

pyPLNmodels: getting started.

Bastien Batardière, Joon Kwon and Julien Chiquet

Université Paris-Saclay, AgroParisTech, INRAE, UMR MIA Paris-Saclay

June 9, 2023

Modelling

- ▶ Dataset:
 - ▶ \mathbf{Y} : $n \times p$ count matrix
 - ▶ \mathbf{X} : $n \times d$ covariates
- ▶ Parameter $\theta = (\beta, \Sigma)$
 - ▶ $\beta \in \mathbb{R}^{d \times p}$ regression parameter
 - ▶ $\Sigma = CC^T$, $C \in \mathbb{R}^{p \times q}$ covariance matrix.
- ▶ Model:

$$W_i \sim \mathcal{N}(\mathbf{0}, I_q)$$

$$Z_i = \beta^T X_i + CW_i$$

$$(Y_{ij} \mid Z_{ij}) \sim \mathcal{P}(\exp(Z_{ij}))$$

- ▶ When $p = q$: unrestricted model: PIn class in pyPLNmodels
- ▶ When $p \gg q$: Σ has low rank: PInPCA class in pyPLNmodels

- ▶ Model:

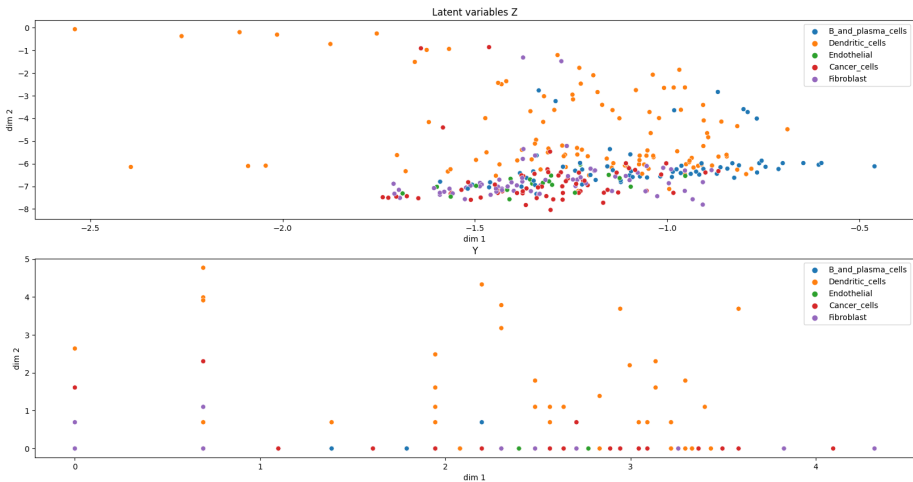
$$W_i \sim \mathcal{N}(\mathbf{0}, I_q)$$

$$Z_i = \beta^\top X_i + CW_i$$

$$(Y_{ij} \mid Z_{ij}) \sim \mathcal{P}(\exp(Z_{ij}))$$

- ▶ Goal of the package: retrieve back the latent representation Z_i :
 - ▶ $W_i \mid Y_i$ is enough for PlnPCA
 - ▶ $Z_i \mid Y_i$ for Pln
- ▶ Problem: both $W_i \mid Y_i$ and $Z_i \mid Y_i$ have unknown distribution.
- ▶ \implies gaussian variational approximation.
 - ▶ PlnPCA: $W_i \mid Y_i \sim \mathcal{N}(M_i, \text{diag}(S_i))$, $M_i, S_i \in \mathbb{R}^q$
 - ▶ Pln: $Z_i \mid Y_i \sim \mathcal{N}(M_i, \text{diag}(S_i))$, $M_i, S_i \in \mathbb{R}^p$

First dimension against second dimension in dimension p = 5



PC 1 agains PC 2 in dimension p = 200.

