# L-shaped data, GLM(M) and double constrained correspondence analysis:

#### from fourth-corner correlation to dc-CA

Cajo J.F. ter Braak, Biometris, WUR

with Petr Šmilauer (Ceske Budejovice), Stéphane Dray (Lyon) and Pedro Peres-Neto (Montréal)

Workshop, 12 January 2018, Erasmus University, Rotterdam









### Abstract

L-shaped data consists of a non-negative central matrix with associated matrices with predictors for rows and columns. Formally, it is (weighted) bigraph with node predictors. Examples are preference data of consumers for products with features of both consumers and products as predictors, supervisory boards of firms with features of supervisors and firms as predictors for the membership, and, in ecology, abundance data of species and environmental variables with traits and environmental variables as predictors. We will discuss the statistical issues of analysing such data and why double constrained correspondence analysis and GLM(M) methods may give very similar results in terms of selecting important features.

An alternative title is:

From the fourth-corner correlation to dc-CA.



## L-shaped data ( $\Gamma$ –shaped data)

■ Central matrix (Y≥0, ) with associated descriptors for rows and columns (E and T)



## **Examples of central table Y**

Data on:

#### Preference of consumers for products

- Which consumer characteristics and product features can predict the preference
- consumer segments, niche markets, niche products

#### Supervisory board memberships of firms

• Which person characteristics determine which type of firm they supervise?

#### Abundance of species in sites

- Which traits (T) of species determine in which type of environments (E, sites by variables) they prosper
- Trait-based ecology, trait-environment relationships



## **Γ**-shaped data

■ Central matrix (Y≥0) with associated descriptors for rows and columns (E and T)





**Γ**-shaped data in ecology: the fourth corner problem

■ Central matrix (Y ≥0) with associated descriptors for rows and columns (E and T)

	Species	Environment
Missing cell or matrix: 'the fourth corner' e.g. 'correlation' between E and T	Y= abundances	E= pH, temp, elevation
	<b>T<sup>t</sup>=</b> bodymass, SLA	'the fourth corner'
WAGENINGEN Legendre et al, Ecology 1997, Dray & L 2008 UNIVERSITY & RESEARCH Brown et al MEE 2014		

### **Issues with Γ –shaped data**

- How to define and test correlations between T and E as there is no common unit of observation?
  - the trait is observed on species,
  - the environment on sites and
  - the mediating abundance on species-site combinations.
- And... observational data only, neither environment nor traits can be randomized as in a designed experiment.
- In ecology, a number of methods such as RLQ (1996) and the fourth-corner correlation (1997) have been proposed to estimate such trait-environment associations.
- What about GLM(M) models for such data?



## An illustrative example

#### **Dutch Dune Meadow data set**

- Abundance (0-9) of 28 plant species in 20 dune meadows
- **Relation/interaction between**
- Trait: SLA (specific leaf area) of species and
- Environmental variable: moisture in the meadow???
- -GLM test on interaction (site bootstrap) \* :  $p \approx 0.03$
- -4<sup>th</sup> corner correlation with default resampling<sup>\*\*</sup>:  $p\approx 0.28$
- Which one cannot be trusted and why???



## Simplest GLM model: log-linear model

Abundance is a count  $y_{ij}$ , assumed to follow a Poisson or neg. bin. distribution with mean specified by

$$\log(\mu_{ij}) = r_i + c_j + \beta_{te} t_j e_i$$
(1)

- r<sub>i</sub> and c<sub>j</sub> row (site) and column (species) main effects (saturated main effects; e⊆ {r<sub>i</sub>}; t ⊆ {c<sub>j</sub>};
- $\beta_{te}$  the coefficient measuring the direction and strength of the t-e interaction
- $H_0:\beta_{te} = 0$  and  $H_1:$  with $\beta_{te} \neq 0$ .



See Gabriel 1998 Generalized bilinear regression Fit of model via GLM

Works via vectorization of Y

gives a single data frame/data set with

- *n×m* rows
- variables: e.g.

species, site, abundance, trait1, trait2, env1 env2, env3

Model:

yield ~ species + site + trait1:env1 + trait2:env1 + ....

Note: allows trait or environmental variable to vary within a species or site



## Selection and testing of traits and environmental variables via GLM with resampling

## Warton et al. use GLM with negative binomial error and adjust for model misspecification via resampling

- 1. Lasso model selection (Brown et al, 2014)
- 2. Statistical testing (Warton, Shipley and Hastie, 2015)

## Advocate: "design-based" i.e. site-based resampling ignoring any randomness of the other entity: species



Methods in Ecology and Evolution 2014, 5, 344-352

environment



doi: 10.1111/2041-210X.12163

The fourth-corner solution – using predictive models to understand how species traits Methods in Ecol

Alexandra M. Brown<sup>1</sup>\*, David I. Warton<sup>1</sup>, Nigel R. Andre and Heloise Gibb<sup>4</sup>



.. ..

Methods in Ecology and Evolution



Methods in Ecology and Evolution 2015, 6, 389-398

doi: 10.1111/2041-210X.12280

SPECIAL FEATURE NEW OPPORTUNITIES AT THE INTERFACE BETWEEN ECOLOGY AND STATISTICS CATS regression – a model-based approach to studying trait-based community assembly

David I. Warton1\*, Bill Shipley2 and Trevor Hastie3

## Allow for variability among species: **GLMM** model

•  $\log(\mu_{ij}) = r_i + c_j + b_j e_i$  (a model without traits...)

- r<sub>i</sub> and c<sub>j</sub> row (site) and column (species) main effects (saturated main effects; e⊆ {r<sub>i</sub>}; t ⊆ {c<sub>j</sub>};
- *b<sub>j</sub>* a species-specific slope with respect to e

Insert the trait information:  $b_j \sim N(\beta_0 + \beta_{te} t_j, \sigma_b^2)^*$  then

$$\log(\mu_{ij}) = r_i + c_j + \beta_{te}t_je_i + \beta_{ze}z_je_i, \qquad (2)$$
  
with  $\beta_{ze} = \sigma_b$  and  $z_j \sim N(0,1) \qquad z_j = latent trait$ 

#### Again:

 $H_0:\beta_{te} = 0$  and  $H_1:$  with $\beta_{te} \neq 0$ .

\* Could be a multivariate normal with covariance matrix depending on the phylogenetic relationships matrix



Pollock et al 2012 replaced  $r_i$  by  $\beta_e e_i$ Jamil et al 2013 assumed  $r_i$  random

## Simulation study (single trait t, single env e)

- The world more likely looks like the GLMM model (2) with negative binomial response, *i.e.* 
  - there are species-specific slopes wrt to e
- But we analyse using the simple Poisson model (1)
- Test null hypothesis of no trait-environment relation
- $H_0:\beta_{te} = 0$  and  $H_1:$  with $\beta_{te} \neq 0$
- with the Poisson or neg. bin. LR/deviance difference as the test statistic in the resampling test of significance



#### Type I error rate in 1000 simulations ter Braak et al 2017, PeerJ

- traitglm (Warton/Hastie) site-based bootstrap (R package mvabund), negative binomial deviance
- sites: site-based permutation of counts, Poisson deviance
- Species: species-based permutation of counts, Poisson deviance
- max r/c: Maximum of the site and species resampling p-values

VAGENINGEN



## **Revisiting the illustrative example:**

In the Dutch Dune Meadow data is: SLA  $\leftrightarrow$ moisture??? -GLM test on interaction (site bootstrap) \* : p  $\approx$  0.03 -4<sup>th</sup> corner correlation with default resampling<sup>\*\*</sup>: p  $\approx$  0.28 Which one cannot be trusted and why???





#### Revisiting the illustrative example: ter Braak 2017, EESt, p.231

- In the Dutch Dune Meadow data is:
- SLA ↔ moisture???
- -GLM test on interaction (site bootstrap) \* :  $p \approx 0.03$ -4<sup>th</sup> corner correlation with default resampling<sup>\*\*</sup>:  $p \approx 0.28$ Which one cannot be trusted and why???
- The slopes wrt moisture are species-specific (GLMM model)
- There is a second ('latent') trait ( z = Seedmass) that has
  - about zero correlation (-0.047) with SLA and
  - interacts with moisture (p<sub>row</sub> <0.001,  $p_{col} \approx 0.01$ )

There is thus no real evidence for SLA ↔ moisture.



#### Failure of site-based only tests ter Braak et al 2017, PeerJ, p.13

The issue is not that of *confounding* or *omitted variable* 

*confounding* is due to an *omitted variable* that is highly correlation with variable of interest and the predictor

In trait-environment problems, the problem :

occurs also if there is an omitted variable that has **zero correlation** with the predictor, and

is due to ignoring

species as a random factor, so as to account for

species-specific response to the environment

(an important random effect)

Conclusion: perform a species-based test wageningen too and take max *p*-values  $\rightarrow$  max r/c test

Alternative for GLM: the fourth-corner correlation?

In such simulations, I also investigated a simpler test statistic than deviance:

the squared fourth-corner correlation

Surprise, surprise.....

fourth-corner correlation gave similar type I error and power as the GLM deviance!!

How does this come about? So, what is this fourthcorner correlation



Fourth corner correlation *f* 

Legendre et al 1997

 $f = cor_{Y}(e, t)] = e^{t} Y t$ if e and t are normalized, See slide 23 for explicit formula i.e. have weighted mean and sd: 0 and 1 using as weights the row- and column-totals of Y

For count Y data:

f = correlation between **e** and **t** in inflated data in which each individual is a row with

values for **e** and **t** of the individual (**e** from its site, **t** from its species)

e.g. a count of 5 in Y gives 5 identical rows in the inflated data WAGENINGEN WAGENINGEN UNIVERSITY & RESEARCH GLM and fourth corner correlation r<sub>4</sub> ter Braak EEST 2017

GLM model: count  $y_{ij}$  follows a Poisson distribution with mean specified by

 $\log(\mu_{ij}) = r_i + c_j + \beta_{te} \, \mathbf{t}_j \mathbf{e}_i \tag{1}$ 

 $f^2 y_{++}$  = squared fourth corner correlation ×  $y_{++}$ 

= Rao score test statistic

for testing the linear-by-linear interaction  $H_0$ :  $\beta_{te} = 0$ 

Asymp. equivalent with LR, much quicker to compute!

**Extension to multiple traits and environmental variables:** 

Score test statistic =  $y_{++}$  × inertia of dc-CA



## Corollary

- T = I<sub>m</sub> gives single constrained correspondence analysis which is canonical correspondence analysis (CCA, ter Braak 1986)
  - Total inertia of CCA ×  $y_{++}$  = Rao's score test statistic
- Used as test statistic in permutation testing since 1990 in Canoco and later in R::vegan

So, we discovered a new property of a much used method!

The result gives a reason for renewed interest in dc-CA



## And is this all a surprise? Hmm...

- T = I<sub>m</sub>, E = I<sub>n</sub> gives (unconstrained) correspondence analysis (CA)
  - Total inertia of CA ×  $y_{++} = y_{++} \sum_a \lambda_a = \chi^2$
  - which is a Rao score test statistic on row-column independence
- T = t, E = e gives the simplest case of dc-CA with  $\lambda_1 = [cor_Y(\mathbf{e}, \mathbf{t})]^2 = f^2$

Recall an original definition of CA (Hirshfield 1935, Fisher 1940)

• CA finds a latent e<sup>\*</sup> and latent t<sup>\*</sup> such that  $\lambda_1 = max_{\{x,u\}}[cor_Y(x,u)]^2 = [cor_Y(e^*,t^*)]^2 = \max f^2$ with e<sup>\*</sup>, t<sup>\*</sup>row- and column scores of CA

→the maximum attainable squared fourth-corner correlation is thus the first CA-eigenvalue!



## History of correspondence analysis (CA)

- CA: Hirschfield 1935, Fisher 1940, Guttman 1941, Benzecri 1969, Hill 1974, Gifi 1990 and many others..
- Single constrained CA (CCA): ter Braak 1986/7, Chessel, Lebreton et al 1987/8, with a precursor: Green 1971!
- Double constrained CA: Bacou & Sabatier 1989, Lavorel & Lebreton 1998/9, Böckenholt & Böckenholt 1990, Takane 2013
- Many different rationales! Relations to PCA, contingency tables, analysis of variance, log-linear models, unfolding, gradient analysis, Gaussian response models,...
- All are special cases of canonical correlation analysis (or, except dc-CA, of discriminant analysis)
- But... it is nontrivial to do the computing via a program for canonical correlation analysis ...so Algorithms for...



#### From fourth corner correlation to dc-CA ter Braak et al EEST 2018

#### fourth-corner correlation f between trait t and environmental variable e

$$f = cor_{\mathbf{Y}}^2(\mathbf{t}, \mathbf{e}) = \frac{\sum_{i,j} y_{ij} \tilde{t}_j \tilde{e}_i}{\{\sum_j y_{+j} \tilde{t}_j^2 \sum_i y_{i+} \tilde{e}_i^2\}^{1/2}}$$
(1)

with

$$\tilde{t}_{j} = t_{j} - \sum_{j} y_{+j} t_{j} / y_{++}$$
 and  $\tilde{e}_{i} = e_{i} - \sum_{i} y_{i+} e_{i} / y_{++}$  (2)  
**Definition**:

dc-CA is a method that finds linear combinations of traits and of environmental variables that maximize their fourth corner correlation



## **Derivation of dc-CA**

Assume traits and environmental variables are centered

 $\mathbf{1}_n^T \mathbf{R} \mathbf{E} = \mathbf{0}_p$  and  $\mathbf{1}_m^T \mathbf{K} \mathbf{T} = \mathbf{0}_q$ 

with R = diag( $\{y_{i+}\}$ ) and K= diag( $\{y_{+j}\}$ ).

The definition of dc-CA leads to the following maximization problem:  $max_{b,c} \mathbf{x}^T \mathbf{Y} \mathbf{u}$  with  $\mathbf{x} = \mathbf{Eb}$ ,  $\mathbf{u} = \mathbf{Tc}$ ,  $\mathbf{x}^T \mathbf{Rx} = 1$  and  $\mathbf{u}^T \mathbf{Ku} = 1$  (3) or

 $max_{b,c} b^T E^T YTc$  subject to  $b^T E^T REb = 1$  and  $c^T T^T KTc = 1$ . (4) Lagrange multiplier method leads to

$$\lambda_{b} \mathbf{b} = (\mathbf{E}^{T} \mathbf{R} \mathbf{E})^{-1} \mathbf{E}^{T} \mathbf{Y} \mathbf{T} \mathbf{c} \qquad (6)$$

$$\lambda_{c} \mathbf{c} = (\mathbf{T}^{T} \mathbf{K} \mathbf{T})^{-1} \mathbf{T}^{T} \mathbf{Y}^{T} \mathbf{E} \mathbf{b} \qquad (7)$$

$$\rightarrow \lambda(\mathbf{E}^{T} \mathbf{R} \mathbf{E}) \mathbf{b} = \mathbf{E}^{T} \mathbf{Y} \mathbf{T} (\mathbf{T}^{T} \mathbf{K} \mathbf{T})^{-1} \mathbf{T}^{T} \mathbf{Y}^{T} \mathbf{E} \mathbf{b} \qquad (8)$$

$$\mathbf{G} \mathbf{E} \mathbf{N} \mathbf{G} \mathbf{E} \mathbf{N}$$

## Transition formulae of dc-CA

- **1.**  $\lambda^{\alpha} u_k^* = \sum_i y_{ik} x_i / y_{+k}$  or in matrix notation,  $\lambda^{\alpha} \mathbf{u}^* = \mathbf{K}^{-1} \mathbf{Y}^T \mathbf{x}$
- $2. \mathbf{c} = (\mathbf{T}^T \mathbf{K} \mathbf{T})^{-1} \mathbf{T}^T \mathbf{K} \mathbf{u}^*$
- $\mathbf{\mathcal{J}}_{\mathbf{u}} = \mathbf{T}\mathbf{c}$
- **4.**  $\lambda^{1-\alpha} x_i^* = \sum_k y_{ik} u_k / y_{i+}$  or in matrix notation,  $\lambda^{1-\alpha} \mathbf{x}^* = \mathbf{R}^{-1} \mathbf{Y} \mathbf{u}$
- $5. \mathbf{b} = (\mathbf{E}^T \mathbf{R} \mathbf{E})^{-1} \mathbf{E}^T \mathbf{R} \mathbf{x}^*$
- $\boldsymbol{\textit{6.}} \quad \mathbf{x} = \mathbf{E}\mathbf{b}$
- $\lambda$  = eigenvalue, c and b are canonical weights,  $\alpha \in [0,1]$  user-defined.

Two sets of row scores  $\{x_i\}$  and  $\{x_i^*\}$  & columns scores,  $\{u_k\}$  and  $\{u_k^*\}$ 

**1&4** 
$$\to$$
 **CA with**  $\{u_k^* = u_k\}$  and  $\{x_i^* = x_i\}$  or  $\{\mathbf{E} = \mathbf{I}_n, \mathbf{T} = \mathbf{I}_m\}$ 

**1,4,5&6**  $\rightarrow$  **CCA** with { $u_k^* = u_k$ } or **T** = **I**<sub>m</sub>

iterative algorithm based on this: power algorithm, slow but can be accelerated

## Algorithm based on a SVD

#### Similar to canonical correlation. Define

 $\mathbf{D} = (\mathbf{E}^T \mathbf{R} \mathbf{E})^{-1/2} \mathbf{E}^T \mathbf{Y} \mathbf{T} (\mathbf{T}^T \mathbf{K} \mathbf{T})^{-1/2}$ 

SVD of D:

 $\mathbf{D} = \mathbf{P} \mathbf{\Delta} \mathbf{Q}^{\mathrm{T}}$ 

with P and Q orthonormal matrices and  $\Delta$  a diagonal matrix with singular values in decreasing order.

Then the singular values are the maximized fourth corner correlations of the dc-CA axes and the columns of

$$\mathbf{B} = (\mathbf{E}^T \mathbf{R} \mathbf{E})^{-1/2} \mathbf{P} \Delta^{\alpha} \text{ and } \mathbf{C} = (\mathbf{T}^T \mathbf{K} \mathbf{T})^{-1/2} \mathbf{Q} \Delta^{\alpha - 1}$$

satisfy the transition formulae.

**X** = **EB** and **U** = **TC**, are **R**- and **K**-orthogonal.

The scaling factor  $\Delta^{\alpha}$  ensures that  $X^T R X = \Lambda^{\alpha}$  and  $U^T K U = \Lambda^{1-\alpha}$ , where  $\Lambda = \Delta^2$ 

 $tr(\mathbf{D}^T\mathbf{D}) = \sum_a \lambda_a$  is the Rao score test statistic/ $y_{++}$ 



#### Comparison with dc-PCA Douglas Carrol et al 1980, two-way CANDELINC

A weighted dc-PCA can be obtained from an SVD of

 $\mathbf{D}_{dc-pca} = (\mathbf{E}^T \mathbf{R} \mathbf{E})^{-1/2} \mathbf{E}^T \mathbf{R} \mathbf{Y} \mathbf{K} \mathbf{T} (\mathbf{T}^T \mathbf{K} \mathbf{T})^{-1/2}$ 

**Compare:** 

$$\mathbf{D}_{\mathbf{dc-ca}} = (\mathbf{E}^T \mathbf{R} \mathbf{E})^{-1/2} \mathbf{E}^T \mathbf{Y} \mathbf{T} (\mathbf{T}^T \mathbf{K} \mathbf{T})^{-1/2}$$

 $\rightarrow$  dc-CA is a weighted dc-PCA of the contingency ratios  $y_{++} \mathbf{R}^{-1} \mathbf{Y} \mathbf{K}^{-1}$ 

with weight matrices with  $\mathbf{R} = \text{diag}(\{y_{i+}\})$  and  $\mathbf{K} = \text{diag}(\{y_{+j}\})$ .

All very similar... dc-CA is a natural method for count-like data



Comparison with RLQ (1) (the standard in ecology) Dolédec et al EEST 1996

#### An RLQ can be obtained from an SVD of

 $D_{rlq} = E^T YT$  with E and T R- and K-standardized Compare:

 $\mathbf{D}_{\mathbf{dc-ca}} = (\mathbf{E}^T \mathbf{R} \mathbf{E})^{-1/2} \mathbf{E}^T \mathbf{Y} \mathbf{T} (\mathbf{T}^T \mathbf{K} \mathbf{T})^{-1/2}$ 

 $\rightarrow$  dc-CA uses the correlations among traits & among environmental variables, whereas RLQ does not

 $\rightarrow$  RLQ is more robust to near-collinearity than dc-CA,

dc-CA needs regularization or variable selection to counter this

Another way of saying similar things:

 $\rightarrow$  dc-CA is based on correlation (based on regression)

→ RLQ is based on covariance (based on coinertia analysis, a tiny wageningen bit like PLS)

## **Comparison with RLQ (2)**

**Because its regression base:** 

- dc-CA can reveal trait and environment dimensions that remain hidden in RLQ
  - if trait and/or env. vars. are moderately correlated
- A simulation study, 10,000 simulated data sets with:
  - *n=m* = 100
  - 6 traits, 9 environmental variables ~ AR<sub>1</sub>(0.7)
  - One latent dimension defined by a contrast of the first two traits and the first two environmental variables; a second dimension unrelated to E,T.
  - So: 4 of the traits and 7 of the env. vars are noise



## dc-CA reveals the contrast, RLQ does not



## Algorithm based on combining CCA and RDA

... gives insight in relations with another existing method, called CWM-RDA (combine two tables (Y & T), then use a two-table method):

- **1.** Combine Y with T in a single table of trait means per site  $M = R^{-1}YT$
- **2.** Analyze M ~ E by redundancy analysis (RDA)
- This is essentially an SVD of

 $\mathbf{D}_{\text{cwm-rda}} = (\mathbf{E}^T \mathbf{E})^{-1/2} \mathbf{E}^T \mathbf{M} = (\mathbf{E}^T \mathbf{E})^{-1/2} \mathbf{E}^T (\mathbf{R}^{-1} \mathbf{Y} \mathbf{T})$ 

Lacks R-weighing and trait covariances

Obtain dc-CA by adding R&K-weighing and a prior orthonormalization of T

Can be done by first performing a CCA and then a weighted RDA on its scores ...
Useful in Canoco as it has



Useful in Canoco as it has testing and selection of variables for (weighted) RDA

## Quadriplot of dc-CA: example

#### 5 out of 6 pairs are weighted least-squares biplots of:

- **1.** Fourth-corner correlations:  $E^T YT$
- **2.** E means per species (SNCs)
- **3.** T means per site (CWMs)
- **4.** Contingency ratios
- **5.** Trait data<sup>\*</sup> T

Dune meadow data: n= 20, m = 28 two traits two environmental variables

\* In column-metric preserving scaling and with fixed species points





**Concluding remarks on L-shaped data** 

#### Statistical issues

- Rows, columns and values are random
- Needs GLMMs or
- Simpler models (GLM, fourth-corner, dc-CA) with
  - Combination of row and column resampling as "the noise in the rows is likely different from that in the columns"
- Fourth-corner and dc-CA
  - provide Rao score test statistics of GLM models that are useful in resampling

dc-CA allows easy testing and variable selection scheme

• Combining row and columns analyses (Canoco 5.10)



## Some references

Jamil, T., W. A. Ozinga, M. Kleyer, and C. J. F. ter Braak. 2013. Selecting traits that explain species–environment relationships: a generalized linear mixed model approach. Journal of Vegetation Science 24:988-1000. <a href="http://dx.doi.org/10.1111/j.1654-1103.2012.12036.x">http://dx.doi.org/10.1111/j.1654-1103.2012.12036.x</a>

Peres-Neto, P. R., S. Dray, and C. J. F. ter Braak. 2017. Linking trait variation to the environment: critical issues with community-weighted mean correlation resolved by the fourth-corner approach. Ecography 40:806-816. http://dx.doi.org/10.1111/ecog.02302

ter Braak, C. J. F. 2017. Fourth-corner correlation is a score test statistic in a log-linear trait–environment model that is useful in permutation testing. Environmental and Ecological Statistics 24:219-242. http://dx.doi.org/10.1007/s10651-017-0368-0

ter Braak, C. J. F., P. Peres-Neto, and S. Dray. 2017. A critical issue in model-based inference for studying traitbased community assembly and a solution. PeerJ 5:e2885. <u>https://doi.org/10.7717/peerj.2885</u>

ter Braak, C. J. F., P. Šmilauer, and S. Dray. 2018. Algorithms and biplots for double constrained correspondence analysis. Environmental and Ecological Statistics. <u>https://doi.org/10.1007/s10651-017-0395-x</u>.or.<u>http://rdcu.be/ETPh</u>







