

Zero-inflation in the Multivariate Poisson Lognormal Family

Bastien Batardière, François Gindraud, Julien Chiquet and
Mahendra Mariadassou

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

May 22, 2024

- ▶ In Single cell analysis, it is usual to deal with high-dimensional count data:

$$\mathbf{Y} = \begin{pmatrix} 12 & 0 & \dots & 0 & 9 \\ 2 & 0 & \dots & 0 & 0 \\ \vdots & & & & \vdots \\ 341 & 5 & \dots & 1 & 0 \end{pmatrix}$$

- ▶ Y_{ij} : count of transcript j in cell i
- ▶ **Non-continuous data** \implies Linear Gaussian models do not apply
- ▶ **High percentage of zeros ($\approx 90\%$)** \implies zero-inflation is needed

- ▶ Dataset:
 - ▶ $\mathbf{Y} : n \times p$ count matrix ($n \approx p \approx 10^4$)
 - ▶ $\mathbf{X} : n \times d$ or $d \times p$ covariates ($d \approx 10$)
- ▶ Parameter $\theta = (\mathbf{B}, \mathbf{\Sigma}, \boldsymbol{\pi})$
 - ▶ $\mathbf{B} \in \mathbb{R}^{d \times p}$ regression coefficient.
 - ▶ $\mathbf{\Sigma} \in \mathbb{R}^{p \times p}$ covariance matrix.
 - ▶ $\boldsymbol{\pi} \in \mathbb{R}^{n \times p}$ zero-inflation coefficient.
- ▶ Model:

$$\mathbf{W}_i \sim \mathcal{B}(\boldsymbol{\pi}_i)$$

$$\mathbf{Z}_i \sim \mathcal{N}(\mathbf{X}\mathbf{B}, \mathbf{\Sigma})$$

$$(Y_{ij} \mid Z_{ij}, W_{ij}) \sim (1 - W_{ij}) \mathcal{P}(\exp(Z_{ij}))$$

The zero-inflation can take several forms:

$$\pi_{ij} = \pi \in [0, 1] \quad (\text{non-dependent})$$

$$\pi_{ij} = \sigma(\mathbf{XB}^0)_{ij}, \mathbf{X} \in \mathbb{R}^{n \times d}, \mathbf{B}^0 \in \mathbb{R}^{d \times p} \quad (\text{column-wise dependence})$$

$$\pi_{ij} = \sigma(\mathbf{B}^0\mathbf{X})_{ij}, \mathbf{B}^0 \in \mathbb{R}^{n \times d}, \mathbf{X} \in \mathbb{R}^{d \times p} \quad (\text{row-wise dependence})$$