

Informative model-based clustering via Centered Partition Processes

Sally Paganin

spaganin@hsph.harvard.edu

<https://salleuska.github.io/>

Developments in Bayesian Nonparametrics - YoungStatS

April 21, 2021

Introduction

Clustering is one of the canonical data analysis goal in statistics

- **Distance based methods:** distance metric between data points
- **Model-based clustering:** rely on discrete mixture models

Bayesian perspective : allows to incorporate prior information

What if, we have prior information on the clustering itself?

Motivating application - Birth defects data

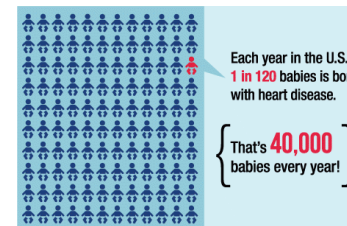
- Relate exposure factors to the development risk of a defect
- **Prior information** available (biology/expert's judgments)

National Birth Defect Prevention Study

- **Population-based case-control study**
 - 300 controls/100 cases per year since 1997
 - monthly n. of controls \propto n. of births previous year
- **Cases** (37 major birth defect)
 - Birth defects surveillance system
 - + clinical genetist review
 - Cases with known etiology were excluded
- **Controls**
 - Non-malformed live birth
 - Birth certificates or hospital delivery records
- **Data collection**
 - CATI (English/Spanish) within 24 months



Congenital Heart Defects



Clinical importance

priority in public health

- most frequent class of defects
- high impact on pediatric mortality

Statistical relevance: challenge in birth defects modeling

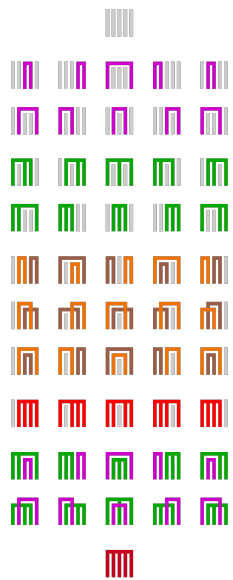
- Some defects are too rare for individual study
- Difficult to determine how best to group birth defects

Experts have provided a **mechanistic classification** of the defects

- relies on biological knowledge and embryologic development
- **translates in a prior guess c_0 for the clustering**

We focus on the **Congenital Heart Defects (CHD)** which are problems in the structure of the heart that are present at birth.

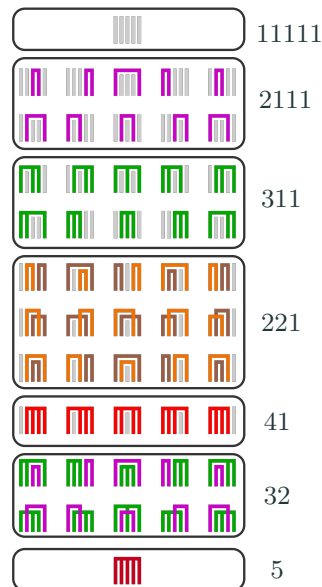
Set partitions



A **set partition** \mathcal{c} of an integer $[n]$ is a collection of non-empty disjoint subsets $\{B_1, B_2, \dots, B_K\}$ such that $\bigcup_{i=1}^K B_i = [n]$

- Number of partitions of $[n]$ into k blocks
→ Stirling numbers $S(n, k)$
- Total number of set partitions
→ Bell number $\mathcal{B}_n = \sum_{k=1}^n S(n, k)$

Set partitions



A **set partition** c of an integer $[n]$ is a collection of non-empty disjoint subsets $\{B_1, B_2, \dots, B_K\}$ such that $\cup_i^K B_i = [n]$

- Number of partitions of $[n]$ into k blocks
→ Stirling numbers $S(n, k)$
- Total number of set partitions
→ Bell number $\mathcal{B}_n = \sum_{k=1}^n S(n, k)$
- **Configuration** $\lambda = \{|B_1|, \dots, |B_K|\}$
→ sequence of block cardinalities
→ individuate an **integer partition**, a set of positive integers $\{\lambda_1, \dots, \lambda_K\}$ such that $\sum_{i=1}^K \lambda_i = n$

Modeling birth defects

- $i = 1, \dots, N$ heart defects, $j = 1, \dots, n_i$ observations
- $y_{ij} = 1$ if observation j has the b.d. i while $y_{ij} = 0$ is a control
- $\mathbf{x}_{ij}^T = (x_{ij1}, \dots, x_{ijp})$ observed values for p dichotomous variables

Grouped logistic regression

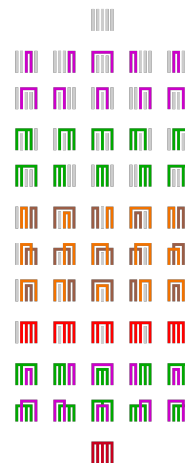
$$\begin{aligned} y_{ij} &\sim \text{Ber}(\pi_{ij}) & \text{logit}(\pi_{ij}) &= \alpha_i + \mathbf{x}_{ij}^T \boldsymbol{\beta}_{c_i}, \quad j = 1, \dots, n_i, \\ \alpha_i &\sim \mathcal{N}(a_0, \tau_0^{-1}) & \boldsymbol{\beta}_{c_i} | \mathbf{c} &\sim \mathcal{N}_{\mathcal{P}}(\mathbf{b}, \mathbf{Q}) \quad i = 1, \dots, N, \end{aligned}$$

Bayesian framework: assign a prior probability $p(c)$

→ *Exchangeable Partition Probability Function (EPPF)*

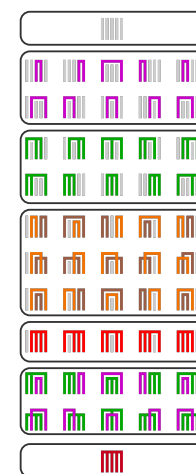
Uniform distribution

$$p(\mathbf{c}) \propto 1/\mathcal{B}_N$$



Dirichlet Process: $p(\mathbf{c}) \propto \prod_{i=1}^K (|B_i| - 1)!$

Pitman-Yor Process: $p(\mathbf{c}) \propto \prod_{i=1}^K (1 - \sigma)_{|B_i|}$



How to account for c_0 ?

Base idea: penalize a baseline EPPF in order to center the prior distribution on the given partition c_0

$$p(c|c_0, \psi) \propto p_0(c) \exp\{-\psi d(c, c_0)\} \quad (1)$$

- $p_0(c)$ indicates a **baseline distribution** (EPPF) on Π_N
- $d(c, c_0)$ a suitable **distance** between partitions
 - ideally a metric on the set partitions lattice
 - Variation of information [Meila (2007)]

$$VI(c, c') = -H(c) - H(c') + 2H(c \wedge c')$$

- ψ **penalization parameter** controlling for the centering
 - $\psi = 0$ $p(c|c_0, \psi) \rightarrow p_0(c)$
 - $\psi \rightarrow \infty$ $p(c|c_0, \psi) = \delta_{c_0}$

Centered Partition Processes

Define sets of partitions with distance δ_l from c_0 and configuration λ_m

$$s_{lm}(c_0) = \{c \in \Pi_N : d(c, c_0) = \delta_l, \Lambda(c) = \lambda_m\}$$

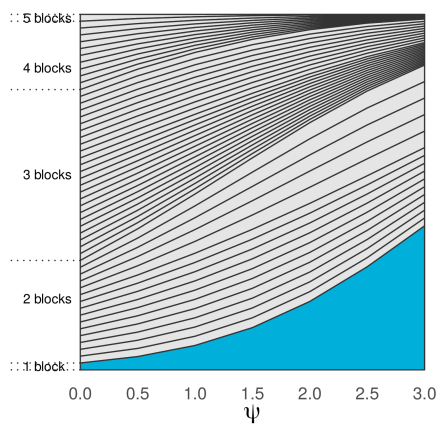
for $l = 0, \dots, L$ and $m = 1, \dots, M$.

Centered Partition Processes - analytic form

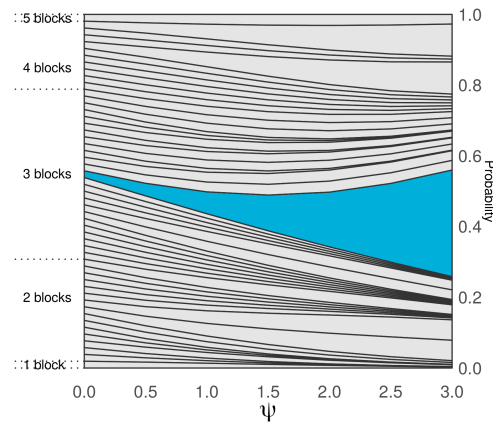
$$p(c|c_0, \psi) = \frac{g(\lambda_m) e^{-\psi \delta_l}}{\sum_{u=0}^L \sum_{v=1}^M |s_{uv}(c_0)| g(\lambda_v) e^{-\psi \delta_u}}, \quad \text{for } c \in s_{lm}(c_0)$$

- $g(\cdot)$ function of the configuration $\Lambda(c)$
 - e.g. Uniform $g(\Lambda(c)) = 1$, DP $g(\Lambda(c)) = \alpha^K \prod_{j=1}^K \Gamma(\lambda_j)$
- $|\cdot|$ cardinality of the set $s_{lm}(c_0)$, not analytically tractable
 - but can be used in Bayesian models relying on MCMC

CP Process - Uniform EPPF

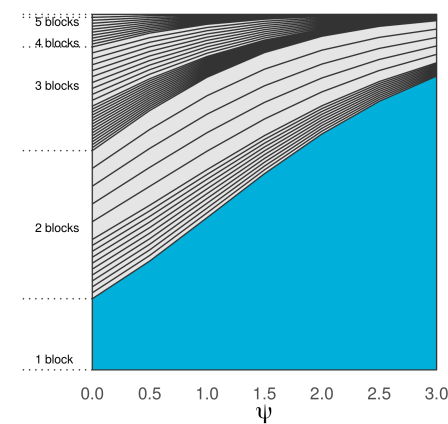


$$c_0 = \{1, 2, 3, 4, 5\}$$

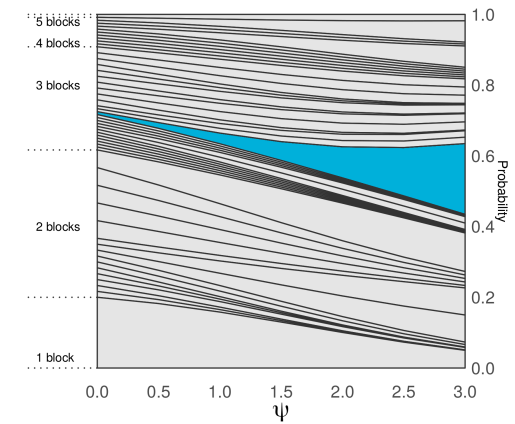


$$c_0 = \{1, 2\}\{3, 4\}\{5\}$$

CP Process - DP EPPF ($\alpha = 1$)



$$c_0 = \{1, 2, 3, 4, 5\}$$



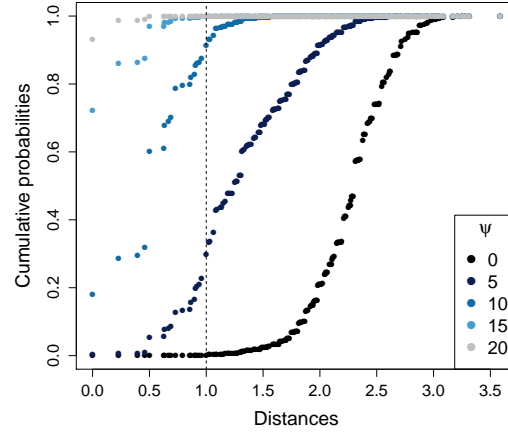
$$c_0 = \{1, 2\}\{3, 4\}\{5\}$$

Prior calibration

We consider to estimate the distribution of **distance** $\delta \in \{\delta_l\}_{l=0}^L$

$$p(\delta = \delta_l) = \frac{\sum_{m=1}^M n_{lm} g(\lambda_m) e^{-\psi \delta_l}}{\sum_{u=0}^L \sum_{v=1}^M n_{uv} g(\lambda_v) e^{-\psi \delta_u}}$$

- **Monte Carlo procedure**
 - uniform sampler on the set partition space Π_N [Stam (1983)]
- **Deterministic local search**
 - for small values of the distance $\delta \in \{\delta_0, \dots, \delta_{L^*}\}$
 - greedy search algorithm



Modeling birth defects

$N = 26$ birth defects, 4,047 cases, 8,125 controls, 90 potential risk factors

$$y_{ij} \sim \text{Ber}(\pi_{ij}) \quad \logit(\pi_{ij}) = \alpha_i + \mathbf{x}_{ij}^T \boldsymbol{\beta}_{c_i}, \quad j = 1, \dots, n_i,$$

$$\alpha_i \sim \mathcal{N}(a_0, \tau_0^{-1}) \quad \boldsymbol{\beta}_{c_i} | \mathbf{c} \sim \mathcal{N}(\mathbf{b}, \mathbf{Q}) \quad i = 1, \dots, N,$$

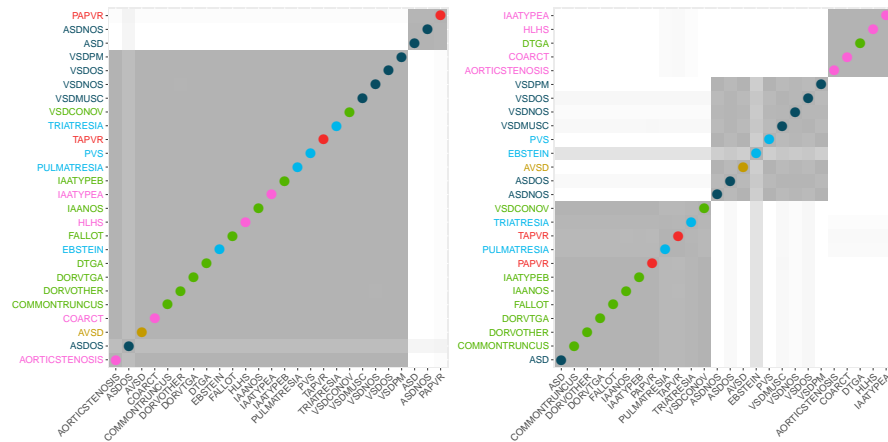
$$p(\mathbf{c}) \sim CP(\mathbf{c}_0, \psi, p_0(\mathbf{c})) \quad p_0(\mathbf{c}) \propto \alpha^K \prod_{k=1}^K (\lambda_k - 1)!$$

from the prior calibration: $\psi = 40$ (90% partitions with $d = 0.8$ ($d_{max} = 4.70$))

Posterior estimation (MCMC)

- A **Polya-gamma data augmentation** for Bayesian logistic regression, introducing latent variables $\omega_i^{(j)} \sim PG(1, \alpha^{(j)} + \mathbf{x}_i^{(j)T} \boldsymbol{\beta}^{c_j})$
- Class allocation step involving prior penalization easily adapt marginal sampling for DP process

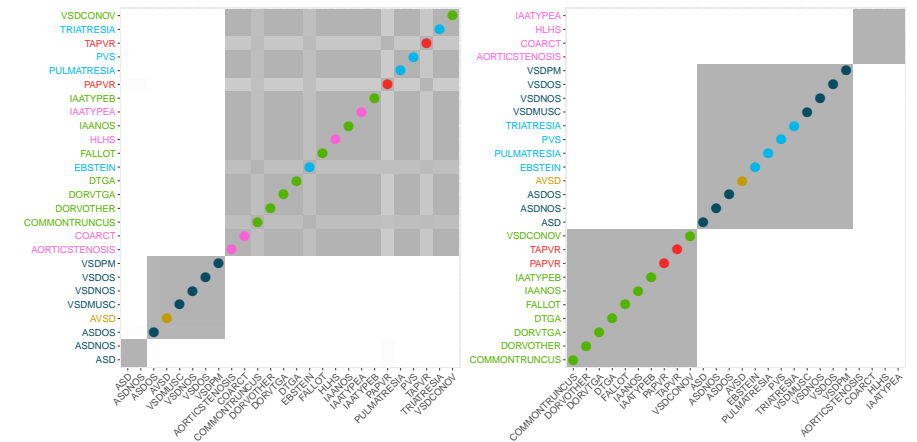
Clustering results



(a) $\psi = 0, VI(\hat{\mathbf{c}}, \mathbf{c}_0) = 2.43$

(b) $\psi = 40, VI(\hat{\mathbf{c}}, \mathbf{c}_0) = 1.78$

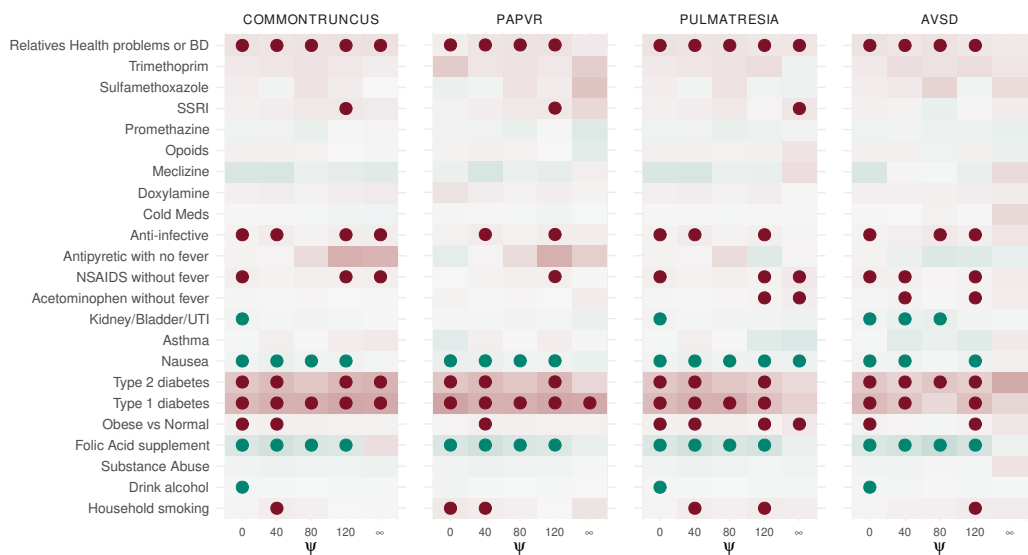
Clustering results



(c) $\psi = 80, VI(\hat{\mathbf{c}}, \mathbf{c}_0) = 1.65$


(d) $\psi = 120, VI(\hat{\mathbf{c}}, \mathbf{c}_0) = 0.86$

Exposure effects



Open questions

→ (NBDPS) data and prior guesses are actually more complicated

- Extensions of the CP process including
 - multiple prior guesses
 - characteristics of the partition (number of clusters)
- Investigate properties of the CP process
 - prediction rules for new observations
- Improve algorithms
 - prior calibration (scalability)
 - posterior sampling (local modes)
 - provide sampling methods via 






Thanks!

Centered Partition Processes: Informative Priors for Clustering.

Paganin S., Herring A. H., Olshan A. F. & Dunson B. D. (2021)

Bayesian Analysis (with Discussion)

References i

-  Hartigan, J. A. (1990).
Partition models
Commun. Statist. A **19**, 2745–2756.
-  Meila M. (2007).
Comparing clusterings - an information based distance.
J. of Mult. Analysis **98**, 873–895.
-  Müller, P., Quintana, F. & Rosner, G. L. (2011).
A Product Partition Model With Regression on Covariates.
J. Comput. Graph. Statist. **20**, 260–278.
-  Neal, R. M. (2000).
Markov chain sampling methods for Dirichlet process mixture models
J. Comput. Graph. Statist. **9**, 249–265.
-  Park, J.-H. & Dunson, D. B. (2010).
Bayesian Generalize Product Partition Models.
Stat. Sin. **20**, 1203–1226.



Rodriguez, A. & David B. D. (2011).

Nonparametric Bayesian models through probit stick-breaking processes

Bayesian analysis (Online) 6.1.



Stam, A.J. (1983).

Generation of a random partition of a finite set by an urn model

J. of Comb. Theory, Series A **35**, 231–240.
