|\boldsymbol{u}_i] = \text{Exp-Fam}(\mu_{ij}, \boldsymbol{\phi}_j)$$
$$\ell(\boldsymbol{\Psi}) = \text{I've never seen anyone try to estimate this using MLE!}$$

30

# Gen 2+: LVMs with all the extras

- There are many other extensions of LVMs, which I do not be cover/know about!

# Some closing remarks/thoughts

# Some closing remarks/thoughts

- JSDMs is a success story of how to translate and sell statistics…
  - Targeted software + relevant interpretations/answers + methods–vs–maths gap

# Some closing remarks/thoughts

- JSDMs is a success story of how to translate and sell statistics...
  - Targeted software + relevant interpretations/answers + methods–vs–maths gap

- Still many gaps in the JSDMs literature to close. Personal examples include:
  - Directional associations (structural equation modeling)?
  - Where do machine learning techniques come into this?
  - Data integration/fusion in JSDMs
  - Gen 3: Replacing latent variables with (spatio–temporal) basis functions
    - https://github.com/fhui28/CBFM

# Some closing remarks/thoughts

- JSDMs is a success story of how to translate and sell statistics…
  - Targeted software + relevant interpretations/answers + methods-vs-maths gap

- Still many gaps in the JSDMs literature to close. Personal examples include:
  - Directional associations (structural equation modeling)?
  - Where do machine learning techniques come into this?
  - Data integration/fusion in JSDMs
  - Gen 3: Replacing latent variables with (spatio-temporal) basis functions
    - https://github.com/fhui28/CBFM

- JSDMs is not the be-all and end-all
  - E.g., Stacked SDMs are still a powerful statistical approach
  - Do not throw the kitchen sink at something that does not need it

# **Thank you for listening!**

## *Any* **questions**?

- francis.hui@anu.edu.au
- https://francishui.netlify.app/

# Multivariate abundance data

- Some common features:
  - Multiple correlated responses (high-dimensional)
  - Non-continuous responses with evident mean-variance relationship
  - Non-linear Y-X relationships

- Other features:
  - Spatio-temporal (high-volume)
  - Multiple data sources
  - Background information

# Question/s of interests

- Depends on the data you have:
  - (a) is a multivariate prediction problem
  - (b) –> how is Y and X related?
  - (c) –> how are the columns of Y related?
  - (d) + (e) –> how do T & C mediate/drive the Y–X relationship?

- Some other applications:
  - Model–based ordination
  - Bioregionalization



JOURNAL ARTICLE   EDITOR'S CHOICE

**Bioregions in Marine Environments: Combining Biological and Environmental Data for Management and Scientific Understanding** [FREE]

Skipton N C Woolley ✉, Scott D Foster, Nicholas J Bax, Jock C Currie, Daniel C Dunn, Cecilie Hansen, Nicole Hill, Timothy D O'Hara, Otso Ovaskainen, Roger Sayre ... Show more

*BioScience*, Volume 70, Issue 1, January 2020, Pages 48–59, https://doi.org/10.1093/biosci/biz133
**Published:** 18 December 2019

**Trends in Ecology & Evolution**

Volume 30, Issue 12, December 2015, Pages 766-779

Review

So Many Variables: Joint Modeling in Community Ecology

David I. Warton [1] ✉, F. Guillaume Blanchet [2], Robert B. O'Hara [3], Otso Ovaskainen [4,5], Sara Taskinen [6], Steven C. Walker [2], Francis K.C. Hui [7]

**MOLECULAR ECOLOGY**

Volume 27, Issue 12
June 2018
Pages 2714-2724

ORIGINAL ARTICLE

Uncovering the drivers of host-associated microbiota with joint species distribution modelling

Johannes R. Björk ✉, Francis K. C. Hui, Robert B. O'Hara, Jose M. Montoya

First published: 14 May 2018 | https://doi.org/10.1111/mec.14718 | Citations: 25

Related     Information

Recommended

38

# Gen 2: Latent variable/factor analytic models

- Generalized linear latent variable models (LVMs)
  - Model residual between-response correlations using rank-reduction
  - Less flexible than MGLMMs, but probably good enough in most scenarios?*
  - Number of parameters scales as m, so can handle (a lot) more species*
  - Less random effects than MGLMMs, scaling as Nd; still quite challenging to fit**
    - *Choice of d remains a complicated and active topic
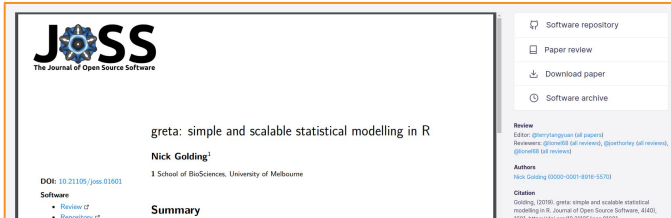    - **Lots of work has been done in this space

# Gen 2: Latent variable/factor analytic models

- Generalized linear latent variable models (LVMs)
  - Model residual between–species correlations using rank–reduction

- LVMs are not new news! Examples include psychometrics

# Gen 2: Latent variable/factor analytic models

- LVMs are not new news, but they took off in ecology!
  - Model–based unconstrained and partial ordination, when d is small

# Gen 3: ???

- Latent variables as an approach to JSDMs is awesome
  - But I think we are pushing the limits of their scalability/computability?

# Gen 3: ???

- Latent variables as an approach to JSDMs is awesome
  - But I think we are <mark>pushing the limits of their scalability/computability?</mark>

- Move the randomness from units to species -> <mark>basis functions</mark>

# Gen 3: CBFMs?

- Community-level basis function models (CBFMs) for spatio-temporal multivariate abundance data
  - Pre-defined spatio-temporal basis functions
  - https://github.com/fhui28/CBFM

Consider a set of species $j = 1, \ldots, m$ recorded at a set of units $i = 1, \ldots, N$, where each unit has a space-time coordinate $(\boldsymbol{s}_i, t_i)$. For a set of pre-defined spatio-temporal basis functions, $\boldsymbol{b}(\boldsymbol{s}, t)$, a (basic) CBFM can be defined as

$$g\{\mu_j(\boldsymbol{s}_i, t_i)\} = \eta_j(\boldsymbol{s}_i, t_i) = \boldsymbol{x}(\boldsymbol{s}_i, t_i)^\top \boldsymbol{\beta}_j + \boldsymbol{b}(\boldsymbol{s}_i, t_i)^\top \boldsymbol{a}_j$$

$$[\boldsymbol{a}] = [(\boldsymbol{a}_1, \ldots, \boldsymbol{a}_m)] = \mathcal{N}(\boldsymbol{0}, \boldsymbol{G} \otimes \boldsymbol{\Sigma})$$

$$\boldsymbol{G} = \boldsymbol{\Lambda}_G \boldsymbol{\Lambda}_G^\top + \kappa_G \boldsymbol{I}_m; \quad \dim(\boldsymbol{\Lambda}_G) = m \times d_m, d_m \ll m$$

$$\Rightarrow m \times m \text{ rank-reduced baseline between-species } \textit{correlation} \text{ matrix}$$

$$\boldsymbol{\Sigma} = \boldsymbol{\Lambda}_\Sigma \boldsymbol{\Lambda}_\Sigma^\top + \kappa_\Sigma \boldsymbol{I}_q; \quad \dim(\boldsymbol{\Lambda}_\Sigma) = q \times d_q, d_q \ll q$$

$$\Rightarrow q \times q \text{ rank-reduced community-level covariance matrix for basis functions.}$$

Note that $\text{Cov}\{\eta_j(\boldsymbol{s}, t), \eta_{j'}(\boldsymbol{s}', t')\} = G_{jj'} \boldsymbol{b}(\boldsymbol{s}, t)^\top (\boldsymbol{\Lambda}_\Sigma \boldsymbol{\Lambda}_\Sigma^\top + \kappa_\Sigma \boldsymbol{I}_q) \boldsymbol{b}(\boldsymbol{s}', t')$, where $G_{jj'} = 1$ if $j = j'$ and $\boldsymbol{\lambda}_{G,j}^\top \boldsymbol{\lambda}_{G,j'}$ otherwise

# Gen 3: CBFMs?

- But why would CBFMs be faster?

Consider a set of species $j = 1, \ldots, m$ recorded at a set of units $i = 1, \ldots, N$, where each unit has a space–time coordinate $(\boldsymbol{s}_i, t_i)$. For a set of pre-defined spatio-temporal basis functions, $\boldsymbol{b}(\boldsymbol{s}, t)$, a (basic) CBFM can be defined as

$$g\{\mu_j(\boldsymbol{s}_i, t_i)\} = \eta_j(\boldsymbol{s}_i, t_i) = \boldsymbol{x}(\boldsymbol{s}_i, t_i)^\top \boldsymbol{\beta}_j + \boldsymbol{b}(\boldsymbol{s}_i, t_i)^\top \boldsymbol{a}_j$$

$$[\boldsymbol{a}] = [(\boldsymbol{a}_1, \ldots, \boldsymbol{a}_m)] = \mathcal{N}(\boldsymbol{0}, \boldsymbol{G} \otimes \boldsymbol{\Sigma})$$

$$\boldsymbol{G} = \boldsymbol{\Lambda}_G \boldsymbol{\Lambda}_G^\top + \kappa_G \boldsymbol{I}_m; \quad \dim(\boldsymbol{\Lambda}_G) = m \times d_m, d_m \ll m$$

$$\Rightarrow m \times m \text{ rank–reduced baseline between-species } \textit{correlation} \text{ matrix}$$

$$\boldsymbol{\Sigma} = \boldsymbol{\Lambda}_\Sigma \boldsymbol{\Lambda}_\Sigma^\top + \kappa_\Sigma \boldsymbol{I}_q; \quad \dim(\boldsymbol{\Lambda}_\Sigma) = q \times d_q, d_q \ll q$$

$$\Rightarrow q \times q \text{ rank–reduced community-level covariance matrix for basis functions.}$$

Note that $\mathrm{Cov}\{\eta_j(\boldsymbol{s}, t), \eta_{j'}(\boldsymbol{s}', t')\} = G_{jj'} \boldsymbol{b}(\boldsymbol{s}, t)^\top (\boldsymbol{\Lambda}_\Sigma \boldsymbol{\Lambda}_\Sigma^\top + \kappa_\Sigma \boldsymbol{I}_q) \boldsymbol{b}(\boldsymbol{s}', t')$, where $G_{jj'} = 1$ if $j = j'$ and $\boldsymbol{\lambda}_{G,j}^\top \boldsymbol{\lambda}_{G,j'}$ other-wise

# Gen 3: CBFMs?

- But why would CBFMs be faster?
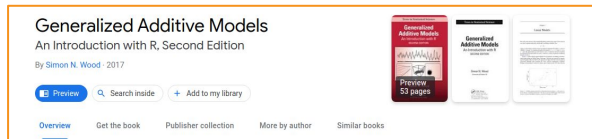  - Although m may not be small, N is still larger in most modern datasets

Consider a set of species $j = 1, \ldots, m$ recorded at a set of units $i = 1, \ldots, N$, where each unit has a space–time coordinate $(s_i, t_i)$. For a set of pre-defined spatio-temporal basis functions, $b(s, t)$, a (basic) CBFM can be defined as

$$g\{\mu_j(s_i, t_i)\} = \eta_j(s_i, t_i) = x(s_i, t_i)^\top \beta_j + b(s_i, t_i)^\top a_j$$

$$[a] = [(a_1, \ldots, a_m)] = \mathcal{N}(0, G \otimes \Sigma)$$

$$G = \Lambda_G \Lambda_G^\top + \kappa_G I_m; \quad \dim(\Lambda_G) = m \times d_m, d_m \ll m$$

$$\Rightarrow m \times m \text{ rank–reduced baseline between–species } \textit{correlation} \text{ matrix}$$

$$\Sigma = \Lambda_\Sigma \Lambda_\Sigma^\top + \kappa_\Sigma I_q; \quad \dim(\Lambda_\Sigma) = q \times d_q, d_q \ll q$$

$$\Rightarrow q \times q \text{ rank–reduced community–level covariance matrix for basis functions.}$$

Note that $\mathrm{Cov}\{\eta_j(s, t), \eta_{j'}(s', t')\} = G_{jj'} b(s, t)^\top (\Lambda_\Sigma \Lambda_\Sigma^\top + \kappa_\Sigma I_q) b(s', t')$, where $G_{jj'} = 1$ if $j = j'$ and $\lambda_{G,j}^\top \lambda_{G,j'}$ otherwise

46

# Gen 3: CBFMs?

- But why would CBFMs be faster?
  - Although m may not be small, N is still larger in most modern datasets
  - It is just a big generalized additive model (GAM)!

Consider a set of species $j = 1, \ldots, m$ recorded at a set of units $i = 1, \ldots, N$, where each unit has a space–time coordinate $(\boldsymbol{s}_i, t_i)$. For a set of pre-defined spatio-temporal basis functions, $\boldsymbol{b}(\boldsymbol{s}, t)$, a (basic) CBFM can be defined as

$$g\{\mu_j(\boldsymbol{s}_i, t_i)\} = \eta_j(\boldsymbol{s}_i, t_i) = \boldsymbol{x}(\boldsymbol{s}_i, t_i)^\top \boldsymbol{\beta}_j + \boldsymbol{b}(\boldsymbol{s}_i, t_i)^\top \boldsymbol{a}_j$$

$$[\boldsymbol{a}] = [(\boldsymbol{a}_1, \ldots, \boldsymbol{a}_m)] = \mathcal{N}(\boldsymbol{0}, \boldsymbol{G} \otimes \boldsymbol{\Sigma})$$

$$\boldsymbol{G} = \boldsymbol{\Lambda}_G \boldsymbol{\Lambda}_G^\top + \kappa_G \boldsymbol{I}_m; \quad \dim(\boldsymbol{\Lambda}_G) = m \times d_m, d_m \ll m$$

$$\Rightarrow m \times m \text{ rank-reduced baseline between-species } correlation \text{ matrix}$$

$$\boldsymbol{\Sigma} = \boldsymbol{\Lambda}_\Sigma \boldsymbol{\Lambda}_\Sigma^\top + \kappa_\Sigma \boldsymbol{I}_q; \quad \dim(\boldsymbol{\Lambda}_\Sigma) = q \times d_q, d_q \ll q$$

$$\Rightarrow q \times q \text{ rank-reduced community-level covariance matrix for basis functions.}$$

Note that $\mathrm{Cov}\{\eta_j(\boldsymbol{s}, t), \eta_{j'}(\boldsymbol{s}', t')\} = G_{jj'} \boldsymbol{b}(\boldsymbol{s}, t)^\top (\boldsymbol{\Lambda}_\Sigma \boldsymbol{\Lambda}_\Sigma^\top + \kappa_\Sigma \boldsymbol{I}_q) \boldsymbol{b}(\boldsymbol{s}', t')$, where $G_{jj'} = 1$ if $j = j'$ and $\boldsymbol{\lambda}_{G,j}^\top \boldsymbol{\lambda}_{G,j'}$ otherwise

# Gen 3: CBFMs?

- Basis functions are not new news
  - GAMs have been known in ecology for a long time. But not so much fixed rank kriging



Home / Annual Review of Statistics and Its Application / Volume 9, 2022 / Cressie

## Basis-Function Models in Spatial Statistics

**Annual Review of Statistics and Its Application**
Vol. 9:- (Volume publication date March 2022)
Review in Advance first posted online on November 18, 2021. (Changes may still occur before final publication.)
https://doi.org/10.1146/annurev-statistics-040120-020733

**Noel Cressie, Matthew Sainsbury-Dale, and Andrew Zammit-Mangion**
School of Mathematics and Applied Statistics, University of Wollongong, Wollongong, New South Wales 2522, Australia; email: ncressie@uow.edu.au

Download PDF | Article Metrics | Permissions | Reprints | Download Citation | Citation Alerts

**Abstract**



## ECOLOGY
ECOLOGICAL SOCIETY OF AMERICA

Volume 98, Issue 3
March 2017
Pages 632-646

Concepts & Synthesis | 🔓 Full Access

## The basis function approach for modeling autocorrelation in ecological data

Trevor J. Hefley ✉, Kristin M. Broms, Brian M. Brost, Frances E. Buderman, Shannon L. Kay, Henry R. Scharf, John R. Tipton, Perry J. Williams, Mevin B. Hooten

First published: 09 December 2016 | https://doi.org/10.1002/ecy.1674 |
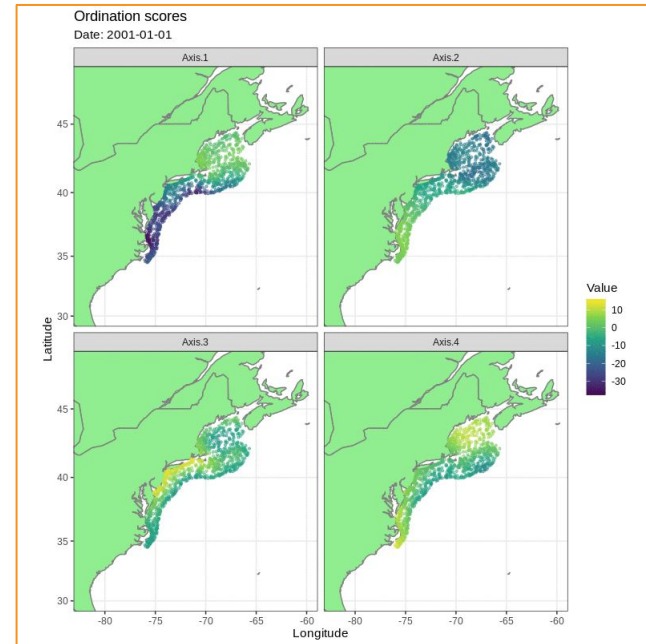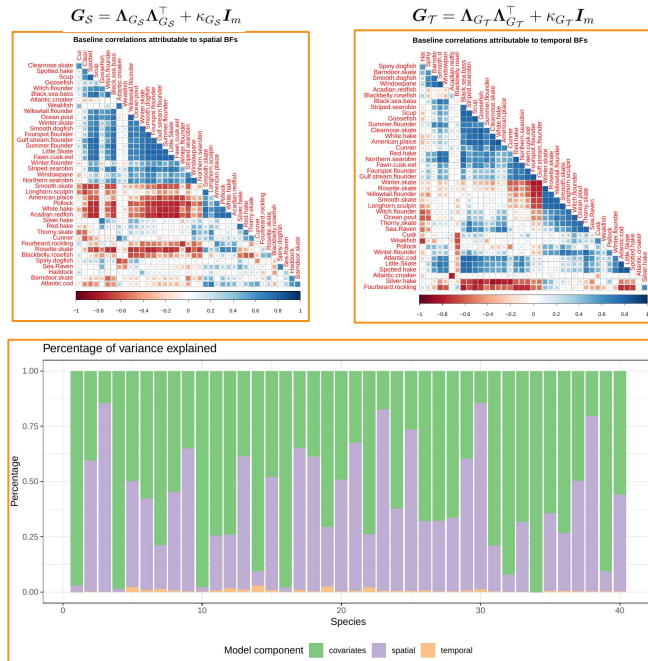Citations: 30

Figures References Related Information

Recommended

# Gen 3: CBFMs?

- Basis functions are not new news
  - GAMs have been known in ecology for a long time. But not so much fixed rank kriging
  - It takes a while to translate statistical methods to other disciplines (properly)…
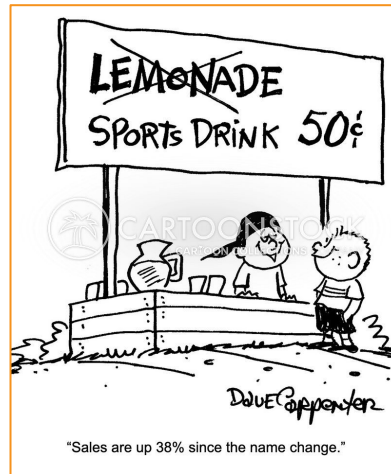
# Gen 3: CBFMs?

- Basis functions are not new news
    - GAMs have been known in ecology for a long time. But not so much fixed rank kriging
    - It takes a while to translate statistical methods to other disciplines (properly)...

- Is Gen 3 ⊂ Gen2+? Isn't basis functions just an approximation of spatio–temporal LVMs?
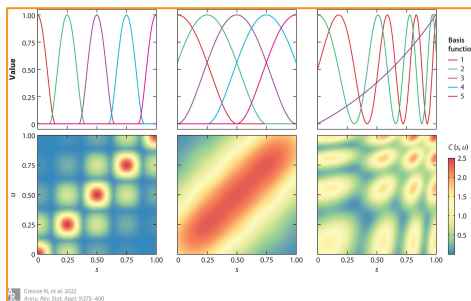


"Sales are up 38% since the name change."

# Gen 3: CBFMs?

- Basis functions are not new news
  - GAMs have been known in ecology for a long time. But not so much fixed rank kriging
  - It takes a while to translate statistical methods to other disciplines (properly)...

- Is Gen 3 ⊂ Gen2+? Isn't basis functions just an approximation of spatio–temporal LVMs?
  - Depends on how you want to approach basis functions: *"one person's mean is another person's covariance"* (Cressie, 1993)
  - A "basis function" mindset can opens up new opportunities

# Estimation, inference and all that jazz

- CBFM = Leveled up FRK = A very big generalized additive model (GAM)
  - Penalized quasi–likelihood (PQL) estimation for all coefficients, dispersion parameters; amenable to parallelization

Let $\boldsymbol{A}$ is the $m \times q$ matrix formed by stacking the $\boldsymbol{a}_j$'s as row vectors. Then given $\boldsymbol{G}$ and $\boldsymbol{\Sigma}$, update $(\boldsymbol{\beta}_j, \boldsymbol{a}_j)$'s, and $\phi_j$'s using

$$\ell_{\mathsf{PQL}} = \sum_{i=1}^{N} \sum_{j=1}^{m} \log\{ f(y_j(\boldsymbol{s}_i, t_i); \mu_j(\boldsymbol{s}_i, t_i), \boldsymbol{\phi}_j) \} - \frac{1}{2}\mathsf{tr}\left( \boldsymbol{G}^{-1} \boldsymbol{A} \boldsymbol{\Sigma}^{-1} \boldsymbol{A}^{\top} \right).$$

# Multivariate abundance data

- NorthEast Fisheries Science Center (NEFSC) fall bottom trawl survey
  - https://www.fisheries.noaa.gov/inport/item/22560
  - Subset of 2000–2019

# Multivariate abundance data

- Four example covariates in fall bottom trawl survey:
  - There are more covariates (between 20–30)...

# Multivariate abundance data

- Eight example demersal fish species in fall bottom trawl survey
  - Around 150ish taxa in total
  - High–dimensional, correlated responses

# Multivariate abundance data

- Eight example demersal fish species in fall bottom trawl survey
  - Around 150ish taxa in total
  - High-dimensional, correlated responses

- Some other noteworthy points:
  - You never visit the same location more than once
  - About 6,000 space-time locations visited between 2000-2019

# Multivariate abundance data

- Responses are:
  - Sparse, non-continuous
  - Strong mean-variance relationship (various reasons behind this)

# Estimation, inference and all that jazz

- CBFM = Leveled up FRK = <mark>A very big generalized additive model (GAM)</mark>
  - Penalized quasi–likelihood (PQL) estimation for all coefficients, dispersion parameters; amenable to parallelization
  - Maximum restricted Laplace–approximated likelihood estimation for the loadings and nugget effects

Let $\mathcal{X}$ and $\mathcal{B}$ be appropriately defined model matrices based on the $\boldsymbol{x}(\boldsymbol{s}_i, t_i)$ and $\boldsymbol{b}(\boldsymbol{s}_i, t_i)$'s respectively, and $\hat{\boldsymbol{W}}$ by a diagonal matrix of weights. Then given $(\boldsymbol{\beta}_j, \boldsymbol{a}_j)$'s, and $\phi_j$'s, update the loadings and nugget effects characterizing $\boldsymbol{G}$ and $\boldsymbol{\Sigma}$ using

$$
\ell_{\text{REML}} = \frac{q}{2} \log \det(\boldsymbol{G}^{-1}) + \frac{m}{2} \log \det(\boldsymbol{\Sigma}^{-1}) - \frac{1}{2} \text{tr} \left( \boldsymbol{G}^{-1} \hat{\boldsymbol{A}} \boldsymbol{\Sigma}^{-1} \hat{\boldsymbol{A}}^\top \right)
$$
$$
- \frac{1}{2} \log \det \left( \mathcal{B}^\top \left( \hat{\boldsymbol{W}} - \hat{\boldsymbol{W}} \mathcal{X} \left( \mathcal{X}^\top \hat{\boldsymbol{W}} \mathcal{X} \right)^{-1} \mathcal{X}^\top \hat{\boldsymbol{W}} \right) \mathcal{B} + \boldsymbol{G}^{-1} \otimes \boldsymbol{\Sigma}^{-1} \right).
$$

58

# Estimation, inference and all that jazz

- CBFM = Leveled up FRK = A very big generalized additive model (GAM)
  - Penalized quasi–likelihood (PQL) estimation for all coefficients, dispersion parameters; amenable to parallelization
  - Maximum restricted Laplace–approximated likelihood estimation for the loadings and nugget effects
  - Approximate large sample distributions for coefficients/linear predictors etc...

$$\begin{pmatrix} \hat{\boldsymbol{\beta}} \\ \hat{\boldsymbol{a}} \end{pmatrix} \approx \mathcal{N} \left\{ \begin{pmatrix} \boldsymbol{\beta}_0 \\ \boldsymbol{a}_0 \end{pmatrix}, \begin{pmatrix} \mathcal{X}^\top \hat{\boldsymbol{W}} \mathcal{X} & \mathcal{X}^\top \hat{\boldsymbol{W}} \mathcal{B} \\ \mathcal{B}^\top \hat{\boldsymbol{W}} \mathcal{X} & \mathcal{B}^\top \hat{\boldsymbol{W}} \mathcal{B} + \hat{\boldsymbol{G}}^{-1} \otimes \hat{\boldsymbol{\Sigma}}^{-1} \end{pmatrix}^{-1} \right\},$$

where $\boldsymbol{\beta}_0$ and $\boldsymbol{a}$ denote the true parameter values of the regression coefficients.

# Estimation, inference and all that jazz

- CBFM = Leveled up FRK = A very big generalized additive model (GAM)
  - Penalized quasi-likelihood (PQL) estimation for all coefficients, dispersion parameters; amenable to parallelization
  - Maximum restricted Laplace-approximated likelihood estimation for the loadings and nugget effects
  - Approximate large sample distributions for coefficients/linear predictors etc...
  - Adapt GAM tools for residual analysis, model selection, prediction etc...; variance-partitioning; space-time ordination using SVD-type ideas, and so on

$$\begin{pmatrix} \hat{\boldsymbol{\beta}} \\ \hat{\boldsymbol{a}} \end{pmatrix} \approx \mathcal{N} \left\{ \begin{pmatrix} \boldsymbol{\beta}_0 \\ \boldsymbol{a}_0 \end{pmatrix}, \begin{pmatrix} \mathcal{X}^\top \hat{\boldsymbol{W}} \mathcal{X} & \mathcal{X}^\top \hat{\boldsymbol{W}} \mathcal{B} \\ \mathcal{B}^\top \hat{\boldsymbol{W}} \mathcal{X} & \mathcal{B}^\top \hat{\boldsymbol{W}} \mathcal{B} + \hat{\boldsymbol{G}}^{-1} \otimes \hat{\boldsymbol{\Sigma}}^{-1} \end{pmatrix}^{-1} \right\},$$

where $\boldsymbol{\beta}_0$ and $\boldsymbol{a}$ denote the true parameter values of the regression coefficients.