

# Finding structures in observations: consistent(?) clustering analysis.

Clara Grazian

School of Mathematics and Statistics, University of Sydney

*clara.grazian@sydney.edu.au*

*OBayes Conference*

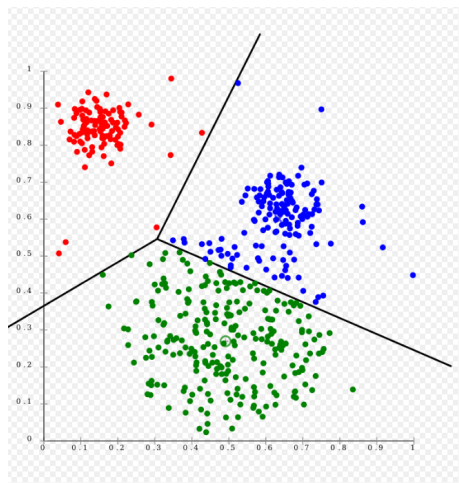
*(Santa Cruz)*

10 September 2022

# Clustering

## Clustering:

- “unsupervised learning”
- requires data, but no labels
- **detect patterns**, e.g.
  - online search results
  - customer shopping patterns
  - effect of pollution
  - animal behaviours
  - cells, tissues, etc
  - regions of images
- common initial analysis: useful when you have no idea
- how to interpret results?



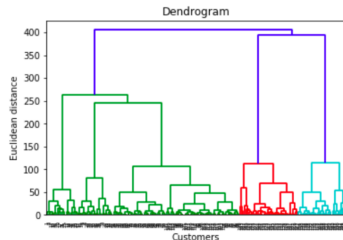
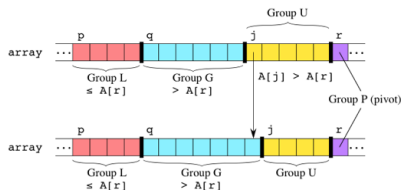
# Clustering algorithms

## Partitioning algorithms:

- k-means
- mixture models
- spectral clustering

## Hierarchical algorithms:

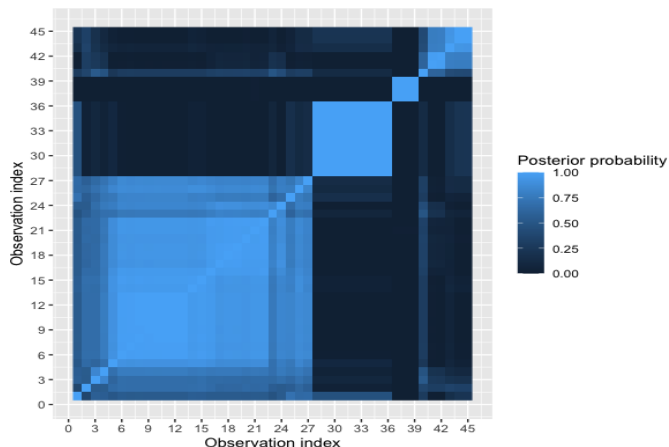
- bottom-up, agglomerative
- top-down, divisive



# Examples of clustering

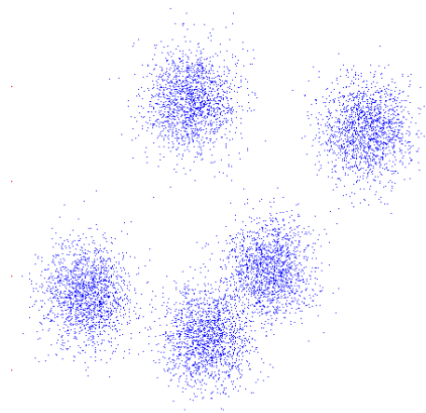
## Clustering gene expression data:

Find clusters of cells with similar biological expression



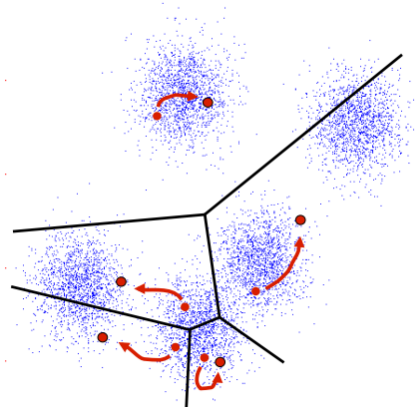
## An iterative algorithm:

- **Initialise:** pick  $K$  random points as cluster centers
- **Alternate:**
  - assign data points to closest cluster center
  - change the cluster center to the average of its assigned points
- **Stop:** when there is no change in assignments

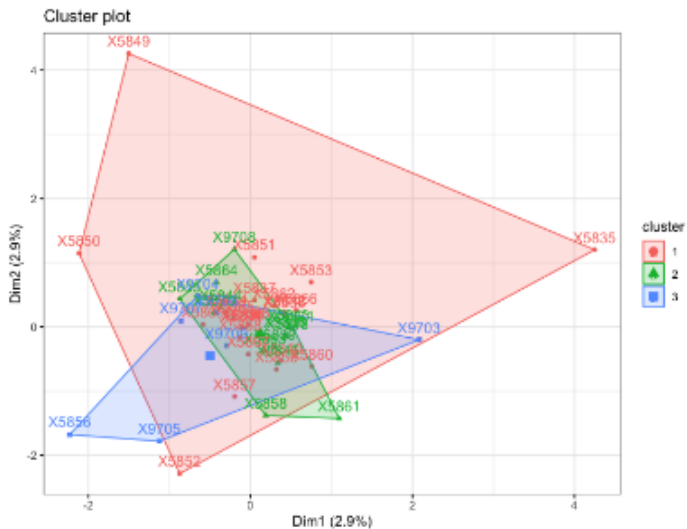


## An iterative algorithm:

- **Initialise:** pick  $K$  random points as cluster centers
- **Alternate:**
  - assign data points to closest cluster center
  - change the cluster center to the average of its assigned points
- **Stop:** when there is no change in assignments



# K-means



Consider the following mixture model

$$g(y; \psi) = \sum_{j=1}^K p_j f_j(y; \theta_j)$$

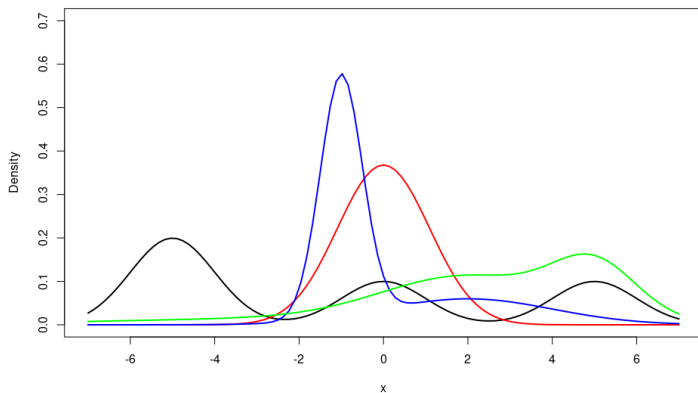
where

- $\psi = (\theta_1, \dots, \theta_K, p_1, \dots, p_K)$
- $p_j \geq 0$  for  $j = 1, \dots, K$
- $\sum_j p_j = 1$ .
- $f_j(\cdot)$  is any probability distribution

These models provide a flexible tool for statistical inference (even in a nonparametric setting, see Lindsay (1995), Roeder (1992) and Roeder and Wasserman (1997)).



# Mixture models



# We need a prior distribution!

Often in Bayesian inference, we want to reduce the effect of the prior on the posterior distribution, in case we do not have strong prior information.

Sometimes so-called “improper” priors are used

$$\int_{\Theta} \pi(\theta) d\theta = \infty.$$

This is not a pdf (or a pmf), therefore the Bayes' theorem cannot be applied.

However they are used in practice as limit of proper prior distributions, when they assure a proper posterior distribution.

It is delicate to produce a noninformative prior for the parameters of a mixture model, since they are often *improper*.

**Why can't we use improper priors?**

**Example:**

Consider independent improper priors

$$\pi(\theta_1, \dots, \theta_K) \propto \prod_{j=1}^K \pi(\theta_j)$$

such that  $\int_{\Theta} \pi(\theta_j) d\theta_j = \infty$

# Classical noninformative solutions - contd

The mixture model is a classical example of latent variable model, then it can be rewritten as

$$g(y; \psi) = \sum_{S \in \mathcal{S}_k} \prod_{j=1}^K f(y; S, \theta_j) \pi(\theta_j) \pi(S | \rho) \pi(\rho)$$

where the summation runs over all  $k^N$  possible classifications  $S$ .

Then the complete-data likelihood is non-informative if there is an empty component (let's say the  $j$ -th)

$$\int \prod_{i: S_i=j} f(y_i; \theta_j) \pi(\theta_j) d\theta_j \propto \int \pi(\theta_j) d\theta_j = \infty$$

# Random partition models

- A random partition model is a probability distribution over  $\mathcal{P}_n$

$$\{p(\rho_n = (S_1, \dots, S_K)) : \rho_n \in \mathcal{P}_n\}$$

- One main approach is to **to define  $p(\rho_n)$  through discrete random probability measures**

$$f(y_1, \dots, y_n | \theta_1, \dots, \theta_n) = \prod_{i=1}^n f(y_i | \theta_i)$$

$$\theta_1, \dots, \theta_n | G \stackrel{i.i.d.}{\sim} G$$

$$G \sim \text{discrete RPM}$$

For example, if  $G(\cdot) = \sum_{h=1}^K p_h \delta_{\psi_h}$  with  $P(\sum_{h=1}^K p_h = 1) = 1$ , then

$$g(y_i | \{\psi_h\}) = \sum_{h=1}^K p_h f(y_i | \psi_h)$$

# Random partition models

- Discreteness of  $G$  implies existence of ties among  $\theta_1, \dots, \theta_n$
- If  $\psi_1, \dots, \psi_K$  denote the corresponding **unique values** then we can define  $\rho_n$  through indicators given as

$$c_i = j \Leftrightarrow \theta_i = \psi_j \text{ or equivalently } \theta_i = \psi_{c_i}$$

and so  $S_j = \{i \in [n] : \theta_i = \psi_j\}$ , where  $[n] = \{1, \dots, n\}$  is the set of  $n$  indices.

- This is an **induced random partition model**.

# Is it easy to define a prior over the partition space?

- $\rho_n = (S_1, \dots, S_K)$  a partition of  $[n]$  into  $K = |\rho_n| \geq 1$  nonempty (and mutually exclusive) subsets;
- $\mathcal{P}_n$ : set of all partitions of  $[n]$ ;
- the size of  $\mathcal{P}_n$  increases as the Bell number; e.g.  $B_{10} = 115,975$

- Consider  $K \sim p_K(k)$

$$y_i | K = k, p_1, \dots, p_K, \theta_1, \dots, \theta_K \stackrel{i.i.d.}{\sim} \sum_{h=1}^K p_h f(y_i | \theta_h)$$

$$\theta_1, \dots, \theta_k | K = k \stackrel{i.i.d.}{\sim} p_0(\theta)$$

$$(p_1, \dots, p_k) | K = k \sim \text{Dir}(\gamma, \dots, \gamma)$$

- The induced partition model is then

$$p(\rho_n = (S_1, \dots, S_K)) = \left( \sum_{h=1}^{\infty} \frac{h^{(k)}}{(\gamma h)^{(n)}} p_K(h) \right) \left( \prod_{s \in (S_1, \dots, S_K)} \gamma^{|s|} \right)$$

where  $x^{(m)} = x(x+1)\dots(x+m-1)$  and  
 $x_{(m)} = x(x-1)\dots(x-m+1)$ .



# CASE 1: Overfitted mixtures

- One approach consists in fixing  $K$  to a large value and use inference to estimate some of the weights as equal to zero, in order to identify the correct  $k < K$  number of clusters.
- Playing on the prior for  $(p_1, \dots, p_K)$
- [Rousseau and Mengersen (2011)] show the asymptotic behaviour of the posterior distribution in a mixture model for overfitted mixtures: the posterior distribution concentrates on a sparse representation of the true density; this is exhibited by a subset of components that adequately describe the density remaining and any superfluous components becoming empty.
- IMPORTANT: need for a prior on the weights that favour small weights (Dirichlet with parameters  $1/2$ ).

We recall that the Jeffreys prior was introduced by Jeffreys (1939) as a default prior based on the Fisher information matrix

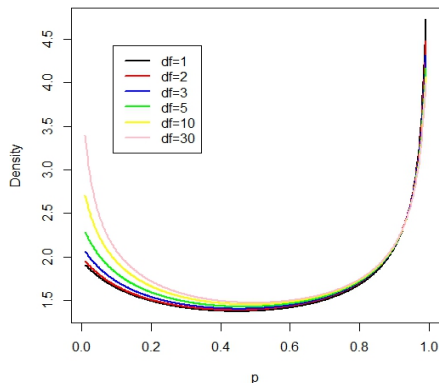
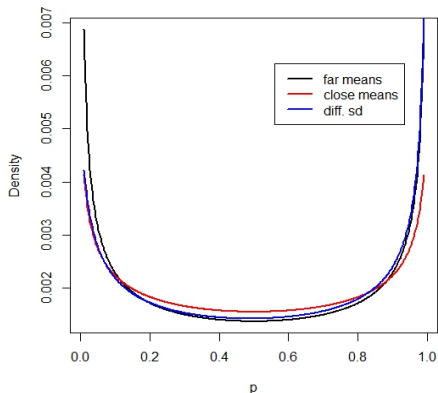
$$\pi^J(\boldsymbol{\theta}, \mathbf{p}) \propto |I(\boldsymbol{\theta}, \mathbf{p})|^{1/2} = \det \left( \mathbb{E}_g \left[ -\frac{d^2}{d\boldsymbol{\psi}^2} \log g(y; \mathbf{p}, \boldsymbol{\theta}) \right] \right)^{1/2}$$

- when using the Jeffreys' prior for all the parameters of the model, the posterior is improper (**OH NO!**)
- but when the Jeffreys' prior is used only for the weights, it can be shown that it leads to the same results as Rousseau & Mengersen (2011)!

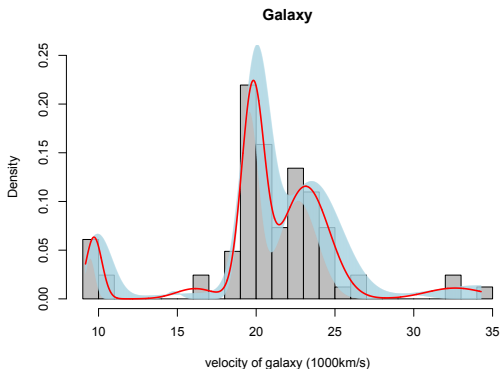
# Jeffreys prior for the weights

Instead, we fix the Jeffreys prior only for the weights conditionally on the other parameters

$$\pi^J(p_1, \dots, p_K | \theta_1, \dots, \theta_K) \propto |I(p_1, \dots, p_K)|^{1/2}$$



# The galaxy dataset



---

**Dataset:** galaxy

---

$p_1$  0.437  
(23.139, 1.507)

$p_2$  0.390  
(19.790, 0.715)

$p_3$  0.080  
( 9.709, 0.503)

$p_4$  0.056  
(32.630, 1.842)

$p_5$  0.037  
(16.138, 1.226)

$\sum_{\ell=6}^{10} p_{\ell}$  0.000

## CASE 2: A prior on the number of the # of components

The induced partition model is then

$$p(\rho_n = (S_1, \dots, S_K)) = \left( \sum_{h=1}^{\infty} \frac{h^{(k)}}{(\gamma h)^{(n)}} p_K(h) \right) \left( \prod_{s \in (S_1, \dots, S_K)} \gamma^{|s|} \right)$$

where  $x^{(m)} = x(x+1)\dots(x+m-1)$  and  $x_{(m)} = x(x-1)\dots(x-m+1)$ .

# A prior on the number of components

For model

$$g(y; \psi) = \sum_{j=1}^K p_j f_j(y; \theta_j)$$

the prior can be specified as

$$\pi(k, \mathbf{p}, \theta) = p_K(k) \pi(\mathbf{p} | k) \pi(\theta | k).$$

The posterior for  $k$  is then given by

$$p_K(k | y) \propto \int f(y | k, \mathbf{p}, \theta) \times p_K(k) \pi(\mathbf{p} | k) \pi(\theta | k) d\mathbf{p} d\theta.$$

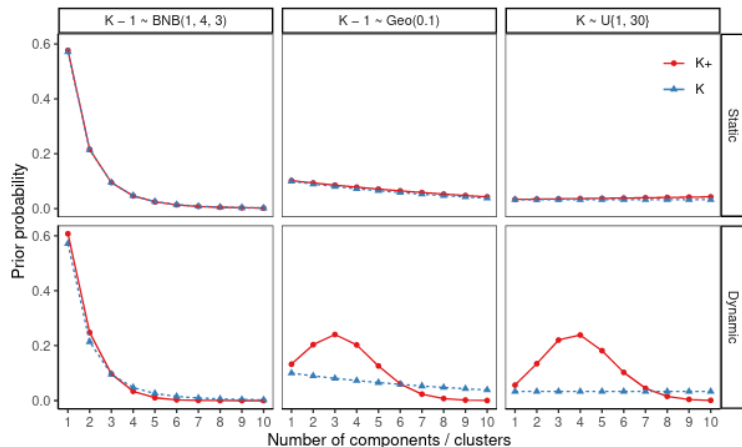
- Although for practical purposes the range of values  $K$  can take is finite, it may be appropriate to define a prior over  $\mathbb{N}$ .
- In fact, by truncating the support of  $K$  there may be possible distortions of the posterior around the boundary, affecting the inferential results.
- **BUT** the prior on  $K$  must be proper, as proved by Nobile (2005).
- Remember the **inconsistency problems** in the nonparametric setting - see Miller and Harrison (2014).

# Possible prior distributions

- $K \sim Unif(0, 30)$  (Richardson and Green, 1997)
- $K \sim Pois(1)$  (Nobile and Fearnside, 2007)
- $K \sim BNB(1, a, b)$   
(Grazian et al. 2020, Frühwirth-Schnatter et al., 2021)
  - Frühwirth-Schnatter et al. (2021) propose to combine a prior on  $K$  with an adaptive prior on the weights  $(p_1, \dots, p_K) \sim Dir(\frac{\alpha}{K}, \dots, \frac{\alpha}{K})$  — “dynamic” version of the mixtures



# Possible prior distributions



Früwirth-Schnatter et al., 2021

**Idea:** inconsistency problems can be prevented by penalising larger values  
→ we can define the prior on  $K$  with a loss-based approach.

To obtain the **loss-based prior** on  $K$ , we define the prior on  $K$  by assigning a prior on the space of models determined by the mixtures with  $k = 1, 2, \dots$  components.

- we can assign a *worth* to each mixture
- we include a component of loss due to the complexity of the model

$$\text{Loss}(k) = \text{Loss}_I(k) + \text{Loss}_C(k)$$

# A loss-based prior: information loss [Grazian et al., 2020]

The quantification of the loss comes from Berk (1966): *if the model is misspecified, the posterior distribution asymptotically tends to accumulate at the most similar model so to minimise the loss in information, in terms of Kullback-Leibler divergence*

If we consider a mixture  $M_s = \{g_s(x|\psi_s), \pi_s(\psi_s)\}$  (where  $\psi_s = (p_s, \theta_s)$ )

$$\text{Loss}_I(k) = \mathbb{E}_{\pi_s} \left\{ \inf_{\psi_m, m \neq j} D_{KL} \left( g_s(x|\psi_s) \parallel g_m(x|\psi_m) \right) \right\},$$

The above loss is linked to the prior mass by means of the *self-information* loss function which associate a loss to a probability statement. As such,

$$p_K(k) \propto \exp \{ \text{Loss}_I(k) \}.$$

**The loss attains its minimum at zero:** Consider mixture  $g_k = \sum_{j=1}^k p_j f_j(x|\theta_j)$  and  $g_{k+1} = \sum_{j=1}^k \check{p}_j f_j(x|\check{\theta}_j) + \check{p}_{k+1} f_{k+1}(x|\check{\theta}_{k+1})$ ; the minimum is obtained when  $\check{p}_j = p_j$  and  $\check{\theta}_j = \theta_j$  and  $\check{p}_{k+1} = 0$ .

## A loss-based prior: complexity loss [Grazian et al., 2020]

To fully describe the *worth* of a mixture model it is also necessary to take into consideration its complexity.

If we keep the mixture model with  $k$  components, the loss would be related to the number of parameters that have to be estimated, and therefore the number of components.

$$\text{Loss}_C(k) = U(\text{keep } k) = -c \cdot k.$$

Therefore,

$$p_K(k) \propto \exp\{-c \cdot k\},$$

where  $c > 0$  is included as loss functions are defined up to a constant.

# A reparametrization

## Theorem

Consider the prior distribution for the number of components of a finite mixture model, where we set  $p = \exp\{-c\}$  and  $k = 1, 2, \dots$ . If we choose  $p \sim \text{Beta}(\alpha, \beta)$ , with  $\alpha, \beta > 0$ , then

$$p_K(k|p) = p^{k-1}(1-p),$$

which is a geometric distribution with parameter  $1-p$ , and

$$p_K(k) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \frac{\Gamma(k + \beta - 1)\Gamma(\alpha + 1)}{\Gamma(k + \alpha + \beta)},$$

which is a beta-negative-binomial distribution where the number of failures before the experiment is stopped is equal to 1, and shape parameters  $\alpha$  and  $\beta$ .

The linear complexity loss is a choices; other choices are also possible.

# A loss-based prior

The prior  $p_K(k)$  just defined

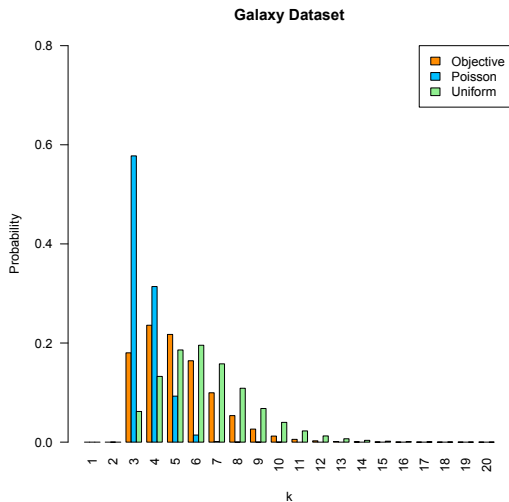
- is defined on the whole support of  $K$ ,  $\mathbb{N}$
- is proper
- has moments

$$\mathbb{E}(K) = \mathbb{E}(\mathbb{E}\{K|p\}) = \mathbb{E}(p^{-1}) = \frac{\alpha + \beta - 1}{\alpha - 1}, \quad \text{for } \alpha > 1,$$

$$\text{Var}(K) = \mathbb{E}(\text{Var}\{K|p\}) + \text{Var}(\mathbb{E}\{K|p\}) = \frac{\alpha\beta(\alpha + \beta - 1)}{(\alpha - 2)(\alpha - 1)^2}, \quad \text{for } \alpha > 2.$$

where  $\beta$  can be used to control how many components we assume a priori and  $\alpha$  can be used to control the variance.

# The galaxy dataset



## CASE 3: Introducing covariates

Suppose that the observations depend on a covariate, e.g. they are time-dependent. → **we can use hidden Markov models!**

Let

- $\{t_1, t_2, \dots, t_T\} \equiv \mathcal{T}$ : set of observed time points
- $\mathbf{y} = \{\mathbf{y}_t\}_{t \in \mathcal{T}}$ : the data
- $\mathbf{c} = \{c_t\}_{t \in \mathcal{T}}$ : a latent variable indicating the cluster each observation belongs to, with  $c_t \in \{1, 2, \dots, K\} \equiv \mathcal{K}$  and  $K$

We assume that the data come from a mixture-type model

$$g(\mathbf{y}|\mathbf{c}, \{\theta_k\}_{k \in \mathcal{K}}) = \prod_{t \in \mathcal{T}} \prod_{k \in \mathcal{K}} f(y_t | \theta_k)^{\mathbb{1}_k(c_t)}$$

i.e. given the latent variables  $\mathbf{c}$ , the observations  $y_t$  are independent.



One possible solution

$$g(\mathbf{y}|\mathbf{c}, \{\theta_k\}_{k \in \mathcal{K}}) = \prod_{t \in \mathcal{T}} \prod_{k \in \mathcal{K}} f(y_t | \theta_k)^{\mathbb{I}_k(c_t)},$$
$$c_t \sim \text{Discrete}(\mathbf{p}_t),$$
$$\mathbf{p}(t) \sim \text{LogitGP}(\mathbf{A}, \mu(t), \mathbf{C}(h))$$

So that the probabilities  $\mathbf{p}_t = \{\mathbf{p}_{t,k}\}_{k \in \mathcal{K}}$  are discrete-time observations of an underlying and non-observed continuous-time process  $\mathbf{p}(t)$ .

We have that

- $\mathbf{A}$  is a co-regionalization matrix
- $\mu(t)$  is a mean function
- $\mathbf{C}(h)$  is a correlation function with  $h$  being a temporal distance

# LogitN distribution and LogitGP

[Aitchison, 1986] proposed the LogitN distribution to model compositional data as an alternative to the Dirichlet distribution.

The vector  $\mathbf{p}_t$  is defined as

$$p_{t,k} = \frac{e^{\omega_{t,k}}}{\sum_{j=1}^K e^{\omega_{t,j}}}, \quad k \in 1, \dots, K$$

where  $\omega_{t,k}$  are real valued variables.

**Remark:** adding a constant to each  $\omega_{t,k}$  produces the same vector of probabilities, and an identifiability constraint is therefore needed; the  $K$ -th element is set to zero ( $\omega_{t,K} = 0$ ) treated as the *reference element*.

$\omega_t$  can be the realisation of a  $K - 1$  dimensional GP  $\omega(t)$ .

# The covariance of the $\mathbf{p}_t$

**Attention must be paid!** The covariance among each element and the sum of all the element is

$$\text{Cov}(p_{t,k}, p_{t,1} + \dots + p_{t,k} + \dots + p_{t,K}) = 0$$

where  $p_{t,1} + \dots + p_{t,k} + \dots + p_{t,K} = 1$ . Therefore we have

$$-\text{Var}(p_{t,k}) = \sum_{\substack{h=1 \\ k \neq h}}^K \text{Cov}(p_{t,k}, p_{t,h}).$$

Aitchison (1986) pointed out that a more consistent measure of dependence between compositional elements can be measure as

$$\tau_{ij,kl}(t, t') = \text{Cov} \left( \log \frac{p_{t,i}}{p_{t,k}}, \log \frac{p_{t',j}}{p_{t',l}} \right), i, j, k, l \in 1, \dots, K,$$

# The covariance of the $\mathbf{p}_t$

Let's keep things simple and suppose that  $\mathbf{p}_t \sim \text{LogitN}(\boldsymbol{\mu}_t, \boldsymbol{\Sigma}_t)$ , where  $\boldsymbol{\mu}_t$  is  $K - 1$  dimensional vector and  $\boldsymbol{\Sigma}_t$  is a  $(K - 1) \times (K - 1)$  square matrix.

- It can be proved that a *LogitN* process has independent components (in term of log-ratio), i.e.  $\tau_{ij,kl}(t, t') = 0$  for arbitrary  $i, j, k$  and  $l$ , at time lag  $|t - t'|$  only if the variance of the Gaussian variable is

$$\boldsymbol{\Sigma}_{t,t'} = \begin{pmatrix} a_1(t, t') + a_K(t, t') & a_K(t, t') & \dots & a_K(t, t') \\ a_K(t, t') & a_2(t, t') + a_K(t, t') & \dots & a_K(t, t') \\ \dots & \dots & \dots & \dots \\ a_K(t, t') & a_K(t, t') & \dots & a_{K-1}(t, t') + a_K(t, t') \end{pmatrix}$$

where the element  $[\boldsymbol{\Sigma}_{t,t'}]_{i,j}$  is  $\tau_{ij,KK}(t, t')$ .

- The elements of  $\mathbf{p}$  are iid if  $\boldsymbol{\mu}_t = \mathbf{0}$  and

$$\boldsymbol{\Sigma}_{t,t'} = \begin{pmatrix} 2a & a & \dots & a \\ a & 2a & \dots & a \\ \dots & \dots & \dots & \dots \\ a & a & \dots & 2a \end{pmatrix}$$

# A new parametrization [Mastrantonio et al., 2019]

We introduce an auxiliary  $K$ -dimensional GP  $\gamma(t)$ , From  $\gamma(t)$ , we construct  $\omega(t)$  as

$$\begin{aligned}\omega_k(t) &= \gamma_k(t) - \gamma_K(t), \\ \gamma(t) &= \mu(t) + \mathbf{A}\gamma^*(t), \\ \gamma_k^*(t) &\sim \text{GP}(0, C_k(h)).\end{aligned}$$

where

- $\mathbf{A}$  is a coregionalization matrix, which we require to be non-negative definite and symmetric
- $\mu(t)$  is a mean function
- $\mathbf{C}(h)$  is a vector of correlation functions

# Covariance structure

Matrix  $\mathbf{A}$  introduces dependence between the elements of  $\gamma(t)$ , and

$$\boldsymbol{\Sigma} = \mathbf{A}\mathbf{A}'$$

is the covariance of  $\gamma(t)$ .

The explicit relation between  $\mathbf{A}$  and  $\boldsymbol{\Sigma}$  is

$$\mathbf{A} = \boldsymbol{\Delta}\boldsymbol{\Xi}^{\frac{1}{2}}\boldsymbol{\Delta}',$$

where

- $\boldsymbol{\Delta}$  is the matrix of the eigenvectors of  $\boldsymbol{\Sigma}$
- $\boldsymbol{\Xi}$  is the diagonal matrix of the eigenvalues of  $\boldsymbol{\Sigma}$

Then

$$\mathbf{p}(t) \sim \text{LogitGP}(\mathbf{A}, \boldsymbol{\mu}(t), \mathbf{C}(h))$$

It is important to highlight that  $\gamma(t)$  is **not identifiable** and any inference about  $\mathbf{p}_t$  is in fact made by looking at  $\omega_t$  through equation.

$$p_{t,k} = \frac{e^{\gamma_{t,k} - \gamma_{t,K}}}{\sum_{j=1}^K e^{\gamma_{t,j} - \gamma_{t,K}}} = \frac{e^{\gamma_{t,k}}}{\sum_{j=1}^K e^{\gamma_{t,j}}}, \quad k \in 1, \dots, K.$$

With respect to the case where  $\omega_{t,K}$  must be set to zero, this equation has a more symmetric form, since all the components of  $\mathbf{p}_t$  are written in terms of exponentials of  $\gamma_k(t)$  and there is no reference element.

It can be proved that this model assures

- invariance from the choice of the reference element;
- invariance with respect to the reordering of the labels;
- the expected structure of the covariance matrix among times, when defined on  $\tau_{ij,kl}(t, t')$  elements.



## CASE 4: infinite mixtures

A mixture model can be extended to consider infinite components:

$$\begin{aligned}y_i | \theta_i &\sim f(y_i | \theta_i) & i = 1, \dots, n \\ \theta_i | G &\sim G \\ G | \alpha, G_0 &\sim DP(\alpha, G_0),\end{aligned}$$

and, since  $G$  is almost surely discrete, this model can be rewritten as

$$y_i \sim \sum_{h=1}^{\infty} \pi_h f(y_i | \psi_h) \quad i = 1, \dots, n$$

where  $\psi_1, \psi_2, \dots$  are independent draws from the base distribution  $G_0$ .

# Partitions for infinite mixtures

The EPPF of the DP is explicitly available; if  $G \sim DP(\alpha, G_0)$ , then

$$p(\rho_n = (S_1, \dots, S_K)) = \frac{\alpha^K \prod_{h=1}^K (n_h - 1)!}{\prod_{i=1}^n (\alpha + i - 1)!}.$$

which is known as Ewens distribution.

(Generalization to other processes, like the PY process are available)

And the conditional EPPF for a DP mixture model, induced for a given number of clusters  $K = k$ , is

$$p_{DP}(\rho_n = (S_1, \dots, S_k) | K = k) = \frac{1}{\text{Const}} \prod_{h=1}^k \frac{1}{n_h}$$

However, it can be shown that this EPPF favours unbalanced partitions with some small values of  $n_h$  (look at the inverse dependence on  $n_h$ ) → **this model is inconsistent for clustering!**

## CASE 4b: Adding covariates - Grazian (2022+)

Suppose  $Y_t(s)$  can be represented as an infinite mixture model:

$$g(y_t(s)|\pi, \theta) = \sum_{k=1}^{\infty} \pi_{t,k}(s) g(y_t(s)|\theta_k)$$

where the mixing probability  $\pi_{t,k}(s)$  is the probability that the location  $s$  belongs to component  $k$  at time  $t$ .

The mixing weights are built similarly to the spatial stick-breaking:

$$F_t(s) = \sum_{k=1}^{\infty} \pi_{t,k}(s) \delta_{\theta_k} \quad s \in \mathcal{D}, t > 0 \quad \text{where}$$

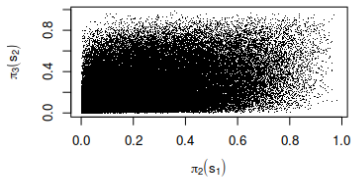
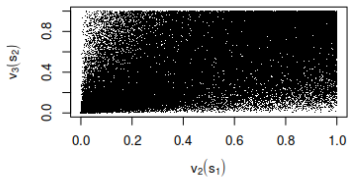
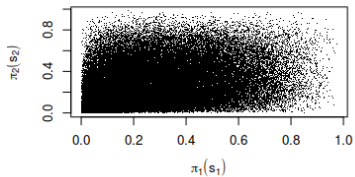
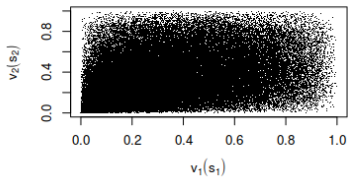
$$\pi_{t,1}(s) = V_{t,1}(s), \quad \pi_{t,k}(s) = V_{t,k}(s) \prod_{j=1}^{k-1} (1 - V_{t,j}(s)) \quad \text{for } k = 2, \dots$$

$$V_{t,k}(s) = w_k(s, \psi, t, \zeta) V_k$$

$$V_k \sim \text{Beta}(a, b)$$

$$\theta_k \sim F_0.$$

# A caveat



Consistently estimating the number of clusters in a Bayesian way is difficult.

## However

- consistency can be found for overfitted mixtures
- the prior may have an important role
- advantage of reducing the number of necessary assumptions and inputs
- easy extension to multivariate setting

- Grazian, C. and Robert, C.P. (2018) Jeffreys priors for mixture estimation: Properties and alternatives. *Computational Statistics & Data Analysis* 121, 149-163.
- Mastrantonio, G., Grazian, C., Mancinelli, S., Bibbona, E. (2019) *New formulation of the logistic normal process to analyze tracking trajectories.. Annals of Applied Statistics*, 13(4),2483–2508.
- Grazian, C., Villa, C. and Liseo, B. (2020) On a loss-based prior for the number of components in mixture models. *Statistics & Probability Letters* 158.
- Grazian, C. (2022+) A review of Bayesian clustering based on mixture models.
- Grazian, C. (2022+) Estimating MIC distributions and cutoffs through mixture models: an application to establish M. Tuberculosis resistance. *bioRxiv* 643429.
- Grazian, C. (2022+) Spatio-temporal stick-breaking.