

Multivariate extended Poisson-Tweedie regression model

Prof. Wagner Hugo Bonat

Laboratory of Statistics and Geionformation - LEG

Paraná Federal University - UFPR

Conference on Multivariate Count Analysis
Besançon/France

July/2018

Motivation

- ▶ Count data are frequent in many research areas:
 1. Medical research, biology.
 2. Environmental and crop sciences.
 3. Economical, social and political sciences.
 4. Computer science, electronic engineering and etc.
- ▶ Statistical challenges:
 1. Overdispersion (mean $>$ variance).
 2. Underdispersion (mean $<$ variance).
 3. Zero-inflation.
 4. Heavy tail.
- ▶ Multiple responses each one with its own set of challenges!
- ▶ Practical interest:
 1. Estimation of the covariance structure.
 2. Multivariate hypothesis tests (MANOVA-type test).
 3. and much more!!

Data set 1: Australian health survey

- ▶ Australian health survey for 1987–1988 (Deb and Trivedi, 1997).
- ▶ Five count outcomes (Ndoc, Nndoc, Nmed, Nhosp, Nadm).
- ▶ Nine covariates concerning social conditions.
- ▶ 5190 respondents and no missing data.
- ▶ Goals: assess covariate effects (specifically and overall) and correlation between outcomes.

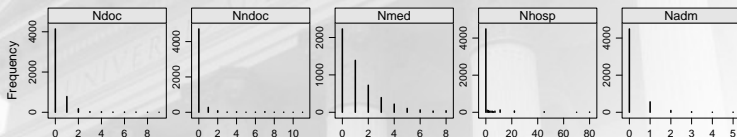


Figure 1. Histograms for each outcome in the Australian health survey data.

Data set 1: Australian health survey

- ▶ Generalized dispersion index GDI (Kokonendji and Puig, 2018)

$$\hat{GDI} = \frac{\sqrt{\bar{\mathbf{y}}^\top \hat{C} \hat{O}V(\mathbf{y}) \sqrt{\bar{\mathbf{y}}}}}{\bar{\mathbf{y}}^\top \bar{\mathbf{y}}}.$$

- ▶ $\hat{GDI} = 17.94$ ($sd = 1.84$).
- ▶ GDI for pairs.

##	Ndoc	Nndoc	Nadm	Nhosp	Nmed
## Ndoc	4.222377	3.283812	2.335228	27.28871	2.257114
## Nndoc	3.283812	8.681901	3.588466	28.07777	2.212992
## Nadm	2.335228	3.588466	2.967461	28.45916	2.075137
## Nhosp	27.288710	28.077771	28.459161	56.16685	17.391456
## Nmed	2.257114	2.212992	2.075137	17.39146	3.977885

- ▶ Overdispersion.

Data set 2: Ant data

- ▶ Abundances of 20 epigeaic ant species across 30 sites in south-eastern Australia.
- ▶ Covariates: Percent cover of bare ground and shrub.
- ▶ Goal: assess the covariates overall effects and describe the relation between species abundances.

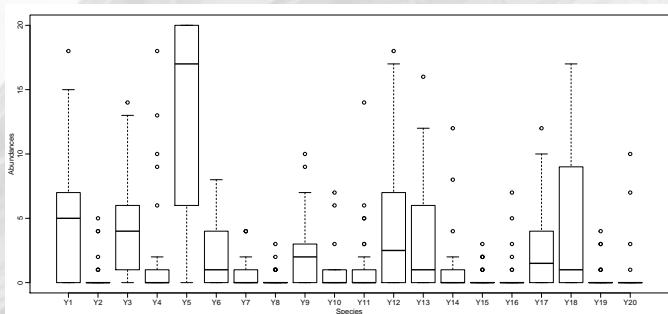
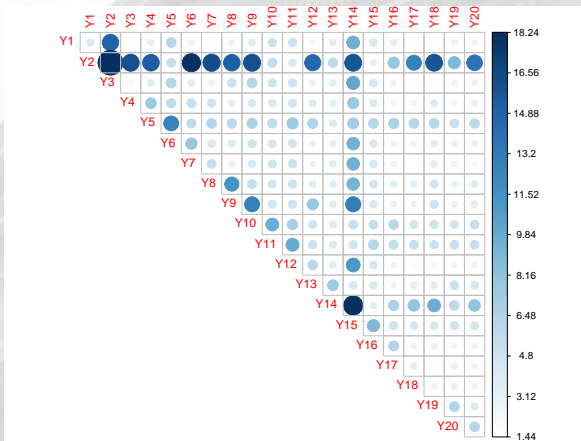


Figure 2. Boxplots of abundances by species.

Data set 2: Ant data

- ▶ $GDI = 6.48(sd = 1.19)$.
- ▶ GDI for every pair.



Review: Univariate approaches

► Models for overdispersion:

1. Poisson-Tweedie and Hinde-Demétrio (Kokonendji, 2004).
2. Negative binomial, Poisson-inverse Gaussian, Sichel, Poisson-SGIG, Delaporte, Poisson-Tweedie (Rigby, et.al., 2008).
3. Normalized tempered stable distribution (Kolossiatis, 2011).
4. Discrete Linnik distribution (Barabesi, 2016).

► Models for under and overdispersion:

1. Conway-Maxwell-Poisson (Sellers, 2010).
2. Gamma-Count (Zeviani, 2014).
3. Discrete Weibull (Kalktawi, 2015).

► Models for heavy tail or zero-inflation:

1. Generalised Poisson-inverse Gaussian family (Zhu, 2009).
2. Hurdle models (Zeileis, 2008).
3. Zero-inflated Poisson and negative binomial (Loeys, 2012).
4. Zero-inflated COM-Poisson (Sellers, 2016).

Review: Multivariate approaches

- ▶ Extensions of univariate distributions (David, 2017):
 1. Marginals are Poisson or negative binomial, etc.
 2. Mixture of Poissons or negative binomials, etc.
 3. Conditional Poissons or negative binomials, etc.
- ▶ Constructing multivariate distributions:
 1. Multivariate dispersion models (Jørgensen, 2000).
 2. Multivariate exponential dispersion models (Jørgensen, 2012).
 3. Convolution and extended convolution method (Jørgensen, 2013).
- ▶ General statistical modelling frameworks:
 1. Multivariate generalized linear mixed models.
 2. Copula based-models.
 3. Multivariate hierarchical generalized linear models.
 4. and ...

Introduction

- ▶ Plethora of distributions/approaches to deal with count data.
- ▶ Multivariate probability distribution is not available in closed-form.
- ▶ Difficult to fit (problems due to badly behaved likelihood function).
- ▶ Difficult to interpret model parameters.
- ▶ Difficult to point, with conviction, the best practical choice.
- ▶ Demand for a unified model that can automatically adapt to the underlying multivariate count distribution.
- ▶ Easy implementation in practice.
- ▶ SOLUTION: Multivariate extended Poisson-Tweedie model!

Multivariate extended Poisson-Tweedie model

- ▶ Extend the Poisson-Tweedie regression model to deal with multiple response variables.
- ▶ Propose an estimating function approach for parameter estimation.
- ▶ Propose an extension of the orthodox MANOVA for dealing with multivariate count data.
- ▶ Provide R code for fitting the models.
- ▶ Illustrative examples to show the flexibility of the multivariate extended Poisson-Tweedie regression model.

Tweedie distribution

- ▶ Tweedie distributions (Jørgensen, 1997)

$$f(z; \mu, \phi, p) = a(z, \phi, p) \exp\{(z\psi - k(\psi))/\phi\},$$

where $\mu = E(Z) = k'(\psi)$ is the mean.

- ▶ $\phi > 0$ and ψ are the dispersion and canonical parameters. $k(\psi)$ is the cumulant function and $a(z, \phi, p)$ is a normalizing constant.
- ▶ $\text{Var}(Z) = \phi\mu^p$ where $p \in (-\infty, 0] \cup [1, \infty)$ is the index determining the distribution.
- ▶ Special cases: Gaussian ($p = 0$), Poisson ($p = 1$), non-central gamma ($p = 1.5$), gamma ($p = 2$), inverse Gaussian ($p = 3$) and stable distributions ($p > 2$).
- ▶ Notation $Z \sim Tw_p(\mu, \phi)$.

Poisson-Tweedie distribution

- Hierarchical specification:

$$Y|Z \sim \text{Poisson}(Z)$$

$$Z \sim \text{Tw}_p(\mu, \phi).$$

- Probability mass function ($p > 1$)

$$f(y; \mu, \phi, p) = \int_0^{\infty} \frac{z^y \exp^{-z}}{y!} a(z, \phi, p) \exp\{(z\psi - k(\psi))/\phi\} dz.$$

- No closed-form available apart of special case - negative binomial.
- It can be approximated by Monte Carlo integration (difficult and time consuming).
- Distribution, mass, quantile, random generation functions available through the functions `p-`, `d-`, `q-`, `rptweedie()` in **RSBD**

Moments and special cases

- ▶ Marginal mean and variance are easily obtained

$$\begin{aligned}E(Y) &= \mu \\ \text{Var}(Y) &= \mu + \phi\mu^p.\end{aligned}$$

- ▶ Special cases: Hermite ($p = 0$), Neyman-Type A ($p = 1$), Pólya-Aeppli ($p = 1.5$), negative binomial ($p = 2$) and Poisson inverse-Gaussian ($p = 3$).
- ▶ Careful - Hermite is a limit case!
- ▶ p is an index that distinguishes between important distributions.
- ▶ Parameter space of p is not trivially defined, i.e. $p \in 0 \cup [1, \infty)$.
- ▶ Estimation of p works as an automatically model selection.
- ▶ Notation $Y \sim PTW_p(\mu, \phi)$.

Shapes - Dispersion index and power parameter

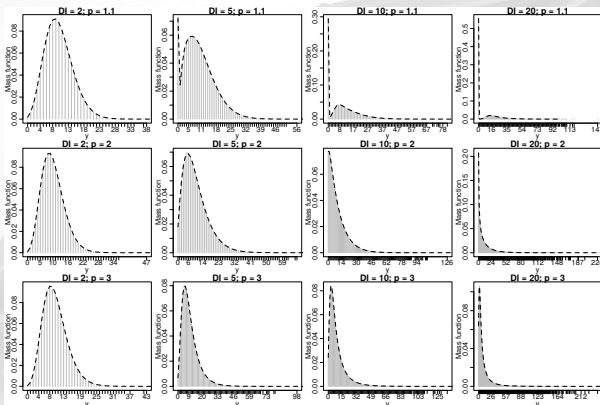


Figure 3. Empirical probability mass function (gray) and approximated probability mass function (black) by dispersion index (DI) and values of the power parameter: Poisson-Tweedie distribution.

Regression models

- ▶ Consider $(Y_i, \mathbf{x}_i), i = 1, \dots, n$, where Y_i 's are iid rv.
- ▶ Full parametric specification:

$$Y_i \sim \text{PTW}_p(\mu_i, \phi).$$

- ▶ Second-moment assumptions:

$$\begin{aligned} E(Y_i) &= \mu_i \\ \text{Var}(Y_i) &= \mu_i + \phi \mu_i^p \end{aligned}$$

where $g(\mu_i) = \eta_i = \mathbf{x}_i^\top \boldsymbol{\beta}$, \mathbf{x}_i and $\boldsymbol{\beta}$ are $(Q \times 1)$ vectors of known covariates and unknown regression parameters.

- ▶ $\text{Var}(Y_i) > 0$, thus $\phi > -\mu_i^{(1-p)} \implies$ under, equi and overdispersion.
- ▶ g link function (log link).
- ▶ Proposed in Bonat et. al. (2017).

Multivariate regression models

- ▶ Consider $(\mathbf{Y}_i, \mathbf{x}_i), i = 1, \dots, n$, where \mathbf{Y}_i 's are i.i.d. random vectors ($R \times 1$).
- ▶ Second-moment assumptions:

$$E(\mathbf{Y}_i) = \boldsymbol{\mu}_i$$

$$\text{Var}(\mathbf{Y}_i) = \boldsymbol{\Sigma}_i = \text{diag}(\boldsymbol{\mu}_i) + V(\boldsymbol{\mu}_i; \boldsymbol{p})^{\frac{1}{2}} \boldsymbol{\Omega}(\boldsymbol{\tau}) V(\boldsymbol{\mu}_i; \boldsymbol{p})^{\frac{1}{2}},$$

where $\boldsymbol{\mu}_i = (g^{-1}(\mathbf{x}_i^\top \boldsymbol{\beta}_1), \dots, g^{-1}(\mathbf{x}_i^\top \boldsymbol{\beta}_R))^\top$ is the $R \times 1$ vector of expected values and \boldsymbol{p} the $R \times 1$ vector of power parameters.

- ▶ $V(\boldsymbol{\mu}_i; \boldsymbol{p}) = \text{diag}(\boldsymbol{\mu}_i^{\boldsymbol{p}})$.
- ▶ $\boldsymbol{\Omega}(\boldsymbol{\tau})$ controls the part of the covariance structure that does not depend on the expected values.
- ▶ $\boldsymbol{\Sigma}_i$ positive definite.
- ▶ Easy interpretation!

Example: Trivariate case

- ▶ Second-moments specification

$$E \begin{bmatrix} Y_{1i} \\ Y_{2i} \\ Y_{3i} \end{bmatrix} = \begin{pmatrix} g(\mu_1) = \mathbf{x}_i^\top \boldsymbol{\beta}_1 \\ g(\mu_2) = \mathbf{x}_i^\top \boldsymbol{\beta}_2 \\ g(\mu_3) = \mathbf{x}_i^\top \boldsymbol{\beta}_3 \end{pmatrix}$$

$$\text{Var} \begin{bmatrix} Y_{1i} \\ Y_{2i} \\ Y_{3i} \end{bmatrix} = \begin{pmatrix} \mu_1 + \mu_1^{p_1} \tau_1 & \sqrt{\mu_1^{p_1} \mu_2^{p_2} \tau_{12}} & \sqrt{\mu_1^{p_1} \mu_3^{p_3} \tau_{13}} \\ \sqrt{\mu_1^{p_1} \mu_2^{p_2} \tau_{12}} & \mu_2 + \mu_2^{p_2} \tau_2 & \sqrt{\mu_2^{p_2} \mu_3^{p_3} \tau_{23}} \\ \sqrt{\mu_1^{p_1} \mu_3^{p_3} \tau_{13}} & \sqrt{\mu_2^{p_2} \mu_3^{p_3} \tau_{23}} & \mu_3 + \mu_3^{p_3} \tau_3 \end{pmatrix}.$$

- ▶ Note that

$$\boldsymbol{\Omega}(\boldsymbol{\tau}) = \begin{pmatrix} \tau_1 & \tau_{12} & \tau_{13} \\ \tau_{12} & \tau_2 & \tau_{23} \\ \tau_{13} & \tau_{23} & \tau_3 \end{pmatrix}.$$

Example: Trivariate case

- Standardized dispersion matrix

$$R(\tau) = \begin{pmatrix} 1 & \frac{\tau_{12}}{\sqrt{\tau_1\tau_2}} & \frac{\tau_{13}}{\sqrt{\tau_1\tau_3}} \\ \frac{\tau_{12}}{\sqrt{\tau_1\tau_2}} & 1 & \frac{\tau_{23}}{\sqrt{\tau_2\tau_3}} \\ \frac{\tau_{13}}{\sqrt{\tau_1\tau_3}} & \frac{\tau_{23}}{\sqrt{\tau_2\tau_3}} & 1 \end{pmatrix},$$

gives us a notion of the correlation between response variables.

- These measures do not depend on the expected values.

Parametrization

- ▶ Let $\mathcal{Y} = (\mathbf{Y}_1^T, \dots, \mathbf{Y}_n^T)^T$ be the stacked vector ($nR \times 1$) of the outcomes.
- ▶ Let $\mathcal{M} = (\boldsymbol{\mu}_1^T, \dots, \boldsymbol{\mu}_n^T)^T$ be the stacked vector ($nR \times 1$) of the expected values.
- ▶ Let $\mathbf{C} = \text{diag}(\boldsymbol{\Sigma}_1, \dots, \boldsymbol{\Sigma}_n)$ denotes a block-diagonal matrix.
- ▶ Thus, the \mathbf{C} matrix is symmetric and $nR \times nR$.
- ▶ Let $\boldsymbol{\beta} = (\boldsymbol{\beta}_1^T, \dots, \boldsymbol{\beta}_R^T)^T$ be a vector $P \times 1$ of regression parameters.
- ▶ Let $\boldsymbol{\lambda} = (p_1, \dots, p_R, \boldsymbol{\tau}^T)^T$ be a $Q \times 1$ vector of dispersion parameters.
- ▶ Multivariate extended Poisson-Tweedie model is specified by two sets of parameters $\boldsymbol{\theta} = (\boldsymbol{\beta}, \boldsymbol{\lambda})$.
- ▶ Quasi-score function for regression parameters.
- ▶ Pearson estimating function for dispersion parameters.

Regression parameters

- ▶ The quasi-score function is defined by,

$$\psi_{\beta}(\beta, \lambda) = \mathbf{D}^{\top} \mathbf{C}^{-1}(\mathcal{Y} - \mathcal{M})$$

where $\mathbf{D} = \nabla_{\beta} \mathcal{M}$ is an $nR \times P$ matrix.

- ▶ The $P \times P$ sensitivity matrix of ψ_{β} is given by

$$S_{\beta} = E(\nabla_{\beta} \psi_{\beta}) = -\mathbf{D}^{\top} \mathbf{C}^{-1} \mathbf{D}. \quad (1)$$

- ▶ The $P \times P$ variability matrix of ψ_{β} is given by

$$V_{\beta} = \text{Var}(\psi_{\beta}) = \mathbf{D}^{\top} \mathbf{C}^{-1} \mathbf{D}. \quad (2)$$

Dispersion parameters

- ▶ The Pearson estimating function is defined by,

$$\psi_{\lambda_i}(\boldsymbol{\beta}, \boldsymbol{\lambda}) = \text{tr}(W_{\lambda_i}(\mathbf{r}^\top \mathbf{r} - \mathbf{C}))$$

where $W_{\lambda_i} = -\frac{\partial \mathbf{C}^{-1}}{\partial \lambda_i}$ and $\mathbf{r} = (\mathcal{Y} - \mathcal{M})$.

- ▶ The entry (i, j) of the $Q \times Q$ sensitivity matrix of $\psi_{\boldsymbol{\lambda}}$ is given by,

$$S_{\lambda_{ij}} = \text{E} \left(\frac{\partial}{\partial \lambda_j} \psi_{\lambda_i} \right) = -\text{tr} (W_{\lambda_i} \mathbf{C} W_{\lambda_j} \mathbf{C}) . \quad (3)$$

- ▶ The entry (i, j) of the $Q \times Q$ variability matrix of $\psi_{\boldsymbol{\lambda}}$ is given by,

$$V_{\lambda_{ij}} = \text{Cov}(\psi_{\lambda_i}; \psi_{\lambda_j}) = 2\text{tr}(W_{\lambda_i} \mathbf{C} W_{\lambda_j} \mathbf{C}) + \sum_{l=1}^{nR} k_l^{(4)} (W_{\lambda_i})_{ll} (W_{\lambda_j})_{ll}$$

where $k^{(4)}$ denotes the fourth cumulant of \mathcal{Y} .

Cross sensitivity and variability matrix

- ▶ The entry (i, j) of the $Q \times P$ cross sensitivity matrix between β and λ is given by,

$$S_{\beta_i \lambda_j} = E \left(\frac{\partial}{\partial \lambda_j} \psi_{\beta_i} \right) = 0. \quad (5)$$

- ▶ The entry (i, j) of the $P \times Q$ cross sensitivity matrix between λ and β is given by,

$$S_{\lambda_i \beta_j} = E \left(\frac{\partial}{\partial \beta_j} \psi_{\lambda_i} \right) = -\text{tr} (W_{\lambda_i} \mathbf{C} W_{\beta_j} \mathbf{C}). \quad (6)$$

Cross sensitivity and variability matrix

- ▶ The entry (i, j) of the $P \times Q$ cross variability matrix between λ and β is given by,

$$V_{\lambda_i \beta_j} = E \left(\sum_{k=1}^{nR} \sum_{l=1}^{nR} \sum_{m=1}^{nR} W_{\lambda_i}^{(lm)} A_m^{(j)} r_k r_l r_m \right), \quad (7)$$

where $A = \mathbf{D}^T \mathbf{C}^{-1}$ and $A^{(j)}$ denotes the j^{th} column of A . In a similar way $W_{\lambda_i}^{(lm)}$ denotes the l^{th} and m^{th} entries of the matrix W_{λ_i} .

Joint sensitivity and variability matrix

- ▶ The joint sensitivity matrix of ψ_β and ψ_λ is given by

$$S_\theta = \begin{pmatrix} S_\beta & S_{\beta\lambda} \\ S_{\lambda\beta} & S_\lambda \end{pmatrix},$$

whose entries are defined in equations (1), (3), (6) and (5).

- ▶ The joint variability matrix of ψ_β and ψ_λ is given by

$$V_\theta = \begin{pmatrix} V_\beta & V_{\beta\lambda} \\ V_{\lambda\beta} & V_\lambda \end{pmatrix},$$

whose entries are defined in equations (2), (4) and (7) above.

Godambe information matrix and asymptotic distribution

- ▶ Denote $\hat{\theta} = (\hat{\beta}, \hat{\lambda})$ the estimating function estimator of θ .
- ▶ The asymptotic distribution of $\hat{\theta}$ is given by

$$\hat{\theta} \sim N_{P+Q}(\theta, J_{\theta}^{-1})$$

where J_{θ}^{-1} is the inverse of the Godambe information matrix,

$$J_{\theta}^{-1} = S_{\theta}^{-1} V_{\theta} S_{\theta}^{-T},$$

where $S_{\theta}^{-T} = (S_{\theta}^{-1})^T$.

Algorithms

- ▶ The *modified chaser*

$$\begin{aligned}\beta^{(i+1)} &= \beta^{(i)} - S_{\beta}^{-1} \psi_{\beta}(\beta^{(i)}, \lambda^{(i)}) \\ \lambda^{(i+1)} &= \lambda^{(i)} - \alpha S_{\lambda}^{-1} \psi_{\lambda}(\beta^{(i+1)}, \lambda^{(i)}).\end{aligned}$$

- ▶ The *reciprocal likelihood* algorithm

$$\begin{aligned}\beta^{(i+1)} &= \beta^{(i)} - S_{\beta}^{-1} \psi_{\beta}(\beta^{(i)}, \lambda^{(i)}) \\ \lambda^{(i+1)} &= \lambda^{(i)} - [\alpha \psi_{\lambda}(\beta^{(i+1)}, \lambda^{(i)})^{\top} \psi_{\lambda}(\beta^{(i+1)}, \lambda^{(i)}) V_{\lambda}^{-1} S_{\lambda} + S_{\lambda}]^{-1} \psi_{\lambda}(\beta^{(i+1)}, \lambda^{(i)})\end{aligned}$$

where α is a *tuning constant*.

- ▶ Easy implementation through the `mcglm` package (Bonat, 2018).
- ▶ Special case of a multivariate covariance generalized linear model (Bonat and Jørgensen, 2016).

Multivariate hypothesis test

- ▶ The general linear hypothesis may be stated as

$$H_0 : L\beta = \mathbf{0},$$

where $L = G \otimes F$.

- ▶ The $G (R \times R)$ states between responses hypotheses.
- ▶ The $F (s \times p)$ states between treatments hypotheses.
- ▶ We assume equal linear predictor for all responses, thus L is a $(sR \times P)$ matrix with s denoting the number of linear constraints.
- ▶ The alternative hypothesis may be stated in the form,

$$H_1 : L\beta = \mathbf{n},$$

where \mathbf{n} is not the null vector.

Multivariate hypothesis test

- ▶ Wald statistics given by

$$W_s = (\mathbf{L}\boldsymbol{\beta})^\top (\mathbf{L}\mathbf{J}_\beta^{-1}\mathbf{L}^\top)^{-1}(\mathbf{L}\boldsymbol{\beta}),$$

under the null hypothesis is asymptotically chi-squared distributed with sR degrees of freedom.

- ▶ Performing test for all response variables as well as between combination of them.
- ▶ All possible contrasts between treatment levels (multivariate multiple hypothesis tests).
- ▶ In the Gaussian case the Wald test corresponds to the Hotelling-Lawley statistics.

Data set 1: Model specification

- ▶ Multivariate count regression model

$$E(\mathbf{Y}_i) = \boldsymbol{\mu}_i = \{\exp(\mathbf{x}_{i1}^\top \boldsymbol{\beta}_1), \dots, \exp(\mathbf{x}_{i5}^\top \boldsymbol{\beta}_5)\}$$

$$\text{Var}(\mathbf{Y}_i) = \boldsymbol{\Sigma}_i = \text{diag}(\boldsymbol{\mu}_i) + V(\boldsymbol{\mu}_i; \boldsymbol{\rho})^{\frac{1}{2}} \boldsymbol{\Omega}(\boldsymbol{\tau}) V(\boldsymbol{\mu}_i; \boldsymbol{\rho})^{\frac{1}{2}}$$

- ▶ Dispersion matrix

$$\boldsymbol{\Omega}(\boldsymbol{\tau}) = \begin{bmatrix} \tau_1 & \tau_{12} & \tau_{13} & \tau_{14} & \tau_{15} \\ \tau_{12} & \tau_2 & \tau_{23} & \tau_{24} & \tau_{25} \\ \tau_{13} & \tau_{23} & \tau_3 & \tau_{34} & \tau_{35} \\ \tau_{14} & \tau_{24} & \tau_{34} & \tau_4 & \tau_{45} \\ \tau_{15} & \tau_{25} & \tau_{35} & \tau_{45} & \tau_5 \end{bmatrix}.$$

- ▶ Log link function.
- ▶ Poisson-Tweedie dispersion function.

Data set 1: Dispersion structure

- Dispersion and power parameter estimates.

Table 1. Dispersion parameter estimates and standard errors (SE).

	Ndoc	Nndoc	Nmed	Nhosp	Nadm
$\hat{\rho}_i$	Estimate (SE) 1.9042 (0.1121)	Estimate (SE) 1.7495 (0.1309)	Estimate (SE) 1.3153 (0.2442)	Estimate (SE) 1.4184 (0.3702)	Estimate (SE) 1.7595 (0.3179)
$\hat{\tau}_i$	1.2496 (0.2052)	7.7676 (1.3134)	0.2928 (0.0526)	21.3497 (4.4126)	1.1121 (0.5192)

- Nmed and Nhosp present excess of zeros ($\hat{\rho}$ near 1).
- Weak over-dispersion Ndoc, Nmed and Nadm.
- High over-dispersion Nhosp and Nndoc.
- Neyman-type A and negative binomial.

Data set 1: Multivariate hypothesis tests

- ▶ MANOVA-type test for multivariate count data.

Effects	Df	Hotelling-Lawley	Chi-squared	p-value
Intercept	5	0.3689	1914.80	<0.0001
sex	5	0.0322	167.5387	<0.0001
age	5	0.0410	213.1572	<0.0001
Levyplus	5	0.0046	23.9239	0.0002
freepoor	5	0.0012	6.5851	0.2533
freerepa	5	0.0054	28.2852	<0.0001
illness	5	0.1320	685.1295	<0.0001
actdays	5	0.1025	531.9883	<0.0001
hscore	5	0.0080	41.7531	<0.0001

Data set 1: Standardized dispersion matrix

- ▶ Standardized dispersion matrix (lower) and standard errors (upper).

$$\hat{R} = \begin{bmatrix} - & 0.0141 & 0.01420 & 0.0152 & 0.0142 \\ 0.0426 & - & 0.0145 & 0.0159 & 0.0144 \\ 0.1198 & 0.0788 & - & 0.0154 & 0.0140 \\ 0.0609 & 0.0680 & 0.0664 & - & 0.055 \\ 0.0893 & 0.0619 & 0.0690 & 0.5142 & - \end{bmatrix}$$

Data set 2: Model specification

- ▶ Linear predictor for each response variable

$$y_i \sim \text{Bare.ground} + \text{Shrub.cover.}$$

- ▶ 60 regression parameters.
- ▶ 20 dispersion parameters.
- ▶ 190 cross-dispersion parameters.
- ▶ Total number of parameters: 270.

Data set 2: R code

► R code

```
# Loading extra packages
require(mcgln)
require(mvabund)
require(Matrix)
require(corrplot)

# Loading data set
data(antTraits)
y <- antTraits$abund[,1:20] # Selecting response variables
names(y) <- paste0("y",1:20)
X <- antTraits$env[,c(3,5)] # Selecting covariates
data <- data.frame(y, X)

# Linear predictor
lp <- paste(names(y), "~", "Bare.ground + Shrub.cover")
form <- lapply(lp, formula)

# Matrix linear predictor
Z0 <- mc_id(data)
```

► Fitting the model to data.

```
# Fitting multivariate model
fit1 <- mcglm(linear_pred = c(form), matrix_pred = rep(list(Z0), 20),
             link = rep("log", 20), variance = rep("poisson_tweedie", 20),
             control_algorithm = list(max_iter = 100),
             data = data)
```

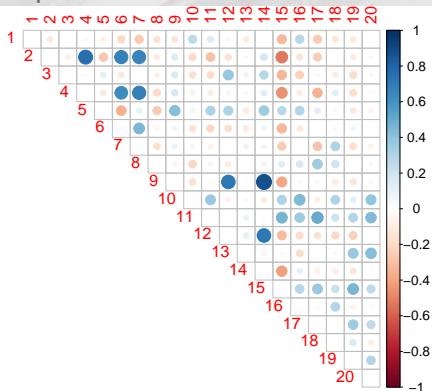
Data set 2: Results

- Multivariate hypothesis test.

```
manova.mcglm(fit1)
```

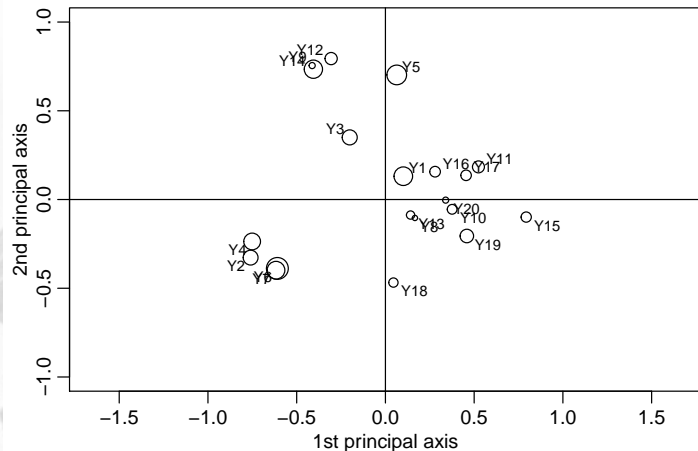
##	Effects	Df	Hotelling.Lawley	Chi.squared	p_value
## 1	Intercept	20	60.876617	1826.29851	0e+00
## 2	Bare.ground	20	2.110956	63.32867	2e-06
## 3	Shrub.cover	20	2.636528	79.09584	0e+00

- Standardized dispersion matrix.



Data set 2: Principal components

Points proportional to logarithm of body length.



Main results

- ▶ Flexible multivariate statistical model to deal with count data.
- ▶ Second-moment assumptions.
- ▶ General framework for estimation based on estimating function.
- ▶ Efficient algorithms for estimation.
- ▶ Asymptotic theory (Godambe Information).
- ▶ General software implementation in R (`mcglm` package).
- ▶ Extension of orthodox MANOVA to deal with count data.

Topics for research

- ▶ MANOVA-type tests

1. Hotelling-Lawley statistics → Wald test.
2. Pillai statistics → Score test.
3. Wilks statistics → likelihood ratio test.

- ▶ Can we have analogous to

1. principal components,
2. factor analysis,
3. correspondence analysis,
4. canonical correlation analysis and
5. Redundancy analysis

in the context of count data?

- ▶ Perhaps, all of them based on the standardized dispersion matrix. Does it make sense?
- ▶ How to deal with high-dimensional data in both outcomes and covariates?
- ▶ How to deal with missing and/or censored data?

Main References

Bonat, W. H. ; Jørgensen, B., *Multivariate Covariance Generalized Linear Models*. Journal of the Royal Statistical Society: Series C (Applied Statistics) 65(5): 649–675, 2016.

Bonat, W. H. ; Jørgensen, B. ; Kokonendji, C., C. ; Hinde, J. ; Demétrio, C. G. *Extended Poisson-Tweedie: properties and regression models*. Statistical Modelling, 2017.

Bonat, W. H. *Multiple regression models in R: The mcglm package*. Journal of Statistical Software, 2018.

Further references

- Kokonenji, C. C.; Dossou-Gbété, S.; Demétrio, C. G. B. *Some discrete exponential dispersion models: Poisson-Tweedie and Hinde-Demétrio classes*. Statistics and Operations Research Transactions, v.28, p.201-214, 2004.
- Rigby, R. A.; Stasinopoulos, D. M.; Akantziliotou, C. *A framework for modelling overdispersed count data, including the Poisson-shifted generalized inverse Gaussian distribution*. Computational Statistics and Data Analysis, v.53, p.381-393, 2008.
- Kolossiatis, M.; Griffin, J.E.; Steel, M.F.J. *Modeling overdispersion with the normalized tempered stable distribution*. Computational Statistics and Data Analysis, v.55, p.2288-2301, 2011.
- Barabesi, L.; Becatti, C.; Marcheselli, M. *The tempered discrete Linnik distribution*. ArXiv e-prints, 1605.02326, 2016.
- Sellers, K. F. and Shmueli, G. *A flexible regression model for count data*. Ann. Appl. Stat., v.4, p.943-961, 2010.
- Zeviani, W. M.; Ribeiro Jr, P. J.; Bonat, W. H.; Shimakura, S. E.; Muniz, J. A. *The Gamma-count distribution in the analysis of experimental underdispersed data*. Journal of Applied Statistics, v.41, p.2616-2626, 2014.
- Kalktawi, H. S.; Vinciotti, V.; Yu, K.A *Simple and Adaptive Dispersion Regression Model for Count Data*. ArXiv e-prints, 1511.00634, 2015.
- Zhu, R.; Joe, H. *Modelling heavy-tailed count data using a generalised Poisson-inverse Gaussian family*. Statistics & Probability Letters, v.79, 2009.
- Zeileis, A.; Kleiber, C.; Jackman, S. *Regression Models for Count Data in R*. Journal of Statistical Software, v.27, p.1-25, 2008.
- Loeys, T.; Moerkerke, B.; De Smet, O.; Buysse, A. *The analysis of zero-inflated count data: Beyond zero-inflated Poisson regression*. British Journal of Mathematical and Statistical Psychology, v.65, p.163-180, 2008.
- Sellers, K. F.; Raim, A. *A flexible zero-inflated model to address data dispersion*. Computational Statistics and Data Analysis, v.99, p.68-80, 2016.
- Inouye, D.; Yang, E.; Allen, G.; Ravikumar, P. *A review of multivariate distributions for count data derived from the Poisson distributions*. Wiley Interdiscip. v.9, e1398, 2018.
- Jørgensen, B.; Lauritzen, S. L. *Multivariate Dispersion Models*. Journal of Multivariate Analysis, v.74, p.267-281, 2000.
- Jørgensen, B.; Martinez, J. R. *Multivariate Exponential Dispersion Models*. WSPC - Proceedings, 2012.
- Jørgensen, B. *Construction of multivariate dispersion models*. Brazilian Journal of Probability and Statistics, v.27, p.285-309, 2013.
- Jørgensen, B. *The Theory of Dispersion models*. Chapman and Hall, London, 1997.

Contact

- ▶ Work in progress ...
- ▶ Joint work with Célestin C. Kokonendji.
- ▶ Name: Wagner Hugo Bonat
- ▶ e-mail: wbonat@ufpr.br
- ▶ Webpages
 1. <https://cran.r-project.org/web/packages/mcglm>
 2. <https://github.com/wbonat/mcglm>
 3. www.leg.ufpr.br/papercompanions
 4. <http://www.leg.ufpr.br/~wagner>
- ▶ Thank you !