

# Mixture of Poisson PCA for joint clustering and dimension reduction of count-data

---

Nicolas Jouvin, Julien Chiquet & others ( Disclaimer: work-in-progress)

Statlearn - Montpellier, jeudi 06 avril 2023



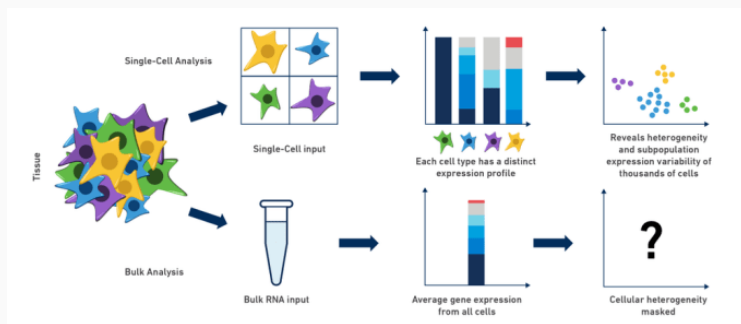
Motivations: clustering &  
visualization of count-data

---

# Count-data arise in many modern scientific field

## Example 1: biology & single-cell RNAseq

Group similar cells based on their *gene* expression profile



Source: 10x Genomics

# Count-data arise in many modern scientific field

## Example 2: ecology

Group ecological sites based on their *species* abundance (thx @ Eleni Matechou)

Site	Psy	Hym	Ath	Cea	...	Set
1	27	2	0	0	...	4
2	220	15	0	0	...	0
3	1173	0	1	2	...	71
4	2671	12	1	3	...	49
⋮						

# Count-data arise in many modern scientific field

## Example 3: document clustering

Group similar texts based on their *words* profile

### MICROBIOPSIE SOUS ECHOGRAPHIE DU SEIN DROIT

#### MACROSCOPIE

Cinq fragments de 5 à 15 mm

#### MICROSCOPIE

Les prélèvements examinés correspondent à des fragments de tissu mammaire remanié par une prolifération tumorale dont les caractères morphologiques sont ceux d'un adénocarcinome canalaire infiltrant. Cette lésion est peu différenciée, d'architecture essentiellement trabéculaire. Les cellules néoplasiques comportent des atypies nucléaires marquées. L'index mitotique est élevé (22 mitoses sur 10 champs au grossissement 400). Deux fragments de 8 et 15 mm. Adénocarcinome mammaire de type canalaire infiltrant peu différencié. Grade histopronostique (EE) : III Index mitotique élevé.

### MICROBIOPSIE SOUS ECHOGRAPHIE DU SEIN DROIT

#### MACROSCOPIE

Cinq fragments de 5 à 15 mm

#### MICROSCOPIE

L'examen histologique met en évidence des lésions tumorales dont les caractères morphologiques sont ceux d'un carcinome canalaire infiltrant moyennement différencié. La lésion est d'architecture trabéculaire et glandulaire. Les cellules sont caractérisées par des atypies cytonucléaires modérées. L'activité mitotique est faible : deux mitoses ont été dénombrées sur dix champs au grossissement 400. Ces lésions sont associées à un stroma dense fibreux. Elles infiltrent le tissu adipeux. Deux séries de prélèvements ont été confisquées : A - 1er tour : onze cylindres biopsiques mesurant 10 à 30 mm de long. B - 2ème tour : onze cylindres biopsiques mesurant 5 à 30 mm de long. Adénocarcinome mammaire de type canalaire infiltrant. Grade histopronostique (EE) I. Index mitotique faible.

### MACROBIOPSIE DU SEIN GAUCHE

#### MACROSCOPIE

3 fragments de 7 à 15 mm

#### MICROSCOPIE

Tous les prélèvements ont un aspect histologique similaire. Ils correspondent à des fragments de tissu mammaire remanié par des lésions de mastose fibreuse commune. Présence d'un discret infiltrat inflammatoire. On retrouve également quelques microcalcifications. L'un des prélèvements cryo-prélevés sera analysé histologiquement et un compte rendu complémentaire adressé ultérieurement. Trois fragments de 7 à 15 mm. Lésions de mastose fibreuse. Le prélèvement paraît peu significatif. Une analyse complémentaire sur le prélèvement cryo-prélevé sera réalisée.

Doc 1	"Lésions cancéreuses (...) carcinome canalaire"
Doc 2	"Lésions cancéreuses (...) carcinome lobulaire"
...	...
Doc n	"Lésions bénignes (...) métaplasie"

# Statistical context & problematics

Multivariate data  $\mathbf{Y} = \{\mathbf{y}_1, \dots, \mathbf{y}_n\}$

- ▶ discrete:  $\mathbf{y}_i \in \mathbb{N}^p$
- ▶ possibly highly-dimensional ( $p \gg n$ ) or with small sample size.
- ▶ No Gaussianity, sparse data, over-dispersion

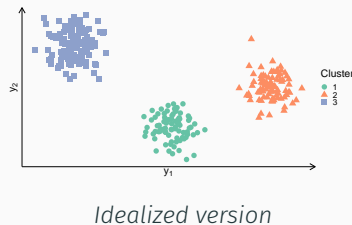
# Statistical context & problematics

Multivariate data  $\mathbf{Y} = \{\mathbf{y}_1, \dots, \mathbf{y}_n\}$

- ▶ discrete:  $\mathbf{y}_i \in \mathbb{N}^p$
- ▶ possibly highly-dimensional ( $p \gg n$ ) or with small sample size.
- ▶ No Gaussianity, sparse data, over-dispersion

**Our goal:** unsupervised data analysis

- ① Clustering  $\rightsquigarrow$  partition
- ② Dimension reduction  $\rightsquigarrow$  visualization



**How ?** design a statistical model on  $\mathbf{Y}$  integrating ① & ②

Count data modeling: the Poisson Log Normal (PLN) family

A mixture of PLN-PCA for joint clustering and dimension reduction

Inference

Conclusion



# Count data modeling: the Poisson Log Normal (PLN) family

---



# What If ? Everything was Gaussian

In a Gaussian world, we would love to use the GLM framework

$$\mathbf{y}_i = \underbrace{\mathbf{x}_i^\top \mathbf{B}}_{\text{covariates}} + \underbrace{\mathbf{o}_i}_{\text{offset}} + \boldsymbol{\varepsilon}_i, \quad \boldsymbol{\varepsilon}_i \sim \mathcal{N}_p(\mathbf{0}_p, \boldsymbol{\Sigma})$$

# What If ? Everything was Gaussian

In a Gaussian world, we would love to use the GLM framework

$$\mathbf{y}_i = \underbrace{\mathbf{x}_i^\top \mathbf{B}}_{\text{covariates}} + \underbrace{\mathbf{o}_i}_{\text{offset}} + \boldsymbol{\varepsilon}_i, \quad \boldsymbol{\varepsilon}_i \sim \mathcal{N}_p(\mathbf{0}_p, \boldsymbol{\Sigma})$$

Pros:

- account for offset  $\mathbf{o}_i$  and covariates  $\mathbf{x}_i$  when available
- $\boldsymbol{\Sigma}$  capture all the remaining covariance
- flexible, generalizes

# What If ? Everything was Gaussian

In a Gaussian world, we would love to use the GLM framework

$$\mathbf{y}_i = \underbrace{\mathbf{x}_i^\top \mathbf{B}}_{\text{covariates}} + \underbrace{\mathbf{o}_i}_{\text{offset}} + \epsilon_i, \quad \epsilon_i \sim \mathcal{N}_p(\mathbf{0}_p, \Sigma)$$

Pros:

- account for offset  $\mathbf{o}_i$  and covariates  $\mathbf{x}_i$  when available
- $\Sigma$  capture all the remaining covariance
- flexible, generalizes

However... for multivariate counts

- Data transformation (log, normalization): quick and dirty
- Non-Gaussian multivariate distributions: do not scale to data dimension yet

## The vanilla PLN model: a Gaussian love story

Gaussian latent layer encoding for Poisson (log-)intensities (Aitchison et al. 1989)

$$\begin{aligned}\boldsymbol{\eta}_i &\sim \mathcal{N}_p(\boldsymbol{o}_i + \boldsymbol{x}_i^\top \boldsymbol{B}, \boldsymbol{\Sigma}), && \text{(param)} \\ \boldsymbol{y}_i \mid \boldsymbol{\eta}_i &\sim \otimes_j \mathcal{P}(\exp(\eta_{ij})) && \text{(emission)}\end{aligned} \tag{PLN}$$

# The vanilla PLN model: a Gaussian love story

Gaussian latent layer encoding for Poisson (log-)intensities (Aitchison et al. 1989)

$$\begin{aligned}\boldsymbol{\eta}_i &\sim \mathcal{N}_p(\mathbf{o}_i + \mathbf{x}_i^\top \mathbf{B}, \boldsymbol{\Sigma}), && \text{(param)} \\ \mathbf{y}_i \mid \boldsymbol{\eta}_i &\sim \otimes_j \mathcal{P}(\exp(\eta_{ij})) && \text{(emission)}\end{aligned} \tag{PLN}$$

## Parameters $\theta$

- $\mathbf{B}$ , the regression parameters
- $\boldsymbol{\Sigma}$ , the variance-covariance matrix

## Observations

- $\mathbf{Y}$ : the count-data matrix  $n \times p$
- $\mathbf{X}$ : the covariates matrix  $n \times d$
- $\mathbf{O}$ : the offsets matrix  $n \times p$

## Over-dispersion

- mean:  $\mathbb{E}(Y_{ij}) = \exp(o_{ij} + \mathbf{x}_i^\top \mathbf{B}_{.j} + \sigma_{jj}/2) > 0$
- variance:  $\mathbb{V}(Y_{ij}) = \mathbb{E}(Y_{ij}) + \mathbb{E}(Y_{ij})^2 (e^{\sigma_{jj}} - 1) > \mathbb{E}(Y_{ij})$
- covariance:  $\text{Cov}(Y_{ij}, Y_{ik}) = \mathbb{E}(Y_{ij})\mathbb{E}(Y_{ik}) (e^{\sigma_{jk}} - 1)$

**Underlying assumption:** correlations are captured in the latent layer

**Flexible**, with many extensions: **clustering**, **dimension reduction**, network inference, ...  
(Chiquet et al. 2021) & an R package **PLNmodels**



# Extension to clustering: mixture modeling

Goal find a partition  $\mathbf{Z}$  into  $K$  groups

How ? Use a Gaussian mixture in the latent layer

$$\boldsymbol{\eta}_i \sim \sum_{k=1}^K \pi_k \mathcal{N}_p(\mathbf{o}_i + \mathbf{x}_i^\top \mathbf{B} + \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

$$\mathbf{z}_i \sim \mathcal{M}_K(1, \boldsymbol{\pi}), \quad (\text{membership})$$

$$\boldsymbol{\eta}_i \mid \{\mathbf{z}_i = k\} \sim \mathcal{N}_p(\mathbf{o}_i + \mathbf{x}_i^\top \mathbf{B} + \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k), \quad (\text{param}) \quad (\text{PLNmixture})$$

$$\mathbf{y}_i \mid \boldsymbol{\eta}_i \sim \otimes_j \mathcal{P}(\exp(\eta_{ij})) \quad (\text{emission})$$

$$\boldsymbol{\theta} = \{\pi_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k, \mathbf{B}\}$$

Clustering ? Via the *posterior* distribution

$$p(\mathbf{Z} \mid \mathbf{Y}, \boldsymbol{\theta})$$

Problem :  $p^2$  parameters in  $\Sigma$ , what if  $p$  is "large" ?

## Extension to dimension reduction with the PLN-PCA (Chiquet et al. 2018)

**Problem** :  $p^2$  parameters in  $\Sigma$ , what if  $p$  is "large" ?

**Regularization** Low-rank factorization of  $\Sigma = \mathbf{C}\mathbf{C}^\top$ , with  $\mathbf{C}$  a  $p \times q$  matrix

## Extension to dimension reduction with the PLN-PCA (Chiquet et al. 2018)

Problem :  $p^2$  parameters in  $\Sigma$ , what if  $p$  is "large" ?

Regularization Low-rank factorization of  $\Sigma = \mathbf{C}\mathbf{C}^\top$ , with  $\mathbf{C}$  a  $p \times q$  matrix

Probabilistic PCA formulation (Tipping et al. 1999b)

$$\begin{aligned} \mathbf{w}_i &\sim \mathcal{N}_p(\mathbf{0}_q, \mathbf{I}_q) && \text{(low-dimensional subspace)} \\ \boldsymbol{\eta}_i &= \mathbf{o}_i + \mathbf{x}_i^\top \mathbf{B} + \mathbf{C}\mathbf{w}_i && \text{(linear transformation)} \\ \mathbf{y}_i \mid \boldsymbol{\eta}_i &\sim \otimes_j \mathcal{P}(\exp(\eta_{ij})) && \text{(emission)} \end{aligned} \quad \text{(PLN-PCA)}$$

$$\theta = \{\mathbf{B}, \mathbf{C}\}$$

Dimension reduction ?  $p(\mathbf{W} \mid \mathbf{Y}, \theta)$ ,  $\mathbf{W} = \{\mathbf{w}_1, \dots, \mathbf{w}_n\}$

- $\mathbf{C}$ : the *loadings* matrix, basis of the latent subspace
- $\mathbf{w}_i$ : the *scores*, coordinates in the subspace

Related to exponential family pPCA (Collins et al. 2001)

A mixture of PLN-PCA for joint  
clustering and dimension  
reduction

---

# Integrating both approaches: mixture of PLN-PCA

Gaussian mixture in the **common** latent  $q$ -dimensional subspace

$$\begin{aligned} \mathbf{z}_i &\sim \mathcal{M}_K(\mathbf{1}, \boldsymbol{\pi}) && \text{(clustering)} \\ \mathbf{w}_i \mid \mathbf{z}_{ik} = 1 &\sim \mathcal{N}_d(\boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k) && \text{(subspace)} \\ \boldsymbol{\eta}_i \mid \mathbf{w}_i &= \mathbf{o}_i + \mathbf{x}_i^\top \mathbf{B} + \mathbf{C}\mathbf{w}_i && \text{(linear transform)} \\ \mathbf{y}_i \mid \boldsymbol{\eta}_i &\sim \otimes_j \mathcal{P}(\exp(\eta_{ij})) && \text{(emission)} \end{aligned} \quad \text{(mPLN-PCA)}$$

With parameters  $\theta = \{\boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k, \mathbf{C}, \boldsymbol{\pi}\}$

The latent layer could be summarized<sup>1</sup>

$$\boldsymbol{\eta}_i \sim \sum_{k=1}^K \pi_k \mathcal{N}_p(\mathbf{o}_i + \mathbf{x}_i^\top \mathbf{B} + \mathbf{C}\boldsymbol{\mu}_k, \mathbf{C}\boldsymbol{\Lambda}_k\mathbf{C}^\top)$$

<sup>1</sup>Analogy with mixture of factors models (Tipping et al. 1999a; McNicholas et al. 2008; McParland et al. 2019)

## Identifiability with common loadings

- (Scale invariance)  $\mathbf{C}^\top \mathbf{C} = \text{Id}_q$
- (Rotational invariance):  $\mathbf{C}$  and  $\mathbf{\Lambda}_k$  only up to a rotation of the latent space

Let  $\mathbf{R} \in \text{Rot}(q)$ , likelihood invariant under  $(\mathbf{C}, \mathbf{\Lambda}_k) \mapsto (\mathbf{C}\mathbf{R}, \mathbf{R}^\top \mathbf{\Lambda}_k \mathbf{R})$   
 $\rightsquigarrow$  can always align the axis with the principal directions of one cluster

**General model** A more general model could be written with one  $\mathbf{C}_k$  per cluster

- one subspace per cluster  $\rightsquigarrow$  no common projection
- $\mathbf{C}_k \mathbf{\Lambda}_k \mathbf{C}_k^\top \rightsquigarrow \mathbf{\Lambda}_k$  should be diagonal

# Inference

---



# Intractable likelihood

Goal estimating  $\theta = \{\pi_k, \boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k, \mathbf{C}, B\}$  via  $\arg \max_{\theta} p_{\theta}(\mathbf{Y})$

**EM algorithm** Standard for latent variable models, use decomposition

$$\log p_{\theta}(\mathbf{Y}) = \mathbb{E}_{p_{\theta}(\mathbf{W}, \mathbf{Z} | \mathbf{Y})} [\log p_{\theta}(\mathbf{Y}, \mathbf{W}, \mathbf{Z})] + \mathcal{H}(p_{\theta}(\mathbf{W}, \mathbf{Z} | \mathbf{Y}))$$

with  $\mathcal{H}(p) = - \int p \log p$  the entropy of  $p$ .

# Intractable likelihood

**Goal** estimating  $\theta = \{\pi_k, \boldsymbol{\mu}_k, \boldsymbol{\Lambda}_k, \mathbf{C}, B\}$  via  $\arg \max_{\theta} p_{\theta}(\mathbf{Y})$

**EM algorithm** Standard for latent variable models, use decomposition

$$\log p_{\theta}(\mathbf{Y}) = \mathbb{E}_{p_{\theta}(\mathbf{W}, \mathbf{Z} | \mathbf{Y})} [\log p_{\theta}(\mathbf{Y}, \mathbf{W}, \mathbf{Z})] + \mathcal{H}(p_{\theta}(\mathbf{W}, \mathbf{Z} | \mathbf{Y}))$$

with  $\mathcal{H}(p) = - \int p \log p$  the entropy of  $p$ .

**Problem(s)** Even for vanilla PLN, intractable

1. likelihood

$$p_{\theta}(\mathbf{y}_i) = \sum_{\mathbf{z}_i} \int_{\mathbb{R}^q} p_{\theta}(\mathbf{y}_i, \mathbf{w}_i, \mathbf{z}_i) d\mathbf{w}_i$$

2. posterior  $p_{\theta}(\mathbf{W}, \mathbf{Z} | \mathbf{Y})$  (or its first moments)

**Solution** use variational inference !

## Variational inference: the Evidence Lower Bound

For any distribution  $q$  on  $(\mathbf{W}, \mathbf{Z})$ , the following inequality holds

$$\log p_{\theta}(\mathbf{Y}) \geq \mathcal{J}(\theta, q) := \mathbb{E}_q [\log p_{\theta}(\mathbf{Y}, \mathbf{W}, \mathbf{Z})] + \mathcal{H}(q)$$

With a quantified gap

$$\log p_{\theta}(\mathbf{Y}) - \mathcal{J}(\theta, q) = \text{KL}(q \parallel p_{\theta}(\mathbf{W}, \mathbf{Z} \mid \mathbf{Y})) \geq 0$$

# Variational inference: the Evidence Lower Bound

For any distribution  $q$  on  $(\mathbf{W}, \mathbf{Z})$ , the following inequality holds

$$\log p_{\theta}(\mathbf{Y}) \geq \mathcal{J}(\theta, q) := \mathbb{E}_q [\log p_{\theta}(\mathbf{Y}, \mathbf{W}, \mathbf{Z})] + \mathcal{H}(q)$$

With a quantified gap

$$\log p_{\theta}(\mathbf{Y}) - \mathcal{J}(\theta, q) = \text{KL}(q \parallel p_{\theta}(\mathbf{W}, \mathbf{Z} \mid \mathbf{Y})) \geq 0$$

**Fix  $\theta$**  Without constraints, minimization leads to  $q = p_{\theta}(\mathbf{W}, \mathbf{Z} \mid \mathbf{Y})$

$\rightsquigarrow$  we constrain  $q = q_{\psi}$  in a parametric class  $\mathcal{Q}$ : the *variational family*

$$\arg \min_{\psi} \text{KL}(q_{\psi} \parallel p_{\theta}(\cdot \mid \mathbf{Y})) = \arg \max_{\psi} \mathcal{J}(\theta, \psi)$$

# Variational inference: the Evidence Lower Bound

For any distribution  $q$  on  $(\mathbf{W}, \mathbf{Z})$ , the following inequality holds

$$\log p_{\theta}(\mathbf{Y}) \geq \mathcal{J}(\theta, q) := \mathbb{E}_q [\log p_{\theta}(\mathbf{Y}, \mathbf{W}, \mathbf{Z})] + \mathcal{H}(q)$$

With a quantified gap

$$\log p_{\theta}(\mathbf{Y}) - \mathcal{J}(\theta, q) = \text{KL}(q \parallel p_{\theta}(\mathbf{W}, \mathbf{Z} \mid \mathbf{Y})) \geq 0$$

Fix  $\theta$  Without constraints, minimization leads to  $q = p_{\theta}(\mathbf{W}, \mathbf{Z} \mid \mathbf{Y})$

$\rightsquigarrow$  we constrain  $q = q_{\psi}$  in a parametric class  $\mathcal{Q}$ : the *variational family*

$$\arg \min_{\psi} \text{KL}(q_{\psi} \parallel p_{\theta}(\cdot \mid \mathbf{Y})) = \arg \max_{\psi} \mathcal{J}(\theta, \psi)$$

**Resulting VEM algorithm** Iteratively solve

$$\text{(VE-step)} \quad \psi^{(t+1)} = \arg \max_{\psi} \mathcal{J}(\theta^{(t)}, \psi)$$

$$\text{(M-step)} \quad \theta^{(t+1)} = \arg \max_{\theta} \mathcal{J}(\theta, \psi^{(t+1)})$$

# The variational family & the ELBO

Mean-field assumption & variational family <sup>2</sup>

$$\mathcal{Q} := \left\{ q_{\psi}(\mathbf{W}, \mathbf{Z}) = \prod_{i=1}^n q_i(\mathbf{w}_i, \mathbf{z}_i) = q_i(\mathbf{w}_i)q_i(\mathbf{z}_i) : \left( \begin{array}{l} \bullet \quad q_i(\mathbf{w}_i) = \mathcal{N}_q(\mathbf{m}_i, \text{diag}(\mathbf{s}_i \odot \mathbf{s}_i)) \\ \bullet \quad q_i(\mathbf{z}_i) = \mathcal{M}_K(\mathbf{1}, \boldsymbol{\tau}_i) \end{array} \right) \right\}$$

---

<sup>2</sup>Alternative choice of variational family, e.g.  $q_{\psi}(\mathbf{w}_i, \mathbf{z}_i) = q_{\psi}(\mathbf{w}_i | \mathbf{z}_i)q_{\psi}(\mathbf{z}_i)$

# The variational family & the ELBO

Mean-field assumption & variational family <sup>2</sup>

$$\mathcal{Q} := \left\{ q_{\psi}(\mathbf{W}, \mathbf{Z}) = \prod_{i=1}^n q_i(\mathbf{w}_i, \mathbf{z}_i) = q_i(\mathbf{w}_i)q_i(\mathbf{z}_i) : \begin{pmatrix} \bullet & q_i(\mathbf{w}_i) = \mathcal{N}_q(\mathbf{m}_i, \text{diag}(\mathbf{s}_i \odot \mathbf{s}_i)) \\ \bullet & q_i(\mathbf{z}_i) = \mathcal{M}_K(\mathbf{1}, \boldsymbol{\tau}_i) \end{pmatrix} \right\}$$

$$\mathcal{J}(\theta, \psi) = \sum_{i=1}^n \mathcal{J}_i(\theta, \psi_i), \quad \psi_i = (\mathbf{m}_i, \mathbf{s}_i, \boldsymbol{\tau}_i) \in \mathbb{R}^q \times \mathbb{R}^q \times \Delta_K \quad (\text{ELBO})$$

Let  $\mathbf{A}_i := \mathbb{E}_{\mathbf{w}_i \sim q_i}[\exp(\boldsymbol{\eta}_i)] = \exp(\mathbf{o}_i + \mathbf{x}_i^\top \mathbf{B} + \mathbf{C}\mathbf{m}_i + \frac{1}{2}\mathbf{C}^2\mathbf{s}_i^2)$

$$\begin{aligned} \mathcal{J}_i(\theta, \psi_i) &= \mathbf{y}_i^\top (\mathbf{o}_i + \mathbf{x}_i^\top \mathbf{B} + \mathbf{C}\mathbf{m}_i) - \mathbf{A}_i^\top \mathbf{1}_p + \log(\mathbf{s}_i^2)^\top \mathbf{1}_q + \text{cst} \\ &\quad - \frac{1}{2} \sum_{k=1}^K \tau_{ik} \left( 2 \log \frac{\tau_{ik}}{\pi_k} - \log |\boldsymbol{\Lambda}_k| + \text{Tr} \left[ \boldsymbol{\Lambda}_k^{-1} (\mathbf{m}_i \mathbf{m}_i^\top + \text{diag}(\mathbf{s}_i^2)) \right] \right) \end{aligned}$$

---

<sup>2</sup>Alternative choice of variational family, e.g.  $q_{\psi}(\mathbf{w}_i, \mathbf{z}_i) = q_{\psi}(\mathbf{w}_i | \mathbf{z}_i)q_{\psi}(\mathbf{z}_i)$

# Properties of the ELBO

The ELBO  $\mathcal{J}(\theta, \psi)$  is

- bi-concave wrt  $\boldsymbol{\tau}$  and  $(\mathbf{M}, \mathbf{S})$
- concave wrt  $(\mathbf{C}, \mathbf{B})$
- bi-concave wrt to  $\boldsymbol{\mu}_k$  and  $\boldsymbol{\Lambda}_k$

but not jointly concave.

Closed-form M-step when  $\psi$  is fixed

- GMM part

$$\hat{\pi}_k = \frac{\tilde{n}_k}{n}, \quad \hat{\boldsymbol{\mu}}_k = \frac{1}{\tilde{n}_k} \boldsymbol{\tau}_{\cdot k}^\top \mathbf{M}, \quad \hat{\boldsymbol{\Lambda}}_k = \frac{1}{\tilde{n}_k} (\mathbf{M} - \mathbf{1}_n \hat{\boldsymbol{\mu}}_k^\top)^\top \text{diag}(\boldsymbol{\tau}_{\cdot k}) (\mathbf{M} - \mathbf{1}_n \hat{\boldsymbol{\mu}}_k^\top) + \text{diag}(\boldsymbol{\tau}_{\cdot k}^\top \mathbf{S})$$

- Covariables:  $\hat{\mathbf{B}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X} \mathbf{M} \mathbf{C}^\top$

But not for  $\mathbf{C}$  or VE-step



## Inference: two possible strategies

① Standard variational EM: alternate between  $\max_{\psi} \mathcal{J}$  &  $\max_{\theta} \mathcal{J}$ .

② Joint optimization of  $J$  w.r.t.  $(\theta, \psi)$

- ▶ Scalability with  $n$  and  $p$  (work of **Bastien Batardière** on PLN-PCA)

**WIP** : currently working on torch implementation (**R** & **Python**)

- automatic differentiation framework to compute  $\nabla_{\theta}$  &  $\nabla_{\psi}$
- stochastic optimization (e.g. ADAM)
- Amortized VI:  $q_{\psi} = q(\mathbf{w}_i, \mathbf{z}_i \mid g_{\psi}(\mathbf{y}_i))$ ,  $g_{\psi}$  neural net with weights  $\psi$

Choice of  $(K, q)$  is a model selection problem

Clustering context: integrated classification likelihood (ICL, Biernacki et al. 2000)

BIC-like: (very) temporarily adopt a Bayesian POV on parameters  $\theta$

$$\log p(\mathbf{Y}, \mathbf{Z}) = \int_{\theta} \int_{\mathbf{W}} p(\mathbf{Y}, \mathbf{Z}, \mathbf{W}, \theta)$$

# Model selection

Choice of  $(K, q)$  is a model selection problem

Clustering context: integrated classification likelihood (ICL, Biernacki et al. 2000)

BIC-like: (very) temporarily adopt a Bayesian POV on parameters  $\theta$

$$\log p(\mathbf{Y}, \mathbf{Z}) = \int_{\theta} \int_{\mathbf{W}} p(\mathbf{Y}, \mathbf{Z}, \mathbf{W}, \theta)$$

Laplace + Stirling approximation + ELBO proxy leads to an approximate ICL

$$vICL(K, q) = \mathcal{J}(\hat{\theta}, \hat{\psi}) - \frac{1}{2} \left( pq - \frac{q(q+1)}{2} + K - 1 + Kq + K \frac{q(q+1)}{2} \right) \log(n) \quad (1)$$

$\mathcal{J}(\hat{\theta}, \hat{\psi})$  serves as a proxy for  $\log p(\mathbf{Y}, \hat{\mathbf{Z}} | \hat{\theta})$

## Conclusion

---

## Objectives

① a partition (*clustering*)

$$\mathbf{Z} = \{z_1, \dots, z_n\}$$

# Summary of the methodology

## Objectives

① a partition (*clustering*)

$$\mathbf{Z} = \{z_1, \dots, z_n\}$$

② low-dimensional representation

$$\mathbf{W} = \{w_1, \dots, w_n\}$$

# Summary of the methodology

## Objectives

① a partition (*clustering*)

$$\mathbf{Z} = \{z_1, \dots, z_n\}$$

② low-dimensional representation

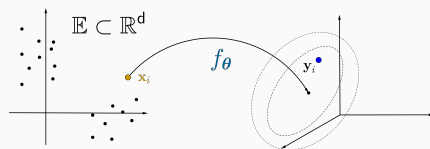
$$\mathbf{W} = \{w_1, \dots, w_n\}$$

## Method

$$\mathbf{Z}, \mathbf{W} \sim p_\theta \quad (\textit{latent})$$

$$\eta = f_\theta(\mathbf{W}) \quad (\textit{param})$$

$$\mathbf{Y} \mid \eta \sim p(\cdot \mid \eta) \quad (\textit{obs})$$



# Summary of the methodology

## Objectives

① a partition (*clustering*)

$$\mathbf{Z} = \{z_1, \dots, z_n\}$$

② low-dimensional representation

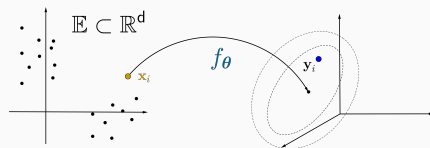
$$\mathbf{W} = \{w_1, \dots, w_n\}$$

## Method

$$\mathbf{Z}, \mathbf{W} \sim p_\theta \quad (\textit{latent})$$

$$\boldsymbol{\eta} = f_\theta(\mathbf{W}) \quad (\textit{param})$$

$$\mathbf{Y} \mid \boldsymbol{\eta} \sim p(\cdot \mid \boldsymbol{\eta}) \quad (\textit{obs})$$



Inference to estimate  $\hat{\theta}$  + (variational) posterior for  $\hat{\mathbf{Z}}, \hat{\mathbf{W}} \approx \arg \max_{\mathbf{W}, \mathbf{Z}} p_{\hat{\theta}}(\mathbf{W}, \mathbf{Z} \mid \mathbf{Y})$



## Remaining things to do

- ▶ Finish the inference algorithm + integration in the **PLNmodels** package
- ▶ Application on single-cell RNAseq data
- ▶ Numerical investigation of the property of M-estimator (work of J. Chiquet et. al. on PLN using (Westling et al. 2015))

$$\hat{\theta}_n = \arg \max_{\theta} \left\{ \mathcal{J}_n(\theta) = \frac{1}{n} \sum_{i=1}^n \mathcal{J}_i(\theta, \hat{\psi}_i(\theta, \mathbf{y}_i)) \right\}, \quad \hat{\theta}_n \xrightarrow[n \rightarrow +\infty]{?} \arg \max_{\theta} \mathbb{E}_{\theta^*} \left[ \mathcal{J}(\theta, \hat{\psi}(\theta, \mathbf{Y})) \right]$$

- ▶ Zero-inflation  $\rightsquigarrow$  introduce a binary *mask* variable  $H_{ij} \sim \mathcal{B}(\text{logit}(\mathbf{x}_i^\top \mathbf{B}_j^0))$

## Non-linear extensions: variational auto-encoders

- ▶  $\boldsymbol{\eta}_i = f_{\theta}(\mathbf{w}_i)$ ,  $f_{\theta}$  neural net with weights  $\theta$  (encoder)
- ▶  $q_{\psi} = q(\mathbf{w}_i, \mathbf{z}_i \mid g_{\psi}(\mathbf{y}_i))$ ,  $g_{\psi}$  neural net with weights  $\psi$  (decoder)






# Thank you for your attention<sup>3</sup>






---

<sup>3</sup>And sorry for the lack of experiments today, though I'd like to share the blame with the French Government...

# References

---

-  Aitchison, J. and C.H. Ho (1989). “The multivariate Poisson-log normal distribution”. In: *Biometrika* 76.4, pp. 643–653.
-  Biernacki, Christophe, Gilles Celeux, and Gérard Govaert (2000). “Assessing a mixture model for clustering with the integrated completed likelihood”. In: *IEEE transactions on pattern analysis and machine intelligence* 22.7, pp. 719–725.
-  Chiquet, Julien, Mahendra Mariadassou, and Stéphane Robin (2018). “Variational inference for probabilistic Poisson PCA”. In: *The Annals of Applied Statistics* 12.4, pp. 2674–2698.
-  Chiquet, Julien, Mahendra Mariadassou, and Stéphane Robin (2021). “The Poisson-lognormal model as a versatile framework for the joint analysis of species abundances”. In: *Frontiers in Ecology and Evolution* 9, p. 588292.
-  Collins, Michael, Sanjoy Dasgupta, and Robert E Schapire (2001). “A generalization of principal components analysis to the exponential family”. In: *Advances in neural information processing systems* 14.

-  McNicholas, Paul David and Thomas Brendan Murphy (2008). “Parsimonious Gaussian mixture models”. In: *Statistics and Computing* 18.3, pp. 285–296.
-  McParland, Damien and Thomas Brendan Murphy (2019). “Mixture modelling of high-dimensional data”. In: *Handbook of Mixture Analysis*. Chapman and Hall/CRC, pp. 239–270.
-  Tipping, Michael E and Christopher M Bishop (1999a). “Mixtures of probabilistic principal component analyzers”. In: *Neural computation* 11.2, pp. 443–482.
-  Tipping, Michael E and Christopher M Bishop (1999b). “Probabilistic principal component analysis”. In: *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 61.3, pp. 611–622.
-  Westling, Ted and Tyler H. McCormick (2015). *Beyond prediction: A framework for inference with variational approximations in mixture models*.

## Questions

---

