# Probabilistic Count Matrix Factorization for Single Cell Expression Data Analysis

Ghislain DURIF

July 4th 2019

Single-cell RNA-seq day, Toulouse

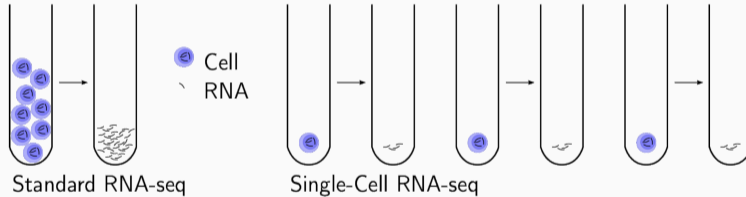CNRS, IMAG, Université de Montpellier, France
ghislain.durif@umontpellier.fr
https://gdurif.perso.math.cnrs.fr/

## Outline

# Introduction

## RNA-seq

- Quantification of gene expression on a genomic scale



## Single-cell level (scRNA-seq)

- gene-to-gene variability: expression dynamics (low expression genes)
- cell-to-cell variability: diversity within a population of cells
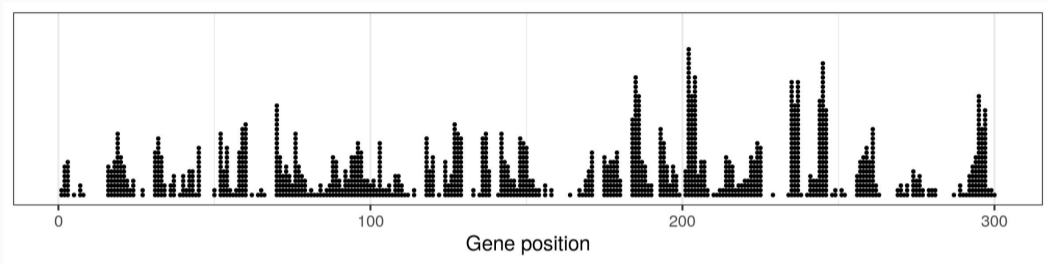
$x_{ij}$ = expression of gene $j$ in sample $i$

$$\mathsf{X}_{n \times p} = \left[ \begin{array}{c} \\ \\ \quad x_{ij} \quad \\ \\ \end{array} \right] \left. \begin{array}{c} 1 \\ \vdots \\ n \end{array} \right\} \text{cells}$$

$$\underbrace{1 \quad \ldots \quad \ldots \quad \ldots \quad \ldots \quad p}_{\text{genes}}$$

· **High dimension:**
  $n$ grows but $\ll p$

$\rightarrow n \sim 100/1000$
and $p \sim 10000$

$x_{ij}$ = height at position $j$

- **Count data** with drop-out events in single-cell RNA-seq (zero-inflation)

$\rightarrow$ number of reads that map to a gene position

## Issues with high dimensional data ($p \gg n$)

*"The curse of high-dimensionality"* (Donoho, 2000)

- **Geometry:** counter-intuitive behavior of metrics (Aggarwal et al., 2001)
  - $\rightarrow$ Representation: how to visualize thousands of variables?

- **Optimization:** numerical singularities due to complex dependencies (colinearity)

- **Computational efficiency and scalability**

- **Geometry:** counter-intuitive behavior of metrics (Aggarwal et al., 2001)

$p \mapsto \|\mathbf{x}_1 - \mathbf{x}_2\|_2$

with $\mathbf{x}_1, \mathbf{x}_2 \in \mathbb{R}^p$

## Issues with high dimensional data ($p \gg n$)

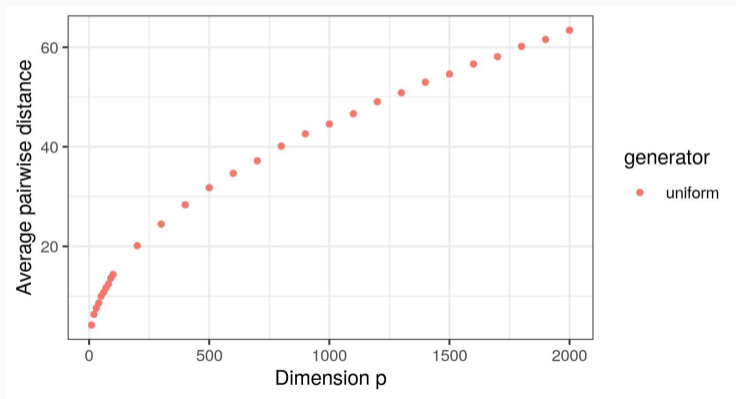*"The curse of high-dimensionality"* (Donoho, 2000)

- **Geometry:** counter-intuitive behavior of metrics (Aggarwal et al., 2001)
  - $\rightarrow$ Representation: how to visualize thousands of variables?

- **Optimization:** numerical singularities due to complex dependencies (colinearity)

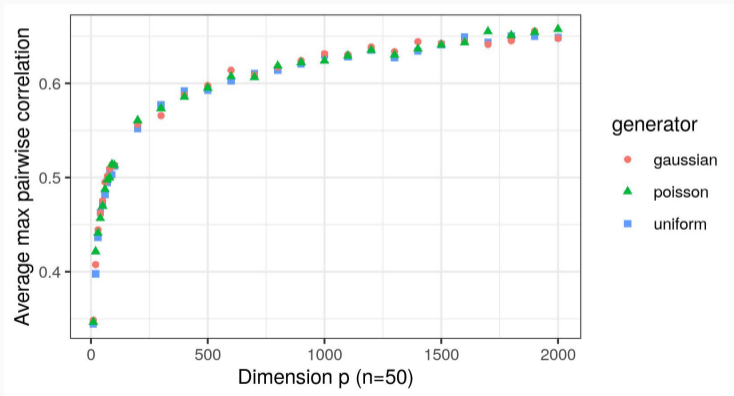- **Computational efficiency and scalability**

- **Optimization:** numerical singularities due to complex dependencies (colinearity)

$$\mathsf{X}_{n \times p} = \left[ \; \mathsf{x}_1 \; \middle| \; \ldots \; \middle| \; \mathsf{x}_p \; \right]$$

$$p \mapsto \max_{j \neq \ell \in \{1, \ldots, p\}} \left| \mathrm{Corr}(\mathsf{x}_j, \mathsf{x}_\ell) \right|$$



6

- **Optimization:** numerical singularities due to complex dependencies (colinearity)

$p \mapsto \text{rank}(X^T X)$

**Note:** if $\text{rank}(X^T X) < p$
then $X^T X$ is not invertible



6

## Issues with high dimensional data ($p \gg n$)

*"The curse of high-dimensionality"* (Donoho, 2000)

- **Geometry:** counter-intuitive behavior of metrics (Aggarwal et al., 2001)
    - $\rightarrow$ Representation: how to visualize thousands of variables?

- **Optimization:** numerical singularities due to complex dependencies (colinearity)

- **Computational efficiency and scalability**

*"The curse of high-dimensionality"* (Donoho, 2000)

- **Geometry:** counter-intuitive behavior of metrics (Aggarwal et al., 2001)
  - $\rightarrow$ Representation: how to visualize thousands of variables?

- **Optimization:** numerical singularities due to complex dependencies (colinearity)

- **Computational efficiency and scalability**

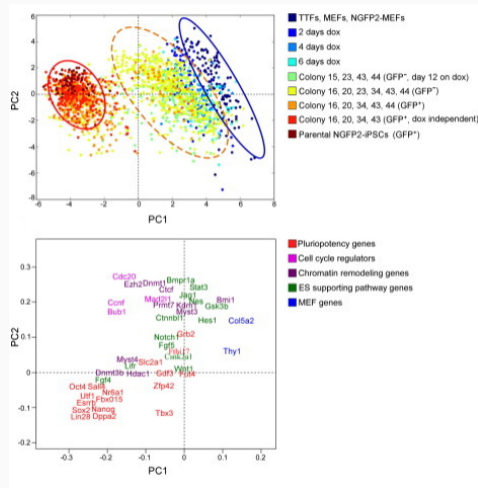$\rightarrow$ Dimension reduction approaches

**Statistical challenges**

- Data exploration

- Visualization / clustering



Buganim et al. (2012)

- Representation of the data in a
  lower dimensional subspace

- Consider sparsity
  - → select the variables (genes) that
    contribute to this representation



Buganim et al. (2012)

# Zero-inflation and drop-out events

- No gene expression

- Transcription is **bursty** (cells are not synchronized)

- Failure of the sequencing (**dropout events** = loss of the information)



Gene expression distribution (Freytag et al., 2018, "goldstandard" dataset)

# Dimension reduction with matrix factorization

## Data dimension

$$\mathsf{X}_{n \times p} = \begin{bmatrix} & & & & & & & \\ \hline & & & x_{ij} & & & & \\ \hline & & & & & & & \end{bmatrix} \begin{matrix} 1 \\ \vdots \\ n \end{matrix} \Bigg\} \text{observations/cells}$$

$$\underbrace{\begin{matrix} 1 & \dots & \dots & \dots & \dots & p \end{matrix}}_{\text{variables/genes}}$$

· $n$ individuals (cells) with $p$ recordings

· $p$ variables (genes) with $n$ observations

# Data dimension

$$X_{n \times p} = \begin{bmatrix} & & & & & \\ \rule{0pt}{0pt} & & & x_{ij} & & \\ & & & & & \end{bmatrix} \begin{matrix} 1 \\ \vdots \\ n \end{matrix} \Bigg\} \text{ observations/cells}$$

$$\underbrace{\quad 1 \quad \ldots \quad \ldots \quad \ldots \quad \ldots \quad p \quad}_{\text{variables/genes}}$$

- $n$ individuals (cells) with $p$ recordings

- $p$ variables (genes) with $n$ observations

- $n$ individuals (cells) with $p$ recordings

- $p$ variables (genes) with $n$ observations

Regression

**2-d**



Source: `wikipedia.org`

**3-d**



Source: `stackoverflow.com`

**Dimension > 3**
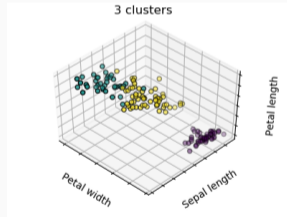
?

## Clustering

### 2-d



Source: `stackoverflow.com`

### 3-d



Source: `scikit-learn.org`

### Dimension > 3

?

## Low dimensional representation

$$
\mathbf{U}_{n \times K} = \begin{pmatrix} & | & & | & \\ \rule[0.5ex]{1em}{0.4pt} & \rule[0.5ex]{1em}{0.4pt} & u_{ik} & \rule[0.5ex]{1em}{0.4pt} & \rule[0.5ex]{1em}{0.4pt} \\ & | & & | & \end{pmatrix} \begin{matrix} 1 \\ \vdots \\ n \end{matrix} \qquad \text{and} \qquad \mathbf{V}_{p \times K} = \begin{matrix} 1 \\ \vdots \\ \vdots \\ p \end{matrix} \begin{pmatrix} & | & & | & \\ \rule[0.5ex]{1em}{0.4pt} & \rule[0.5ex]{1em}{0.4pt} & v_{jk} & \rule[0.5ex]{1em}{0.4pt} & \rule[0.5ex]{1em}{0.4pt} \\ & | & & | & \end{pmatrix}
$$

$$
\mathbf{u}_1 \ \ldots \ \mathbf{u}_K \qquad\qquad\qquad\qquad\qquad\qquad \mathbf{v}_1 \ \ldots \ \mathbf{v}_K
$$

### Visualization/clustering

- Represention of individuals/cells (columns of $\mathbf{U}$) in dimension $K < p$
- Contribution of variables/genes (columns of $\mathbf{V}$) in dimension $K < n$

$$U_{n \times K} = \begin{pmatrix} & & \\ & u_{ik} & \\ & & \end{pmatrix} \begin{matrix} 1 \\ \vdots \\ n \end{matrix} \quad \text{and} \quad V_{p \times K} = \begin{matrix} 1 \\ \vdots \\ \vdots \\ p \end{matrix} \begin{pmatrix} & & \\ & v_{jk} & \\ & & \end{pmatrix}$$

$$u_1 \ldots u_K \qquad\qquad\qquad v_1 \ldots v_K$$

## Visualization/clustering

- Represension of individuals/cells (columns of U) in dimension $K < p$
- Contribution of variables/genes (columns of V) in dimension $K < n$

$$\mathsf{U}_{n \times K} = \begin{pmatrix} & | & | & \\ & u_{ik} & & \\ & | & | & \end{pmatrix} \begin{matrix} 1 \\ \vdots \\ n \end{matrix} \qquad \text{and} \qquad \mathsf{V}_{p \times K} = \begin{matrix} 1 \\ \vdots \\ \vdots \\ p \end{matrix} \begin{pmatrix} | & | \\ v_{jk} & \\ | & | \end{pmatrix}$$

$$\mathsf{u}_1 \ \ldots \ \mathsf{u}_K \qquad\qquad\qquad \mathsf{v}_1 \ \ldots \ \mathsf{v}_K$$
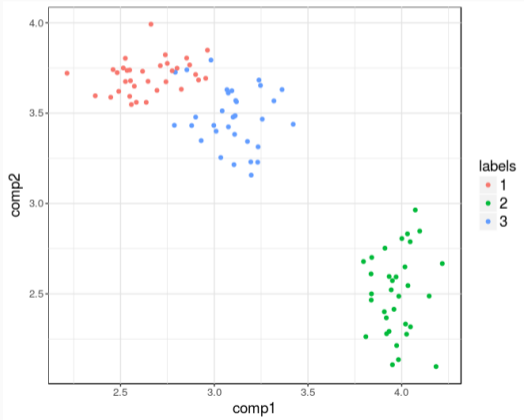
## Visualization/clustering

- Represention of individuals/cells (columns of $\mathsf{U}$) in dimension $K < p$
- Contribution of variables/genes (columns of $\mathsf{V}$) in dimension $K < n$

**Data visualization** with $K = 2$

· scatter plot of $(u_{i1}, u_{i2})_{i=1:n}$



· Expose hidden/latent structures

· Principal Component Analysis (PCA)?

11

- Linear **projection** of X onto a lower dimensional space (of dim. $K$)

$$u_{ik} = \sum_{j=1}^{p} x_{ij} \, v_{jk} \qquad\qquad \mathbf{u}_k = \mathbf{X} \, \mathbf{v}_k \ \text{ with } \begin{cases} \mathbf{u}_k \in \mathbb{R}^n \\ \mathbf{v}_k \in \mathbb{R}^p \end{cases}$$

$$\mathbf{v}_k = \underset{\mathbf{v} \in \mathbb{R}^p}{\mathrm{argmax}}\ \mathrm{Var}(\mathbf{Xv})$$

- Additional constraints: $\|v\|_2 = 1$ and $\mathbf{t}_k$ orthogonal to $\mathbf{u}_1, \ldots, \mathbf{u}_{k-1}$

$$X_{n \times p} \times V_{p \times K} = U_{n \times K}$$

- **Enforce sparsity:** only a few variables contribute to the model

- **Objective:** drop non relevant variables from the model

$$\underset{\mathbf{v} \in \mathbb{R}^p}{\operatorname{argmax}} \operatorname{Var}(\mathbf{Xv}) + \lambda \sum_{j=1}^{p} |v_j|$$

- Lasso principle Tibshirani (1996): $\ell_1$ penalty on the weights $\mathbf{v}$

# Why matrix factorization?

→ generalization of PCA principle

$$x_{ij} = \sum_{k=1}^{K} u_{ik}\, v_{jk}$$

$$X_{n \times p} \approx U_{n \times K} \quad UV^T_{n \times p}$$

$$V^T_{K \times p}$$

$$x_{ij} = \sum_{k=1}^{K} u_{ik} v_{jk}$$

- Individuals (cells): $U \in \mathbb{R}^{n \times K}$
- Variables (genes): $V \in \mathbb{R}^{p \times K}$



$$x_{ij} = \sum_{k=1}^{K} u_{ik} \, v_{jk}$$

$X_{n \times p}$ $\approx$ $U_{n \times K}$ $UV^T_{n \times p}$

$V^T_{K \times p}$

- $K$ = latent dimension (hopefully small)
- $UV^T$ = low-rank representation of $X$

Sense of the approximation?

$X_{n \times p}$ $\approx$ $U_{n \times K}$ $UV^T_{n \times p}$ $V^T_{K \times p}$

17

Least Square Approximation?

$$\operatorname*{argmin}_{\substack{U \in \mathbb{R}^{n \times K} \\ V \in \mathbb{R}^{p \times K}}} \left\| X - UV^T \right\|_F^2$$

Least Square Approximation?

$$\underset{\substack{U \in \mathbb{R}^{n \times K} \\ V \in \mathbb{R}^{p \times K}}}{\mathrm{argmin}} \left\| X - UV^T \right\|_F^2$$

- Solution given by Singular Value Decomposition (SVD Eckart and Young, 1936)

- PCA = SVD of $\widetilde{X}$ where $\widetilde{x}_{ij} = x_{ij} - \bar{x}_j$

Sparse SVD (Shen and Huang, 2008; Witten et al., 2009)

$$\underset{\substack{\mathbf{u}\in\mathbb{R}^n \\ \mathbf{v}\in\mathbb{R}^p}}{\mathrm{argmin}} \left\{ \left\| \mathbf{X} - \mathbf{u}\mathbf{v}^T \right\|_F^2 + \lambda \sum_{j=1}^{p} |v_j| \right\}$$

- $\ell_1$ penalty shrinks contributions of non pertinent variables to zero

Relation between geometry and underlying model

$\| \cdot \|_2 \leftrightarrow$ Gaussian distribution

Gaussian SVD?

- $X \sim \mathcal{N}\left( UV^T, \ \mathbf{\Sigma}^2 \right)$     i.e.     $X_{ij} \sim \mathcal{N}\left( \sum_k u_{ik} v_{jk}, \ \sigma_{ij}^2 \right)$

- $\log \mathcal{L}(U, V) = \left\| X - UV^T \right\|_F^2$

## Relation between geometry and underlying model

$$\| \cdot \|_2 \leftrightarrow \text{Gaussian distribution}$$

- Count = not Gaussian

- First idea: $X_{ij} \sim \mathscr{P}(\lambda)$

- Highly expressed genes

  $\rightarrow$ large $\lambda$

  $\rightarrow$ Gaussian approximation



Empirical distribution, counts drawn from $\mathscr{P}(200)$

19

## Relation between geometry and underlying model

$$\|\cdot\|_2 \leftrightarrow \text{Gaussian distribution}$$

- Lowly expressed genes in single-cell RNA-seq?

- Poisson assumption?

  $\rightarrow$ Why not Negative Binomial?

  $\rightarrow$ $\mathcal{NB}$ suitable for RNA-seq (Anders and Huber, 2010) and for scRNA-seq (Chen et al., 2016)



Empirical distribution, counts drawn from $\mathscr{P}(2)$

## Relation between geometry and underlying model

$\| \cdot \|_2 \leftrightarrow$ Gaussian distribution

- Lowly expressed genes in single-cell RNA-seq?

- Poisson assumption?

  $\rightarrow$ Why not Negative Binomial?

  $\rightarrow$ $\mathcal{NB}$ suitable for RNA-seq (Anders and Huber, 2010) and for scRNA-seq (Chen et al., 2016)



Empirical distribution, counts drawn from $\mathcal{NB}(n = 5, p = 2.5E - 3)$

19

1) Interest for **lowly expressed genes** in single-cell

2) **Over-dispersion** in RNA-seq data
   $\rightarrow \mathsf{Var}(X_{ij}) > \mathbb{E}[X_{ij}]$

3) Single-cell data: **zero-inflation**
   $\rightarrow \mathbb{P}(X_{ij} = 0) > e^{-\lambda}$

## Zero-inflated over-dispersed counts

| Dataset | $n$ | $p$[1] | prop. 0 |
|---|---|---|---|
| Baron et al. (2016) | 1886 | 6080 | 80.9% |
| Freytag et al. (2018) goldstandard | 925 | 8580 | 39.5% |
| Freytag et al. (2018) silverstandard 5 | 8352 | 4547 | 86.3% |
| Llorens-Bobadilla et al. (2015) | 141 | 13826 | 64.8% |

---

[1]after pre-filtering

# PCA on zero-inflated count data



High intensity Poisson data

Same data with zero-inflation

Observations scores over first two principal components

# Matrix Factorization for count data

Least Square Approximation with non-negativity constraints

$$\underset{\substack{\mathbf{U}\in\mathbb{R}^{n\times K}\\\mathbf{V}\in\mathbb{R}^{p\times K}}}{\text{argmin}} \left\|\mathbf{X} - \mathbf{U}\mathbf{V}^T\right\|_F^2 \qquad \text{where} \qquad u_{ik}, v_{jk} \geq 0 \text{ for any } i, j, k$$

Embed PCA with a **probabilistic model**

- Replace $\| \cdot \|_2$ approximation by likelihood-based approaches

- $X_{ij} \sim$ probability distribution in the exponential family

$\rightarrow$ Factorization of $\mathbb{E}[\mathsf{X}]$ rather than $\mathsf{X}$

$$\mathbb{E}[X_{ij}] = \sum_{k=1}^{K} u_{ik} \, v_{jk}$$

$$X_{n \times p} \sim \mathscr{P}(\Lambda) \qquad U_{n \times K} \qquad \Lambda \approx UV^{T}_{n \times p}$$

- $X_{ij} \sim \mathscr{P}(\lambda_{ij})$ with the Poisson rate matrix $\Lambda = [\lambda_{ij}]_{n \times p}$

- Factorization: $\mathbb{E}[X] = \Lambda \approx UV^{T} \quad \leftrightarrow \quad \lambda_{ij} \approx \sum_{k} u_{ik} v_{jk}$

- Maximum Likelihood Estimation under non-negativity constraint over **U** and **V**



$$v_{jk}$$

$$\mathscr{P}\left(\sum_k u_{ik}\, v_{jk}\right)$$

$$X_{ij}$$

$$u_{ik}$$

- **U** and **V** are parameters

- Optimization computationally expensive

- No account for over-dispersion or zero-inflation

**Bregman divergence** between X and $UV^T$

$$D\left(X \mid UV^T\right)$$

- Based on the parametrization in the exponential family

$$\text{Poisson:} \quad D\left(x_{ij} \mid \lambda_{ij}\right) = x_{ij} \log \frac{x_{ij}}{\lambda_{ij}} - x_{ij} + \lambda_{ij}$$

$\rightarrow$ Connected to the likelihood

$\rightarrow$ Choice of the geometry driven by the model

# Matrix factorization for over-dispersed zero-inflated count data?

- Probabilistic matrix factorization

- **Hierarchical model**: prior on factors U and V

- Model inference: likelihood optimization $\rightarrow$ **variational inference**

- Impose sparsity on V: how to select variables?
  $\rightarrow$ **Probabilistic selection**

- Negative Binomial = standard distribution for over-dispersed count

- $X_{ij} \sim \mathcal{NB}(r_{ij}, \pi_{ij}) \rightarrow$ complex optimization of the likelihood

- Negative Binomial = standard distribution for over-dispersed count

- $X_{ij} \sim \mathcal{NB}(r_{ij}, \pi_{ij}) \rightarrow$ complex optimization of the likelihood

- Gamma-Poisson hierarchical model

$$\lambda_{ij} \sim \Gamma(\alpha_1, \alpha_2)$$
$$X_{ij} \mid \lambda \sim \mathscr{P}(\lambda_{ij})$$

$\rightarrow$ Marginal distribution of $X_{ij}$ is a Negative Binomial distribution:

$$X_{ij} \sim \mathcal{NB}\left(\alpha_1, \frac{\alpha_2}{\alpha_2 + 1}\right)$$

- Independent Gamma prior distributions over factors $\mathsf{U}$ and $\mathsf{V}$:

$$U_{ik} \sim \Gamma(\alpha_{k,1}, \alpha_{k,2}) \;\; \text{and} \;\; V_{jk} \sim \Gamma(\beta_{k,1}, \beta_{k,2})$$

- Conditional Poisson distribution over the data $\mathsf{X}$:

$$X_{ij} \,|\, (U_{ik}, V_{jk})_{k=1:K} \;\sim\; \mathscr{P}\left(\textstyle\sum_k U_{ik} V_{jk}\right)$$

- Factors = latent variables

- Recover the posterior
  $\widehat{U} = \mathbb{E}[U \mid X]$ and $\widehat{V} = \mathbb{E}[V \mid X]$

- Marginal distribution is over-dispersed
  $\text{Var}(X_{ij}) > \mathbb{E}[X_{ij}]$

1. "Zero-inflated" Gamma-Poisson factor model

- **Poisson-Dirac mixture:** $\quad X_{ij} \,|\, (U_{ik}, V_{jk})_{k=1:K} \;\sim\; (1 - \pi_j^{\mathrm{D}}) \times \delta_0 + \pi_j^{\mathrm{D}} \times \mathscr{P}(\lambda_{ij})$

- $1 - \pi_j^{\mathrm{D}} \in [0, 1]$ is the dropout rate for gene $j$



- $D_{ij}$ = drop-out event indicator

- $\mathbb{P}(X_{ij} = 0 \,|\, \mathbf{U}, \mathbf{V}) > e^{-\lambda_{ij}}$

- Variable $j$ contributes to factor $k$ if $V_{jk} \neq 0$

- Objective: force the $V_{jk}$'s to be null for non pertinent genes

If **V** was a parameter $\rightarrow \ell_1$ penalty

$$\operatorname*{argmin}_{\substack{\mathbf{u} \in \mathbb{R}^n \\ \mathbf{v} \in \mathbb{R}^p}} \left\{ \left\| \mathbf{X} - \mathbf{u}\mathbf{v}^T \right\|_F^2 + \lambda \sum_{j=1}^{p} |v_j| \right\}$$

· $V_{jk}$ is a random variable $\rightarrow$ necessary to use sparsity-inducing priors

## Spike and slab model:

· Continuous one-group prior: shrinkage to small value near zero (ex: Bayesian Lasso with Laplace prior)



· **Two-group prior:** mixture between a Dirac and a continuous distribution, true mass at zero

$\rightarrow$ to induce a "sparse" posterior with a mass at zero

## 2. Sparse Gamma-Poisson model



- Probabilistic variable selection

- Gamma-Dirac mixture
  $V_{jk} \sim (1 - \pi_j^s)\, \delta_0 + \pi_j^s\, \Gamma(\beta_{k,1}, \beta_{k,2})$

- $\pi_j^s \in [0, 1]$ probability that gene $j$ contributes to the model

- $S_{jk}$ = sparsity indicator

- **Objective:** estimation of the factors $\mathsf{U}$ and $\mathsf{V}$

- $\mathscr{L}(\mathsf{U}\,|\,\mathsf{X})$ and $\mathscr{L}(\mathsf{V}\,|\,\mathsf{X})$ are not explicit

  - $\rightarrow$ **Cannot use** Maximum a Posteriori (MAP) or Expectation-Maximization (EM)
  - $\rightarrow$ Inference of the posterior of latent variables

- Markov Chain Monte Carlo (MCMC) are computationally expensive

- Variational inference[2]: approximation of the posterior

[2]See supplementary slides at the end for more details

- Objective: estimation of the factors $U$ and $V$

- $\mathscr{L}(U\,|\,X)$ and $\mathscr{L}(V\,|\,X)$ are not explicit

    - $\rightarrow$ **Cannot use** Maximum a Posteriori (MAP) or Expectation-Maximization (EM)

    - $\rightarrow$ Inference of the posterior of latent variables

- Markov Chain Monte Carlo (MCMC) are computationally expensive

- Variational inference[2]: approximation of the posterior

[2]See supplementary slides at the end for more details

# Choice of $K$?

- Concerns all low rank methods

- No consensus procedure

- Explained variance
    $\rightarrow \ell_2$ metric criterion

- Bregman divergence:

$$k \mapsto D\Big(\mathsf{X} \mid \widehat{\mathsf{U}}_{1:k}(\widehat{\mathsf{V}}_{1:k})^T\Big)$$

- Visualization: $K = 2$

# Quality of the model

## Percentage of explained deviance

$$\%\text{dev} = \frac{\log p(\mathsf{X} \mid \mathbf{\Lambda} = \widehat{\mathsf{U}}\widehat{\mathsf{V}}^T) - \log p(\mathsf{X} \mid \mathbf{\Lambda} = \bar{\mathsf{X}})}{\log p(\mathsf{X} \mid \mathbf{\Lambda} = \mathsf{X}) - \log p(\mathsf{X} \mid \mathbf{\Lambda} = \bar{\mathsf{X}})}$$

- $\log p(\mathsf{X} \mid \mathbf{\Lambda})$: conditional distribution in the model
- $\mathbf{\Lambda} = \mathsf{X}$: saturated model
- $\mathbf{\Lambda} = \bar{\mathsf{X}}$: moment estimator (column-wise average)

## Percentage of explained deviance

$$\%\text{dev} = \frac{\log p(\mathsf{X} \mid \mathbf{\Lambda} = \widehat{\mathsf{U}}\widehat{\mathsf{V}}^T) - \log p(\mathsf{X} \mid \mathbf{\Lambda} = \bar{\mathsf{X}})}{\log p(\mathsf{X} \mid \mathbf{\Lambda} = \mathsf{X}) - \log p(\mathsf{X} \mid \mathbf{\Lambda} = \bar{\mathsf{X}})}$$

- $\log p(\mathsf{X} \mid \mathbf{\Lambda})$: conditional distribution in the model
- $\mathbf{\Lambda} = \mathsf{X}$: saturated model
- $\mathbf{\Lambda} = \bar{\mathsf{X}}$: moment estimator (column-wise average)

## Example in the Gaussian case

$$\%\text{dev} = \frac{\sum_{k=1}^{K} \sigma_k^2}{\sum_{\ell=1}^{\text{rk}(\mathsf{X})} \sigma_\ell^2} = \% \text{ explained variance from PCA}$$

$\sigma_1 > \ldots > \sigma_{\text{rk}(\mathsf{X})}$ = singular values of $\mathsf{X}$

t-SNE[3] (van der Maaten and Hinton, 2008)

- No measure of the quality of the representation

- How to choose the "perplexity" parameter?

---

[3]C.f. later

# Experiments

## t-SNE (van der Maaten and Hinton, 2008)

t-Stochastic Neighbourhood Embedding

### Dimension $p$

- Observations: $x_1, \ldots, x_n \in \mathbb{R}^p$
- Probabilistic distribution $\mathcal{P}$ on pairwise distance $d(x_i, x_{i'})$

### Dimension 2

- Low-dimensional obervations: $u_1, \ldots, u_n \in \mathbb{R}^2$
- Probabilistic distribution $\mathcal{Q}$ on pairwise distance $d(u_i, u_{i'})$

Dimension $p$

Dimension 2

distrib. $\mathcal{P}$

distrib. $\mathcal{Q}$

$$\underset{\mathbf{U}\in\mathbb{R}^{n\times 2}}{\operatorname{argmin}} \operatorname{KL}\Big(\mathcal{P} \mid \mathcal{Q}\Big)$$

Dimension $p$

Dimension 2

distrib. $\mathcal{P}$

distrib. $\mathcal{Q}$

$$\operatorname*{argmin}_{\mathsf{U} \in \mathbb{R}^{n \times 2}} \mathsf{KL}\Big(\mathcal{P} \mid \mathcal{Q}\Big)$$

Perplexity parameter?

- "You see what you want"

- What happens if you don't know what you are looking for?

- See `https://distill.pub/2016/misread-tsne/`

Gaussian PCA with zero-inflated compartment

$$X_{ij} \sim (1 - \pi_j)\delta_0 + \pi_j \mathcal{N}\left(\sum_k U_{ik} v_{jk}, \sigma_{ij}^2\right)$$

Latent factors:

- $U_{ik} \sim \mathcal{N}(\cdot, \cdot)$
- $v_{jk} =$ parameter

Gaussian PCA with zero-inflated compartment

$$X_{ij} \sim (1 - \pi_j)\delta_0 + \pi_j \mathcal{N}\left(\sum_k U_{ik} v_{jk}, \sigma_{ij}^2\right)$$

Latent factors:

- $U_{ik} \sim \mathcal{N}(\cdot, \cdot)$
- $v_{jk} =$ parameter

u$_1$ vs u$_2$ (individual representation)

**v**$_1$ vs **v**$_2$ (variable representation)

Clustering of individuals          Quality of the model

Clustering of variables

Quality of the model

$u_1$ vs $u_2$ (individual representation)



Freytag et al. (2018) goldstandard dataset

$v_1$ vs $v_2$ (variable representation)



Freytag et al. (2018) goldstandard dataset

# scRNA-seq: quantitative results

| | prop. 0 | ngroup | | (s)pCMF | PCA | ZIFA | t-SNE |
|---|---|---|---|---|---|---|---|
| Baron et al. (2016) | 80.9% | 13 | adj. RI | 21.2% | 14.3% | 15.4% | 14.2% |
| | | | %dev | 73.2% | 41.6% | 53.5% | / |
| Freytag et al. (2018) goldstandard | 39.5% | 3 | adj. RI | 81.3% | 60.1% | 56.8% | 60.5% |
| | | | %dev | 55.7% | 65.6% | 48.6% | / |
| Freytag et al. (2018) silverstandard 5 | 86.3% | 11 | adj. RI | 24.2% | 16.2% | 19.8% | 24.8% |
| | | | %dev | 70.0% | 55.1% | / | / |
| Llorens-Bobadilla et al. (2015) | 64.8% | 6 | adj. RI | 40.1% | 25.3% | 38.3% | 29.8% |
| | | | %dev | 64.4% | 34.8% | 42.6% | / |

## To conclude

- Optimization algorithm $\rightarrow$ variational inference

- Algorithm initialization
  - $\rightarrow$ especially sparse compartment (gene pre-filtering strategy)

- Optimization algorithm → variational inference

- Algorithm initialization
  - → especially sparse compartment (gene pre-filtering strategy)

- Data visualization
    - PCA
        - → hidden PCA?
        - → representation with low of explained variance ($< 10\%$)?

    - t-SNE: clustering vs visualization a posteriori

- Dimension reduction (unsupervised)

- probabilistic Count Matrix Factorization

- Data visualization
    - PCA
        - → hidden PCA?
        - → representation with low of explained variance ($< 10\%$)?

    - t-SNE: clustering vs visualization a posteriori

- Dimension reduction (unsupervised)

- probabilistic Count Matrix Factorization

- Data visualization
    - PCA
        - $\rightarrow$ hidden PCA?
        - $\rightarrow$ representation with low of explained variance ($< 10\%$)?

    - t-SNE: clustering vs visualization a posteriori

- Dimension reduction (unsupervised)

- probabilistic Count Matrix Factorization

- Data visualization

- Dimension reduction (unsupervised)
    - → Latent space projection
    - → Variable selection (sparsity)

- probabilistic Count Matrix Factorization

## Take-home message

- Data visualization

- Dimension reduction (unsupervised)
  - $\rightarrow$ Latent space projection
  - $\rightarrow$ Variable selection (sparsity)

- probabilistic Count Matrix Factorization
  - $\rightarrow$ Model-based
  - $\rightarrow$ Data-driven (count, over-dispersion, zero-inflation)

Count matrix Factorization

- Model selection criterion (choice of *K*)

- Calibration of the sparsity hyper-parameter

- Stochastic procedure to improve the optimization

- Extension to account for covariates in the model

# Thanks for your attention

Collaborators

- Franck Picard (CNRS, LBBE)
- Sophie Lambert-Lacroix (Université Grenoble Alpes, TIMC)
- Laurent Modolo (CNRS, LBMC)
- Jeff Mold (Karolinska Institutet, Stockholm)

Institutions

- LBBE, Lyon 1 University
- Inria Grenoble, Thoth team

# Thanks for your attention

Paper            https://doi.org/10.1093/bioinformatics/btz177
                 https://arxiv.org/abs/1710.11028

Software (R)     https://github.com/gdurif/pCMF
                 https://github.com/gdurif/pCMF_experiments

# References

Aggarwal, C. C., A. Hinneburg, and D. A. Keim (2001). On the surprising behavior of distance metrics in high dimensional space. In *International Conference on Database Theory*, pp. 420–434. Springer.

Anders, S. and W. Huber (2010). Differential expression analysis for sequence count data. *Genome Biology 11*(10), R106.

Baron, M., A. Veres, S. L. Wolock, A. L. Faust, R. Gaujoux, A. Vetere, J. H. Ryu, B. K. Wagner, S. S. Shen-Orr, A. M. Klein, D. A. Melton, and I. Yanai (2016, October). A Single-Cell Transcriptomic Map of the Human and Mouse Pancreas Reveals Inter- and Intra-cell Population Structure. *Cell systems 3*(4), 346–360.e4.

Beal, M. J. and Z. Ghahramani (2003). The variational Bayesian EM algorithm for incomplete data: With application to scoring graphical model structures. *Bayesian statistics 7*, 453–464.

Buganim, Y., D. A. Faddah, A. W. Cheng, E. Itskovich, S. Markoulaki, K. Ganz, S. L. Klemm, A. van Oudenaarden, and R. Jaenisch (2012, September). Single-Cell

Expression Analyses during Cellular Reprogramming Reveal an Early Stochastic and a Late Hierarchic Phase. *Cell 150*(6), 1209–1222.

Cemgil, A. T. (2009, January). Bayesian Inference for Nonnegative Matrix Factorisation Models. *Intell. Neuroscience 2009*, 4:1–4:17.

Chen, H.-I. H., Y. Jin, Y. Huang, and Y. Chen (2016, August). Detection of high variability in gene expression from single-cell RNA-seq profiling. *BMC Genomics 17*(Suppl 7).

Collins, M., S. Dasgupta, and R. E. Schapire (2001). A generalization of principal components analysis to the exponential family. In *Advances in Neural Information Processing Systems*, pp. 617–624.

Donoho, D. L. (2000). High-dimensional data analysis: The curses and blessings of dimensionality. *AMS Math Challenges Lecture*, 1–32.

Durif, G., L. Modolo, J. E. Mold, S. Lambert-Lacroix, and F. Picard (2019).

Probabilistic Count Matrix Factorization for Single Cell Expression Data Analysis. *Bioinformatics In press*.

Eckart, C. and G. Young (1936). The approximation of one matrix by another of lower rank. *Psychometrika 1*(3), 211–218.

Freytag, S., L. Tian, I. Lönnstedt, M. Ng, and M. Bahlo (2018, August). Comparison of clustering tools in R for medium-sized 10x Genomics single-cell RNA-sequencing data. *F1000Research 7*.

Hoffman, M. D., D. M. Blei, C. Wang, and J. Paisley (2013, May). Stochastic Variational Inference. *J. Mach. Learn. Res. 14*(1), 1303–1347.

Hotelling, H. (1933). Analysis of a complex of statistical variables into principal components. *Journal of educational psychology 24*(6), 417.

Lee, D. D. and H. S. Seung (1999, October). Learning the parts of objects by non-negative matrix factorization. *Nature 401*(6755), 788–791.

Llorens-Bobadilla, E., S. Zhao, A. Baser, G. Saiz-Castro, K. Zwadlo, and A. Martin-Villalba (2015, September). Single-Cell Transcriptomics Reveals a Population of Dormant Neural Stem Cells that Become Activated upon Brain Injury. *Cell Stem Cell 17*(3), 329–340.

Pearson, K. (1901, November). LIII. On lines and planes of closest fit to systems of points in space. *Philosophical Magazine Series 6 2*(11), 559–572.

Pierson, E. and C. Yau (2015). ZIFA: Dimensionality reduction for zero-inflated single-cell gene expression analysis. *Genome Biology 16*, 241.

Shen, H. and J. Z. Huang (2008). Sparse principal component analysis via regularized low rank matrix approximation. *Journal of multivariate analysis 99*(6), 1015–1034.

Tibshirani, R. (1996). Regression Shrinkage and Selection via the Lasso. *Journal of the Royal Statistical Society. Series B (Methodological) 58*(1), 267–288.

van der Maaten, L. and G. Hinton (2008). Visualizing Data using t-SNE. *Journal of Machine Learning Research 9*(Nov), 2579–2605.

Witten, D. M., R. Tibshirani, and T. Hastie (2009, July). A penalized matrix decomposition, with applications to sparse principal components and canonical correlation analysis. *Biostatistics 10*(3), 515–534.

Zou, H., T. Hastie, and R. Tibshirani (2006). Sparse principal component analysis. *Journal of computational and graphical statistics 15*(2), 265–286.

## Variational inference (see Hoffman et al., 2013)

- $p(\mathsf{U}, \mathsf{V} \,|\, \mathsf{X})$ approximated by the **variational distribution** $q(\mathsf{U}, \mathsf{V})$

- Regarding the **Kullback-Leibler divergence**
  $\rightarrow$ quantify the "proximity" between two distributions of probability

$$q(\mathsf{U}, \mathsf{V}) \;=\; \underset{\text{distribution } \widetilde{q}}{\mathrm{argmin}} \; \mathsf{KL}\Big( \widetilde{q}(\mathsf{U}, \mathsf{V}) \,\big|\, p(\mathsf{U}, \mathsf{V} \,|\, \mathsf{X}) \Big)$$

- **Constraints on $q$:**
  $\rightarrow$ $q$ is factorizable: independence between the factors

$$q(\mathsf{U}, \mathsf{V}) = \prod_{i,k} q(u_{ik} \,;\, \mathsf{a}_{ik}) \times \prod_{j,k} q(v_{jk} \,;\, \mathsf{b}_{jk})$$

  $\rightarrow$ $q$ respects the Gamma-Poisson conjugacy in the exponential family

## The Evidence Lower Bound (ELBO)

**Objective:** $\qquad J(q) = \mathbb{E}_q[\log p(\mathsf{X}, \mathsf{U}, \mathsf{V})] - \mathbb{E}_q[\log q(\mathsf{U}, \mathsf{V})]$

- A **lower bound** on the marginal log-likelihood: $\log p(\mathsf{X}) \geq J(q)$
  (by Jensen's inequality)

- **Maximizing** $J(q)$ equivalent to **minimizing** $\mathrm{KL}\Big( q(\mathsf{U}, \mathsf{V}) \,|\, p(\mathsf{U}, \mathsf{V} \,|\, \mathsf{X}) \Big)$
  $\rightarrow$ because $J(q) = \log p(\mathsf{X}) - \mathrm{KL}\Big( q(\mathsf{U}, \mathsf{V}) \,|\, p(\mathsf{U}, \mathsf{V} \,|\, \mathsf{X}) \Big)$

- $J(q)$ is optimized regarding the variational parameters $\mathsf{a}_{ik}$ and $\mathsf{b}_{jk}$

## Optimization of the ELBO

- Gradient of $J(q)$ regarding the variational parameters:

$$\left.\begin{array}{c} \nabla_{a_{ik}} \\ \nabla_{b_{jk}} \\ \nabla_{(r_{ijk})_k} \end{array}\right| J(q)$$

- Expression of the ELBO regarding the variational parameters:

$$\widetilde{J}(a_{ik}) = \mathbb{E}_q[\log p(u_{ik} \mid -)] - \mathbb{E}_q[\log q(u_{ik}\,;\, a_{ik})] + \text{cst}$$

$$\widetilde{J}(b_{jk}) = \mathbb{E}_q[\log p(v_{jk} \mid -)] - \mathbb{E}_q[\log q(v_{jk}\,;\, b_{jk})] + \text{cst}$$

$$\widetilde{J}((r_{ijk})_k) = \mathbb{E}_q\left[\log p((z_{ijk})_k \mid -)\right] - \mathbb{E}_q\left[\log q((z_{ijk})_k\,;\, (r_{ijk})_k)\right] + \text{cst}$$

$\rightarrow$ Explicit coordinates of the point that sets the gradient to zero

$\rightarrow$ **Iterative optimization through a fixed-point algorithm**

## Variational EM algorithm (Beal and Ghahramani, 2003)

1) Variational E-step:

   $\rightarrow$ Estimation of the variational parameters $\mathbf{a}$ and $\mathbf{b}$

2) M-step:

   $\rightarrow$ Estimation of the prior parameters (Gamma parameters $\boldsymbol{\alpha}$ and $\boldsymbol{\beta}$)

   $$\text{EM: } \underset{\boldsymbol{\alpha}, \boldsymbol{\beta}}{\operatorname{argmax}} \ \mathbb{E}_{\mathsf{U}, \mathsf{V} \mid \mathsf{X}}[\log \mathcal{L}(\mathsf{X}, \mathsf{U}, \mathsf{V} \, ; \, \boldsymbol{\alpha}, \boldsymbol{\beta})]$$

Output: estimation of the factors by the variational expectation

$$\widehat{\mathsf{U}} = \mathbb{E}_q[\mathsf{U}] \ \text{ and } \ \widehat{\mathsf{V}} = \mathbb{E}_q[\mathsf{V}]$$

## Variational EM algorithm (Beal and Ghahramani, 2003)

1) Variational E-step:

$\rightarrow$ Estimation of the variational parameters $a$ and $b$

2) M-step:

$\rightarrow$ Estimation of the prior parameters (Gamma parameters $\alpha$ and $\beta$)

$$\text{vEM:} \ \underset{\alpha,\beta}{\text{argmax}} \ \mathbb{E}_q[\log \mathcal{L}(X, U, V \, ; \, \alpha, \beta)]$$

Output: estimation of the factors by the variational expectation

$$\widehat{U} = \mathbb{E}_q[U] \ \text{ and } \ \widehat{V} = \mathbb{E}_q[V]$$