

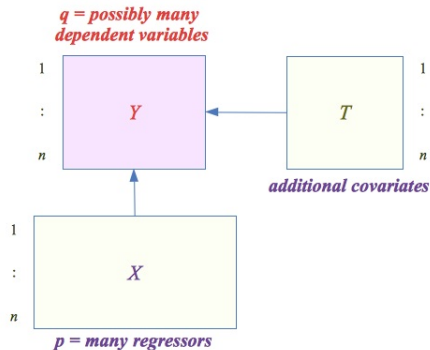


SCGLR: a component-based multivariate regression method to model species distributions.

Frédéric Mortier
with Xavier Bry, Guillaume Cornu & Catherine Trottier

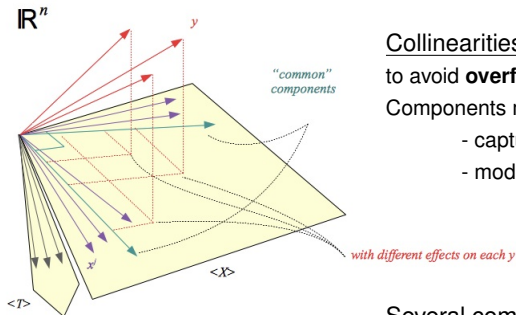
International Statistical Ecology Conference
Montpellier
July 2014

Problematic



⇒ Question: What in X may predict what in Y ?

↪ Approach: *Dimension reduction by construction of components*



Collinearities:

to avoid **overfitting**, search for **components**.

Components must:

- capture enough variance in X ,
- model and predict y .

Several components:

to avoid **redundancy**, search for uncorrelation.

→ constraint of construction: **orthogonality**.

Multiple y :

same components,

but each y with its **own regression coefficients**.

Exponential family distributed

→ **generalized linear regression**.

- **First component** is a **compromise** between the direction of X that best predicts y and the first principal component (PC) of X .

↪ *Criterion:* $\max_{\|u\|^2=1} [\text{cov}(y, Xu)]$

$$\max_{\|u\|^2=1} \left[\sqrt{\text{var}(y)} \sqrt{\text{var}(Xu)} \text{corr}(y, Xu) \right]$$

↪ *Program to solve:* $P_1 : \max_{\|u\|^2=1} [\langle y, Xu \rangle_w]$

- **Further components:** W -orthogonality of components is ensured using the part of X that is not yet used, i.e. the residuals of X regressed on previous components.

- **First component** can be obtained using several equivalent programs:

$$\hookrightarrow P_2 : \max_{\|u\|^2, \|v\|^2=1} [\langle Xu, Yv \rangle_W]$$

$$\hookrightarrow P_3 : \max_{\|u\|^2=1} [\sum_{k=1}^q \langle Xu, y^k \rangle_{W_k}^2]$$

P_3 is adapted to the case of multiple weighting:

$$\hookrightarrow P_4 : \max_{\|u\|^2=1} [\sum_{k=1}^q \langle Xu, y^k \rangle_{W_k}^2]$$

\implies *Solution: eigenvector associated to largest eigenvalue of:*

$$A = X' \Omega X \text{ with } \Omega = \sum_{k=1}^q W_k y^k y'^k W_k$$

- **Further components:** idem, subject to constraint of orthogonality to previous components.

In the GLM, linear predictors are constrained to be collinear to one another:

$$\forall k = 1, q: \quad \eta^k = X\beta_k + T\delta_k = X\gamma_k u + T\delta_k$$

→ **modified Fisher Scoring Algorithm:**

u and $\gamma = (\gamma_k)_{k=1,q}$ estimated iterating an alternated least squares two steps sequence:

- Given γ , working data $(z^k)_k$ is regressed on matrix $[\gamma \otimes X, 1_q \otimes T]$ with respect to working matrix $W = \text{diag}[W_k]_k$

- coefficient vectors $\hat{u}, \hat{\delta} = (\hat{\delta}_k)_k$

- \hat{u} made unit norm → updated u

- (2) Given Xu , each working vector z^k is regressed on $[Xu, T]$ with respect to working matrix W_k

- updated γ_k, δ_k

Step t of the FSA:

$$\begin{aligned}
 & \min_{\gamma, u: u'=1} \left[\sum_k \|z^{k[t]} - X\gamma_k u\|_{W_k^{[t]}}^2 \right] \\
 & \Leftrightarrow \min_{u: u'=1} \left[\sum_k \|z^{k[t]} - \Pi_{Xu} z^{k[t]}\|_{W_k^{[t]}}^2 \right] \\
 & \Leftrightarrow \max_{u: u'=1} \left[\sum_k \|z^{k[t]}\|_{W_k^{[t]}}^2 \cos^2_{W_k^{[t]}}(z^{k[t]}, Xu) \right]
 \end{aligned}$$

is **replaced** by: $\max_{u: u'=1} \left[\sum_k \|z^{k[t]}\|_{W_k^{[t]}}^2 \cos^2_{W_k^{[t]}}(z^{k[t]}, Xu) \quad \|Xu\|_{W_k^{[t]}}^2 \right]$

equivalent to:

$$\max_{u: u'=1} \left[\sum_k \langle z^{k[t]}, Xu \rangle_{W_k^{[t]}}^2 \right]$$

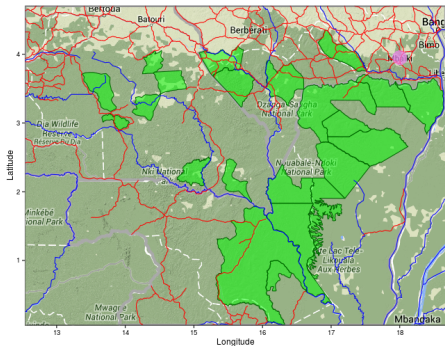
= local extended PLS2

\Rightarrow Solution: eigenvector associated to largest eigenvalue of:

$$A = X' \Omega^{[t]} X \text{ with } \Omega^{[t]} = \sum_{k=1}^q W_k^{[t]} z^{k[t]} z'^{k[t]} W_k^{[t]}$$

Application

Abundance of tropical tree species (CoForChange project)



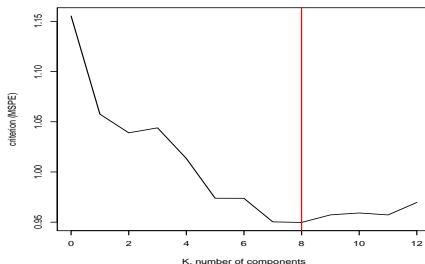
- all trees with diameter higher than 30 cm
- more than 120,000 plots of 0.5 ha
- more than 200 genera
- soil, rainfall, human disturbances, vegetation activity (EVI) maps available

Application II

Select number of component: a cross-validation approach

```
> library(SCGLR)
> genus.cv <- scglrCrossVal(formula=form, data=genus, family=fam, K=12,
+   offset=genus$surface)
>
> mean.crit <- t(apply(genus.cv, 1, function(x) x/mean(x)))
> mean.crit <- apply(mean.crit, 2, mean)
> K.cv <- which.min(mean.crit)-1
> cat("Best number of components: ", K.cv)
```

Best number of components: 8

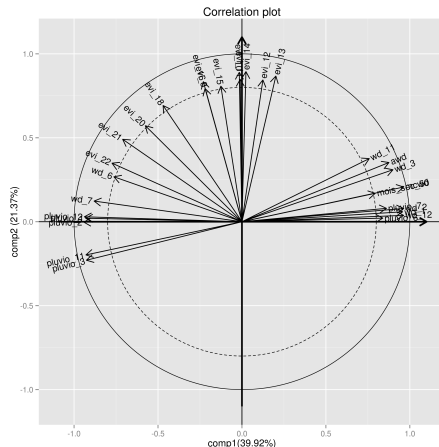
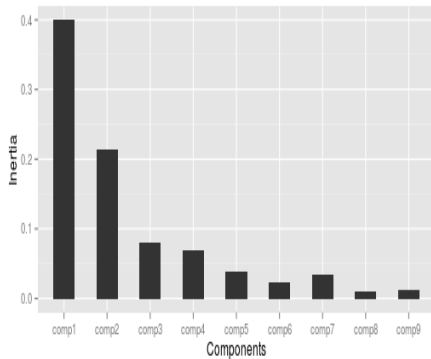


Application III

Fitting and Plots

```
> genus.scglr <- scglr(formula=form, data=genus, family=fam, K=K.cv,  
+ offset=genus$surface)
```

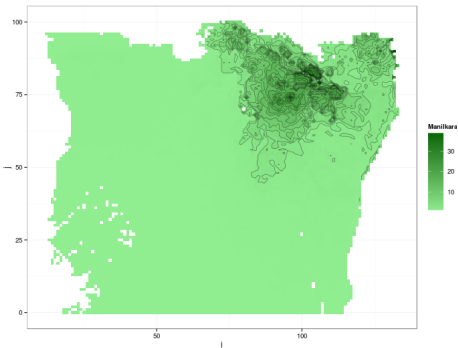
Inertia per component



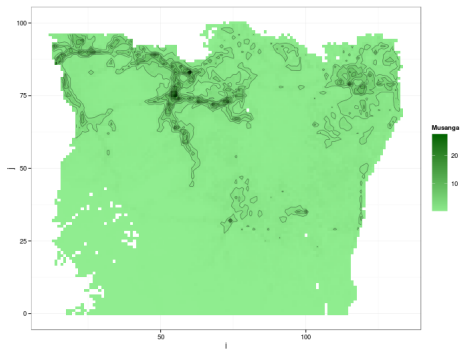
Application IV

Prediction of two genera

Manilkara mabokeensis



Musanga cecropioides



Ongoing works

- New alternate optimization algorithms: Iterative Normalized Gradient
- Multi-table (Theme) support
- SCGLR packages new versions
 - Enhancements for plot customization
 - New distribution families (Negative-Binomial, Exponential, Inverse Gaussian)
 - Multi-theme
 - ...

- 1 X. Bry, C. Trottier, T. Verron and F. Mortier (2013). Supervised component generalized linear regression using a PLS-extension of the Fisher scoring algorithm. *Journal of Multivariate Analysis*, 119(0), 47.
- 2 S. Gourlet-Fleury et al. (2009–2014) CoForChange project:
<http://www.coforchange.eu/>.
- 3 F. Mortier, C. Trottier, G. Cornu and X. Bry (2014). SCGLR: Supervised Component Generalized Linear Regression (SCGLR). *R package version 1.2*. <http://CRAN.R-project.org/package=SCGLR>
- 4 F. Mortier, C. Trottier, G. Cornu and X. Bry (2014). SCGLR - An R package for Supervised Component Generalized Linear Regression. *Journal of Statistical Software*. submitted