

Dependent Dirichlet-Multinomial Processes with Random Number of Components

Andrea Cremaschi¹ Beatrice Franzolini²

SISBAYES 2025

Department of Statistical Sciences, University of Padova 4-5 September 2025

¹School of Science and Technology, IE University, Madrid

²Bocconi Institute for Data Science and Analytics, Bocconi University, Milan



Several settings present data with a grouped structure



- Several settings present data with a grouped structure
- ✗ Ignore grouping. Neglects heterogeneity across groups



- Several settings present data with a grouped structure
- Ignore grouping. Neglects heterogeneity across groups
- Assume independence across groups. Does not allow for sharing of information



- Several settings present data with a grouped structure
- Ignore grouping. Neglects heterogeneity across groups
- Assume independence across groups. Does not allow for sharing of information
- Allow dependence within and between groups



In Bayesian Statistics (and Nonparametrics) several modelling approaches exists

of different ways to induce dependence among groups

In Bayesian Statistics (and Nonparametrics) several modelling approaches exists of different ways to induce dependence among groups

FOCUS on of this talk: mixture models with random number of components

They are less studied than the infinite-dimensional approaches, especially in the context of grouped data

In Bayesian Statistics (and Nonparametrics) several modelling approaches exists

of different ways to induce dependence among groups

FOCUS on of this talk: mixture models with random number of components

They are less studied than the infinite-dimensional approaches, especially in the context of grouped data

Most recent and related work:

Colombi et al. (2024) "Hierarchical Mixture of Finite Mixtures". Bayesian Analysis.

We exploit the well-known Dirichlet-Multinomial construction.

(2)

We exploit the well-known Dirichlet-Multinomial construction.

Let G groups, each equipped with a mixture of M components.

For each group g = 1, ..., G:

$$P_g(\cdot) \stackrel{a.s.}{=} \sum_{m=1}^{M} W_{gm} \delta_{\theta_m}(\cdot)$$
 (1)

$$\theta_1,\ldots,\theta_M\mid M\stackrel{\text{i.i.d.}}{\sim} P_0$$

$$W_a = (W_{a1}, \dots, W_{aM}) \mid M \sim \text{Dir}(\alpha, \dots, \alpha)$$
 (3)

where:

- (1) are the mixing measures
- (2) are the atoms (shared across groups)
- (3) are the mixture weights with Dirichlet marginals
 - within each group ${\color{red} \bowtie}$ Dirichlet-Multinomial process $P_g \sim \mathsf{DMP}(\alpha, P_0)$

IDEA \square Can we directly model the matrix of weights $\mathbf{W} = \{\mathbf{W}_g, g = 1, \dots, G\}$ to induce dependence across groups?

<u>IDEA</u> \square Can we directly model the matrix of weights $\mathbf{W} = \{\mathbf{W}_g, g = 1, \dots, G\}$ to induce dependence across groups?

Properties we require:

•
$$W_{gm} \in (0,1)$$
, for $g = 1, ..., G$, $m = 1, ..., M$

•
$$\sum_{m=1}^{M} W_{gm} = 1$$
, for $g = 1, ..., G$

•
$$\textit{W}_g = (\textit{W}_{g1}, \dots, \textit{W}_{gM}) \sim \text{Dir}(\alpha, \dots, \alpha)$$

Let $P_a, P_\ell \sim \mathsf{DMP}(\alpha, P_0)$; $X \mid P_a \sim P_a$ and $Y \mid P_\ell \sim P_\ell$.

For any compatible joint distribution of the weights the following holds:

Proposition (Lower bound correlations)

If $Corr(W_{gm}, W_{\ell m}) \geq 0$ for any m, then:

- (i) $Corr[P_g(A), P_\ell(A)] \ge \frac{M\alpha + 1}{M\alpha + M}$ (ii) $Corr(X, Y) \ge \frac{1}{M}$
- (iii) The lower bounds obtained with independent sequences of weights $\mathbf{W}_a \perp \mathbf{W}_\ell$

Let $P_a, P_\ell \sim \mathsf{DMP}(\alpha, P_0)$; $X \mid P_a \sim P_a$ and $Y \mid P_\ell \sim P_\ell$.

For any compatible joint distribution of the weights the following holds:

Proposition (Lower bound correlations)

If $Corr(W_{gm}, W_{\ell m}) \geq 0$ for any m, then:

- (i) $Corr[P_g(A), P_\ell(A)] \ge \frac{M\alpha + 1}{M\alpha + M}$ (ii) $Corr(X, Y) \ge \frac{1}{M}$
- (iii) The lower bounds obtained with independent sequences of weights $\mathbf{W}_a \perp \mathbf{W}_\ell$

Colombi et al. (2024) with Dirichlet weights of lowest possible correlation

Let $P_a, P_\ell \sim \mathsf{DMP}(\alpha, P_0)$; $X \mid P_a \sim P_a$ and $Y \mid P_\ell \sim P_\ell$.

For any compatible joint distribution of the weights the following holds:

Proposition (Lower bound correlations)

If $Corr(W_{gm}, W_{\ell m}) \geq 0$ for any m, then:

- (i) $Corr[P_g(A), P_\ell(A)] \ge \frac{M\alpha + 1}{M\alpha + M}$ (ii) $Corr(X, Y) \ge \frac{1}{M}$
- (iii) The lower bounds obtained with independent sequences of weights $W_a \perp W_\ell$

Colombi et al. (2024) with Dirichlet weights of lowest possible correlation

Different approach: joint distribution of the weights via matrix-variate random variables

Let $\boldsymbol{U} = (U_1, \dots, U_M)$ be an *M*-dimensional array of $G \times G$ matrices

 ${\it U}$ has a (symmetric) $\underline{\it matrix-variate Dirichlet distribution}$ with parameter α , and we write ${\it U} \mid \alpha \sim {\sf MDir}_{\it M.G}({\it U} \mid \alpha)$, if it has a density:

$$p(\boldsymbol{U} \mid \alpha) = \frac{1}{\beta_{G}^{\alpha}} \prod_{m=1}^{M} |U_{m}|^{\alpha - (G+1)/2} \mathbb{1}_{\boldsymbol{U} \in \mathcal{S}_{M,G}}$$

where:

- α > G/2
- β_G^{α} is the multivariate beta function (Gupta & Nagar, 2018)
- support:

$$\overline{\mathcal{S}_{M,G}} = \left\{ (U_1, \dots, U_M) : U_m \text{ is positive definite for each } m \text{ and } U_M = \mathbb{I}_G - \sum\limits_{m=1}^{M-1} U_m \right\}$$

- Olkin and Rubin (1964)

 study of the distributional and independence properties
- Extensively studied in Gupta and Nagar (2018)
- Admits constructive definition via normalisation of Wishart matrices

- Olkin and Rubin (1964)

 study of the distributional and independence properties
- Extensively studied in Gupta and Nagar (2018)
- Admits constructive definition via normalisation of Wishart matrices

Why is it of interest to us?

- Olkin and Rubin (1964)

 study of the distributional and independence properties
- Extensively studied in Gupta and Nagar (2018)
- Admits constructive definition via normalisation of Wishart matrices

Why is it of interest to us?

Let $W_{gm} \stackrel{d}{=} [U_m]_{gg}$, then:

- $✓ W_{gm}$ ∈ (0,1), for g = 1, ..., G and m = 1, ..., M
- $\bigvee_{m=1}^{M} W_{gm} = 1$, for g = 1, ..., G
- \mathbf{V} $\mathbf{W}_g = (\mathbf{W}_{g1}, \dots, \mathbf{W}_{gM}) \sim \text{Dir}(\alpha, \dots, \alpha)$

• Let $\boldsymbol{U} \mid \alpha \sim \mathsf{MDir}_{M,G}(\boldsymbol{U} \mid \alpha)$

• Let $\boldsymbol{U} \mid \alpha \sim \mathsf{MDir}_{M,G}(\boldsymbol{U} \mid \alpha)$

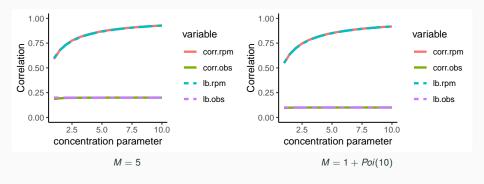
 The diagonal slice of this array has the required properties

Groups (g)

- Let $\boldsymbol{U} \mid \alpha \sim \mathsf{MDir}_{M,G}(\boldsymbol{U} \mid \alpha)$
- The diagonal slice of this array has the required properties
- We set $W_{gm} = [U_m]_{gg}$ for $g = 1, \ldots, G$ and $m = 1, \ldots, M$

What is the dependence among groups when using ${\bf \textit{U}} \mid \alpha \sim {\rm MDir}_{M,G}({\bf \textit{U}} \mid \alpha)$?

What is the dependence among groups when using $\boldsymbol{U} \mid \alpha \sim \text{MDir}_{M,G}(\boldsymbol{U} \mid \alpha)$? Numerical example with G = 2, $\alpha > 1$



ISSUES

- 1. matrix-variate Dirichlet construction does not yield a wide range of correlations
- 2. $\alpha > G/2$ limits span of correlations and flexibility of resulting distributions

ISSUES

- 1. matrix-variate Dirichlet construction does not yield a wide range of correlations
- 2. $\alpha > G/2$ limits span of correlations and flexibility of resulting distributions

SOLUTIONS

- 1. leverage the unnormalised construction of the matrix-variate Dirichlet distribution
- 2. start from *Generalised-Wishart distribution* to allow for $\alpha \leq G/2$ (Srivastava, 2003)

New construction:

Let $S_m \overset{\mathrm{i.i.d.}}{\sim} \mathsf{Gen\text{-}Wishart}_G\left(2\alpha,\Psi\right)$, then for $m=1,\ldots,M$

$$S_g = \sum_{m=1}^M \left[S_m \right]_{gg}$$

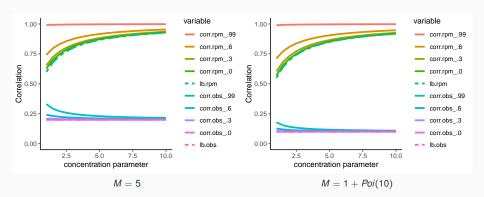
$$\textit{W}_{\textit{gm}} = \left[\textit{S}_{\textit{m}}\right]_{\textit{gg}}/\textit{S}_{\textit{g}}$$

and we write $(P_1, \dots, P_G) \sim W\text{-DMP}(\alpha, P_0, \Psi)$

- no need for matrix inversion or Cholesky decomposition!

☆ Correlation under new construction

Numerical example with
$$G=2$$
, $\alpha>1$, $\Psi=\begin{pmatrix} \psi_{11} & \psi_{12} \\ \psi_{12} & \psi_{22} \end{pmatrix}$



- (i) Closure under marginalisation: $P_g \sim \text{DMP}(\alpha, P_0)$ $\Rightarrow \Psi$ does not affect the marginal distribution of the processes
- Ψ does not affect the marginal distribution of the processe
- (ii) Independence

 $\textit{\textbf{W}}_g \perp \textit{\textbf{W}}_\ell$ (Colombi et al. 2024, Dirichlet weights) recovered when $\Psi = \mathbb{I}_G$

(iii) Full-dependence

$$\psi_{gl}
ightarrow 1$$
 or $\psi_{gl}
ightarrow -1$ $extbf{ extit{ extbf{ iny W}}_g} \stackrel{a.s.}{=} extbf{ extbf{ iny W}_\ell}, P_g \stackrel{a.s.}{=} P_\ell$

Ψ act as correlation matrix btw rpms

 Ψ can be chosen a priori to encode the dependence among rpms

(i) Distribution of the matrix of weights $[W]_{gm}$, for G=2, conditionally on M>0

Let
$$\Psi = \begin{pmatrix} \psi_{11} & \psi_{12} \\ \psi_{12} & \psi_{22} \end{pmatrix}$$
 and $\rho^2 = \frac{\psi_{12}^2}{\psi_{11}\psi_{22}} \in [0,1)$, then:

$$\rho(\mathbf{W}) = \frac{(1 - \rho^2)^{M\alpha}}{\Gamma(\alpha)^M} \sum_{(k_1, \dots, k_M)} \rho^{2K} \Gamma(K + M\alpha)^2 \prod_{m=1}^M \frac{(W_{1m} W_{2m})^{k_m + \alpha - 1}}{k_m! \Gamma(k_m + \alpha)} \mathbb{1}_{\mathbf{W} \in \Delta_M^2}$$

where $k_m \in \mathbb{N} \cup 0$, $K = \sum_m k_m$, Δ_M^2 is an appropriate probability simplex.

(i) Distribution of the matrix of weights $[W]_{qm}$, for G=2, conditionally on M>0

Let
$$\Psi = \begin{pmatrix} \psi_{11} & \psi_{12} \\ \psi_{12} & \psi_{22} \end{pmatrix}$$
 and $\rho^2 = \frac{\psi_{12}^2}{\psi_{11}\psi_{22}} \in [0,1)$, then:

$$p(\mathbf{W}) = \frac{(1 - \rho^2)^{M\alpha}}{\Gamma(\alpha)^M} \sum_{(k_1, \dots, k_M)} \rho^{2K} \Gamma(K + M\alpha)^2 \prod_{m=1}^M \frac{(W_{1m} W_{2m})^{k_m + \alpha - 1}}{k_m! \Gamma(k_m + \alpha)} \mathbb{1}_{\mathbf{W} \in \Delta_M^2}$$

where $k_m \in \mathbb{N} \cup 0$, $K = \sum_m k_m$, Δ_M^2 is an appropriate probability simplex.

(ii) Hierarchical representation of $p(\mathbf{W})$

$$extbf{\emph{W}}_{g\cdot} = (extbf{\emph{W}}_{g1}, \dots, extbf{\emph{W}}_{gM}) \mid extbf{\emph{k}} = (extbf{\emph{k}}_1, \dots, extbf{\emph{k}}_{M}) \overset{iid}{\sim} \operatorname{Dirichlet}(extbf{\emph{k}}_1 + lpha, \dots, extbf{\emph{k}}_{M} + lpha), \quad g = 1, 2$$

$$extbf{\emph{k}} = (extbf{\emph{k}}_1, \dots, extbf{\emph{k}}_{M}) \overset{iid}{\sim} \operatorname{NegBin}(lpha, 1 -
ho^2),$$

Let $X_{gi} \mid (P_1, \dots, P_G) \stackrel{ind}{\sim} P_g$ and $(P_1, \dots, P_G) \sim \text{W-DMP}(\alpha, P_0, \Psi)$. Denote with $(X_m^\star, m = 1, \dots, M^{(obs)})$ the unique values in the data, and with $\mathbf{n}_g = (n_{g1}, \dots, n_{gM})$ the corresponding numerosities, then a posteriori (conditionally on M > 0)

$$\begin{split} P_g(\cdot) &\stackrel{a.s.}{=} \sum_{m=1}^{M^{(obs)}} W_{gm} \delta_{X_m^{\star}}(\cdot) + \sum_{m=M^{(obs)}+1}^{M} W_{gm} \delta_{\theta_m}(\cdot), \quad g=1,2 \\ W_g. &= (W_{g1}, \ldots, W_{gM}) \mid \textbf{\textit{k}} \stackrel{\textit{iid}}{\sim} \text{Dirichlet} (\textbf{\textit{k}} + \textbf{\textit{n}}_{g.} + \alpha) \qquad \text{for } g=1,2 \\ \rho(\textbf{\textit{k}}) &\propto \frac{\prod\limits_{g=1}^{2} \mathsf{B}(\textbf{\textit{k}} + \textbf{\textit{n}}_{g.} + \alpha)}{\mathsf{B}(\textbf{\textit{k}} + \alpha)^2} \prod_{m=1}^{M} \left[\binom{k_m + \alpha - 1}{k_m} \rho^{2k_m} \right] \mathbb{1}_{k_m \in \mathbb{N} \cup 0} \end{split}$$

Let $X_{gi} \mid (P_1, \dots, P_G) \stackrel{ind}{\sim} P_g$ and $(P_1, \dots, P_G) \sim \text{W-DMP}(\alpha, P_0, \Psi)$. Denote with $(X_m^\star, m = 1, \dots, M^{(obs)})$ the unique values in the data, and with $\mathbf{n}_g = (n_{g1}, \dots, n_{gM})$ the corresponding numerosities, then a posteriori (conditionally on M > 0)

$$\begin{split} P_g(\cdot) &\stackrel{a.s.}{=} \sum_{m=1}^{M^{(obs)}} W_{gm} \delta_{X_m^*}(\cdot) + \sum_{m=M^{(obs)}+1}^{M} W_{gm} \delta_{\theta_m}(\cdot), \quad g=1,2 \\ W_{g\cdot} &= (W_{g1}, \dots, W_{gM}) \mid \mathbf{k} \stackrel{iid}{\sim} \text{Dirichlet} (\mathbf{k} + \mathbf{n}_{g\cdot} + \alpha) \qquad \text{for } g=1,2 \\ p(\mathbf{k}) &\propto \frac{\prod\limits_{g=1}^{2} \mathsf{B}(\mathbf{k} + \mathbf{n}_{g\cdot} + \alpha)}{\mathsf{B}(\mathbf{k} + \alpha)^2} \prod_{m=1}^{M} \left[\binom{k_m + \alpha - 1}{k_m} \rho^{2k_m} \right] \mathbb{1}_{k_m \in \mathbb{N} \cup 0} \end{split}$$

Extensions to G > 2 available (more complicated)

ie

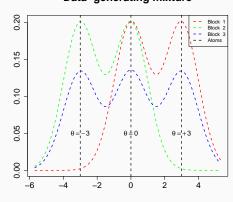
Simulation Study

Simulation setting:

- G = 30 split into 3 blocks of size 10
- Within each block, the true data-generating mixture has different weights but same atoms

$$W_g^{true} = \begin{cases} \left(0, \frac{1}{2}, \frac{1}{2}\right), & g = 1, \dots, 10 & \frac{6}{8} \\ \left(\frac{1}{2}, \frac{1}{2}, 0\right), & g = 11, \dots, 20 \\ \left(\frac{1}{3}, \frac{1}{3}, \frac{1}{3}\right), & g = 21, \dots, 30 & \frac{8}{8} \end{cases}$$

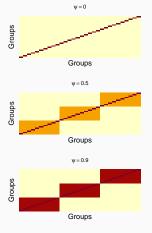
Data-generating mixture





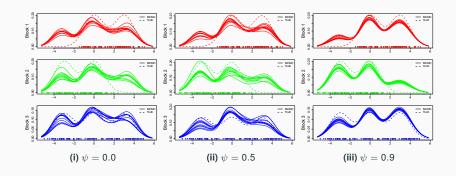
We fit a univariate Normal-Inverse Gamma model

For each simulation setting, the matrix $\Psi=[\psi]_{g_1g_2}$ has different off-diagonal elements for $g_1,g_2=1,\ldots,G$

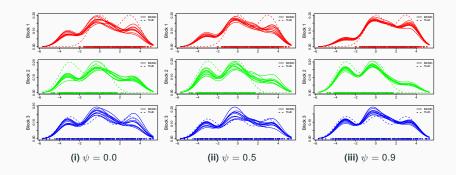


- (i) $\psi = 0.0$
- (ii) $\psi = 0.5$
- (iii) $\psi = 0.9$

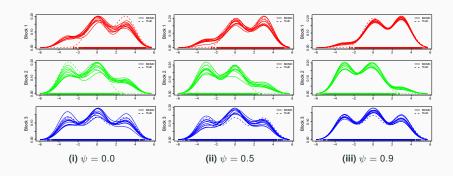
$$n_q = 10$$
, for $g = 1, ..., G$



$$n_q = 25$$
, for $g = 1, ..., G$

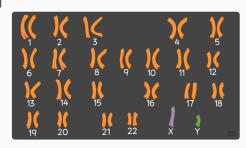


$$n_q = 100$$
, for $g = 1, ..., G$

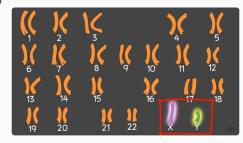


Sex differences in gene expressions from the human brain

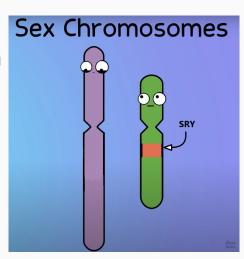
 Humans carry 23 pairs of chromosomes as their genetic material (karyotype)



- Humans carry 23 pairs of chromosomes as their genetic material (karyotype)
- 22 pairs are called <u>autosome</u>
 1 special pair contains the XY chromosomes

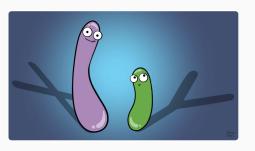


- Humans carry 23 pairs of chromosomes as their genetic material (karyotype)
- 22 pairs are called <u>autosome</u>
 1 special pair contains the XY chromosomes



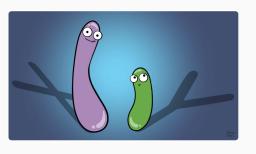
- Humans carry 23 pairs of chromosomes as their genetic material (karyotype)
- 22 pairs are called <u>autosome</u>
 1 special pair contains the <u>XY chromosomes</u>
- Karyotype changes across species





 Some traits are sex-specific and the corresponding genes carried by X or Y are called X-linked and Y-linked genes



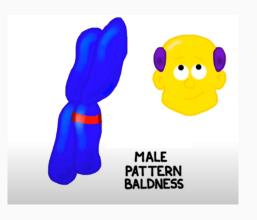


- Some traits are sex-specific and the corresponding genes carried by X or Y are called X-linked and Y-linked genes
- · Some examples:

X- and Y-linked genes

Hemophilia Normal blood vessel Tear Platelets and clotting factors Blood Hemophilia Uncontrolled bledding Cleveland

- Some traits are sex-specific and the corresponding genes carried by X or Y are called X-linked and Y-linked genes
- Some examples:
 - haemophilia r X-linked



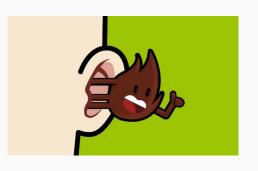
- Some traits are sex-specific and the corresponding genes carried by X or Y are called X-linked and Y-linked genes
- · Some examples:
 - haemophilia x X-linked
 - baldness
 X-linked





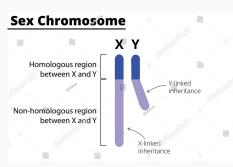
- Some traits are sex-specific and the corresponding genes carried by X or Y are called X-linked and Y-linked genes
- · Some examples:
 - haemophilia x X-linked
 - baldness
 X-linked
 - webbed toes Y-linked





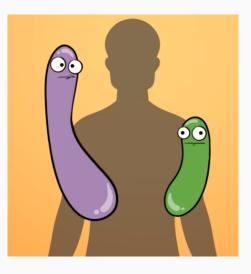
- Some traits are sex-specific and the corresponding genes carried by X or Y are called X-linked and Y-linked genes
- · Some examples:
 - haemophilia r X-linked
 - baldness
 X-linked
 - webbed toes
 Y-linked
 - hairy ears
 thought to be Y-linked but not yet confirmed (Lee et al. 2004)

X- and Y-linked genes



- Some traits are sex-specific and the corresponding genes carried by X or Y are called X-linked and Y-linked genes
- Some examples:
 - haemophilia
 X-linked
 - baldness
 X-linked
 - webbed toes Y-linked
 - hairy ears
 thought to be Y-linked but not yet confirmed (Lee et al. 2004)
- X-Y homologues are genes that are present in the X and Y chromosomes.
 They are very similar (but not identical)...





- Some traits are sex-specific and the corresponding genes carried by X or Y are called X-linked and Y-linked genes
- · Some examples:
 - haemophilia S X-linked
 - baldness 🖙 X-linked
 - webbed toes Y-linked
 - hairy ears thought to be Y-linked but not yet confirmed (Lee et al. 2004)
- X-Y homologues are genes that are present in the X and Y chromosomes.
 They are very similar (but not identical)...
- · ...also found in non-gonadal tissues

There is evidence that X- and Y-linked genes contribute to the development of non-gonadal tissues...what about the brain?

There is evidence that X- and Y-linked genes contribute to the development of non-gonadal tissues...what about the brain?

- Gene expressions measured from post-mortem brain samples Vawter et al. (2004)
 - · three brain regions
 - three laboratories (UC Irvine, UC Davis, UMichigan Ann Arbor)
 - 10 subjects (5 male, 5 females)
- · We select a subset of genes:
 - From Vawter et al. (2024): Y-linked W UTY, USP9Y, SMCY, DBY, RPS4Y
 - From Vawter et al. (2024): X-linked 🖙 XIST
 - \star top 14 genes with highest variability across subjects/labs/brain regions (H=20)

Each subject represents a group, g = 1, ..., G = 10 so model the genes within each group via regression:

ie

$$\begin{aligned} & Y_{gi} \overset{\text{ind.}}{\sim} N(\theta_{cgh_i} + \beta \, \textbf{\textit{x}}_i, \sigma_{cgh_i}^2), \quad \textbf{\textit{x}}_i = [\mathsf{Lab}_i, \mathsf{Brain Region}_i] \\ & \mathbb{P}\left(c_{gh} = m \mid M, \textbf{\textit{W}}\right) = W_{gm}, \quad m = 1, \dots, M, \quad h = 1, \dots, H \\ & \left(\theta_1, \sigma_1^2\right), \dots, \left(\theta_M, \sigma_M^2\right) \mid M \overset{\text{i.i.d.}}{\sim} \mathsf{N-inv-\Gamma}\left(\theta, \sigma^2 \mid m_0, k_0, a_0, b_0\right) \\ & \textbf{\textit{W}} \sim \mathsf{W-DMP}(\alpha = 1, \Psi), \quad M \sim \mathsf{Poi}_1\left(\Lambda\right) \end{aligned}$$

where, within each group, we indicate by h_i the gene measured in the i-th observation.

Each subject represents a group, g = 1, ..., G = 10

ie

model the genes within each group via regression:

$$\begin{split} &Y_{gi} \overset{\text{ind.}}{\sim} \textit{N}(\theta_{\textit{c}gh_{i}} + \beta \, \textbf{\textit{x}}_{i}, \sigma_{\textit{c}gh_{i}}^{2}), \quad \textbf{\textit{x}}_{i} = [\mathsf{Lab}_{i}, \mathsf{Brain} \; \mathsf{Region}_{i}] \\ &\mathbb{P}\left(\textit{c}_{gh} = m \mid \textit{M}, \textbf{\textit{W}}\right) = \textit{W}_{gm}, \quad m = 1, \ldots, \textit{M}, \quad h = 1, \ldots, \textit{H} \\ &\left(\theta_{1}, \sigma_{1}^{2}\right), \ldots, \left(\theta_{M}, \sigma_{M}^{2}\right) \mid \textit{M} \overset{\text{i.i.d.}}{\sim} \; \mathsf{N}\text{-inv-}\Gamma\left(\theta, \sigma^{2} \mid \textit{m}_{0}, \textit{k}_{0}, \textit{a}_{0}, \textit{b}_{0}\right) \\ &\textbf{\textit{W}} \sim \mathsf{W-DMP}(\alpha = 1, \Psi), \quad \textit{M} \sim \mathsf{Poi}_{1}\left(\Lambda\right) \end{split}$$

where, within each group, we indicate by h_i the gene measured in the i-th observation.

we estimate clustering of the genes across subjects and experimental conditions

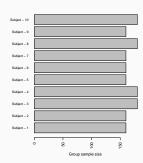
model the genes within each group via regression:

$$\begin{split} & Y_{gi} \overset{\text{ind.}}{\sim} \textit{N}(\theta_{\textit{c}gh_{i}} + \beta \, \textbf{\textit{x}}_{i}, \sigma_{\textit{c}gh_{i}}^{2}), \quad \textbf{\textit{x}}_{i} = [\mathsf{Lab}_{i}, \mathsf{Brain} \, \mathsf{Region}_{i}] \\ & \mathbb{P}\left(\textit{c}_{gh} = m \mid \textit{M}, \textbf{\textit{W}}\right) = \textit{W}_{gm}, \quad m = 1, \ldots, \textit{M}, \quad h = 1, \ldots, \textit{H} \\ & \left(\theta_{1}, \sigma_{1}^{2}\right), \ldots, \left(\theta_{\textit{M}}, \sigma_{\textit{M}}^{2}\right) \mid \textit{M} \overset{\text{i.i.d.}}{\sim} \, \mathsf{N-inv-\Gamma}\left(\theta, \sigma^{2} \mid \textit{m}_{0}, \textit{k}_{0}, \textit{a}_{0}, \textit{b}_{0}\right) \\ & \textbf{\textit{W}} \sim \mathsf{W-DMP}(\alpha = 1, \Psi), \quad \textit{M} \sim \mathsf{Poi}_{1}\left(\Lambda\right) \end{split}$$

where, within each group, we indicate by h_i the gene measured in the i-th observation.

we estimate clustering of the genes across subjects and experimental conditions

Number of observations within each group (total n = 1680)



Each subject represents a group,
$$g = 1, ..., G = 10$$

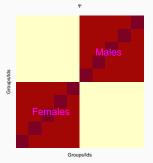
model the genes within each group via regression:

$$\begin{split} & Y_{gi} \overset{\text{ind.}}{\sim} \textit{N}(\theta_{c_{gh_i}} + \beta \, \textbf{\textit{x}}_i, \sigma_{c_{gh_i}}^2), \quad \textbf{\textit{x}}_i = [\mathsf{Lab}_i, \mathsf{Brain} \, \mathsf{Region}_i] \\ & \mathbb{P}\left(c_{gh} = m \mid \textit{M}, \textbf{\textit{W}}\right) = \textit{W}_{gm}, \quad m = 1, \ldots, \textit{M}, \quad h = 1, \ldots, \textit{H} \\ & \left(\theta_1, \sigma_1^2\right), \ldots, \left(\theta_{\textit{M}}, \sigma_{\textit{M}}^2\right) \mid \textit{M} \overset{\text{i.i.d.}}{\sim} \mathsf{N-inv-\Gamma}\left(\theta, \sigma^2 \mid m_0, k_0, a_0, b_0\right) \\ & \textbf{\textit{W}} \sim \mathsf{W-DMP}(\alpha = 1, \Psi), \quad \textit{M} \sim \mathsf{Poi}_1\left(\Lambda\right) \end{split}$$

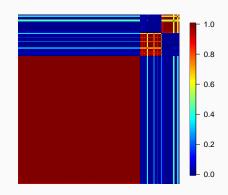
where, within each group, we indicate by h_i the gene measured in the i-th observation.

we estimate clustering of the genes across subjects and experimental conditions

Scale matrix Ψ \blacksquare dependence information specified via the covariate Sex_i



- Posterior co-clustering probabilities for H = 20 unique genes
- We actually cluster H × G = 200 different variables (dimension of the matrix rightarrow)
- · Three main clusters are identified



Expressions of sex-specific genes in brain tissue



 Distinct clustering structures found for Male and Female groups

· Females:

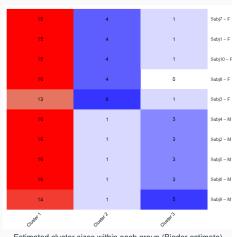
- Cluster 2: contains Y-linked genes
- Cluster 3: contains XIST (X-linked) overexpressed

Males:

- Cluster 2: contains XIST (X-linked)
- Cluster 3: contains Y-linked genes overexpressed

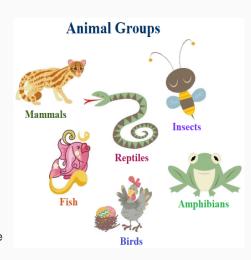
Cluster summary:

- · Cluster 2: low-expressed genes
- · Cluster 3: overexpressed genes
- · Cluster 1: all the rest

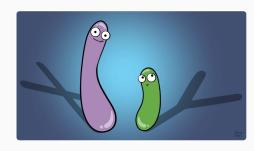


Estimated cluster sizes within each group (Binder estimate)

- New framework for grouped data via mixtures with random number of components
- Joint modelling of weights via matrix construction
- Prior correlation btw rpms encoded through matrix $\boldsymbol{\Psi}$
- Improved estimation accuracy
- Applied to sex-specific gene expression; results align with literature and reveal gene clusters



- New framework for grouped data via mixtures with random number of components
- Joint modelling of weights via matrix construction
- Prior correlation btw rpms encoded through matrix $\boldsymbol{\Psi}$
- Improved estimation accuracy
- Applied to sex-specific gene expression; results align with literature and reveal gene clusters



- Argiento & De Iorio (2022). "Is infinity that far? A Bayesian nonparametric perspective of finite mixture models". The Annals of Statistics, 50(5), 2641-2663.
- Colombi, Argiento, Camerlenghi and Paci (2024) "Hierarchical Mixture of Finite Mixtures" Bayesian Analysis, 1(1), pp.1-29.
- Olkin & Rubin (1964). "Multivariate beta distributions and independence properties of the Wishart distribution". The Annals of Mathematical Statistics, 261-269.
- Gupta & Nagar (2018). "Matrix variate distributions". Chapman and Hall/CRC.
- Vawter et al. (2004) "Gender-specific gene expression in post-mortem human brain: localization to sex chromosomes". Neuropsychopharmacology, 29(2), 373-384.
- Lee et al. (2004) "Molecular evidence for absence of Y-linkage of the Hairy Ears trait".
 European journal of human genetics, 12(12), 1077-1079.

We consider Wishart-type rvs for any integer α , a *Generalised Wishart* distribution.

We consider Wishart-type rvs for any integer α , a *Generalised Wishart* distribution.

Let $S \sim \text{Gen-Wishart}_G(2\alpha, \Psi)$, then:

$$f_{\text{Gen-Wishart}}(S) = \begin{cases} \frac{(\det S)^{\frac{\alpha - G - 1}{2}} \exp\left\{-\frac{1}{2}\operatorname{tr}\left(\Psi^{-1}S\right)\right\}}{2^{\frac{\alpha G}{2}}\Gamma_{G}\left(\frac{\alpha}{2}\right)(\det \Psi)^{\frac{\alpha}{2}}} & \alpha > G/2 \\ \\ \frac{(\det S_{11})^{\frac{\alpha - G - 1}{2}} \exp\left\{-\frac{1}{2}\operatorname{tr}\left[\Psi^{-1}S\right]\right\}}{\pi^{\alpha(G - \alpha)/2}2^{\frac{\alpha G}{2}}\Gamma_{\alpha}\left(\frac{\alpha}{2}\right)(\det \Psi)^{\frac{\alpha}{2}}} & \alpha \leq G/2 \end{cases}$$

where in the singular case $S = \begin{pmatrix} S_{11} & S_{12} \\ S'_{12} & S_{22} \end{pmatrix}$ and $S_{22} = S'_{12}S_{11}^{-1}S_{12}$.

1. <u>Cluster Allocations and Shared Atoms</u> ➤ simple full-conditionals

- 1. <u>Cluster Allocations and Shared Atoms</u> ➤ simple full-conditionals
- 2. Weights $ightharpoonup s_{gm} | \cdot \sim$ mixtures of Gammas depending on n_{gm} , for $g=1,\ldots,2\alpha$ We obtain a single Gamma rv when $n_{gm}=0$ or $\Psi=\mathbb{I}_G$

- 1. <u>Cluster Allocations and Shared Atoms</u> ➤ simple full-conditionals
- 2. Weights $ightharpoonup s_{gm} | \cdot \sim$ mixtures of Gammas depending on n_{gm} , for $g=1,\ldots,2\alpha$ We obtain a single Gamma rv when $n_{gm}=0$ or $\Psi=\mathbb{I}_G$
- Number of components M ➤ unnormalised weights construction.
 Let U = diag(u₁,..., u_G) auxiliary variables, M^(empty) := number of empty components (across groups):

$$q_{M^{(empty)}}(m) = \frac{(m + M^{(obs)})!}{m!} q_{M}(m + M^{(obs)}) \psi(\boldsymbol{u})^{m}, \quad \psi(\boldsymbol{u}) = \det\left(\mathbb{I}_{G} + \frac{1}{2}\Psi U\right)^{-\alpha}$$

When $q_M = Poi_1(\Lambda)$:

$$egin{aligned} q_{M^{(empty)}}(\emph{\emph{m}}) &= \pi_0 \mathsf{Poi}_0\left(\Lambda \psi\left(\emph{\emph{\emph{u}}}
ight)
ight) + \pi_1 \mathsf{Poi}_1\left(\Lambda \psi\left(\emph{\emph{\emph{u}}}
ight)
ight) \\ \pi_0 &= \dfrac{\mathit{\emph{M}}^{(obs)}}{\mathit{\emph{\emph{M}}}^{(obs)} + \Lambda \psi\left(\emph{\emph{\emph{\emph{u}}}}
ight)}, \pi_1 = 1 - \pi_0 \end{aligned}$$

analogous to Argiento and De Iorio (2022) and Colombi et al. (2024) constructions.