







Model-based clustering of categorical data based on the Hamming distance

Lucia Paci

Department of Statistical Sciences, Università Cattolica del Sacro Cuore, Milano



Raffaele Argiento University of Bergamo



Andrea Cremaschi IE University



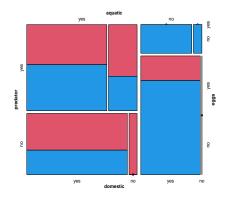
Edoardo Filippi-Mazzola Broadgate Advisers



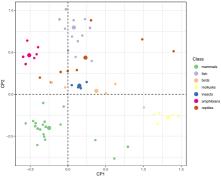
Benedetta Sabina Leone Politecnico di Milano

SISBAYES 2025

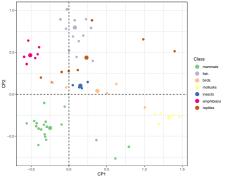
Department of Statistical Sciences, University of Padova 4-5 September 2025



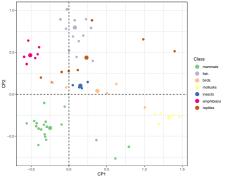
 Multivariate categorical data with no natural ordering (e.g., animal features)



- Multivariate categorical data with no natural ordering (e.g., animal features)
- Goal: identify a clustering structure in the data



- Multivariate categorical data with no natural ordering (e.g., animal features)
- Goal: identify a clustering structure in the data
- Heuristic framework: measure distances between observations
 - * K-modes (Huang, 1998)
 - * Hamming distance-vector algorithm
 (Zhang et al., 2006) finds clustering
 patterns using the Hamming distance
- Model-based clustering: **latent class analysis** (Goodman, 1974; Celeux and Govaert, 2015)



- Multivariate categorical data with no natural ordering (e.g., animal features)
- Goal: identify a clustering structure in the data
- Heuristic framework: measure distances between observations
 - * K-modes (Huang, 1998)
 - * Hamming distance-vector algorithm
 (Zhang et al., 2006) finds clustering
 patterns using the Hamming distance
- Model-based clustering: **latent class analysis** (Goodman, 1974; Celeux and Govaert, 2015)

Our approach

- ✓ Family of probability mass functions built upon the Hamming distance
- ✓ Model-based clustering based on a mixture of finite mixture of Hamming distributions
- ✓ Provide full posterior inference on the number of clusters and their structure

A bit of notation

- $X = (X_1, \dots, X_p)^{\top}$ is a vector of p nominal categorical variables, or **attributes**
- Each variable j, for $j=1,\ldots,p$, assumes m_j possible levels (categories) or **modalities**, over the finite set $A_j=\{a_{j1},\ldots,a_{jh},\ldots,a_{jm_j}\}$
- $x = (x_1, \dots, x_p)^{\top}$ is the vector of observed modalities
- Categorical sample space $\Omega_p = \{x = (x_1, \dots, x_p)^\top | x_1 \in A_1, \dots, x_p \in A_p\} = A_1 \times A_2 \times \dots \times A_p$
- Hamming distance: number of attributes whose modalities are different

Hamming distance between two points in Ω_n

$$d(\mathbf{x}_i, \mathbf{x}_h) = \sum_{j=1}^p 1 - \delta_{x_{ij}} \left(x_{hj} \right)$$

where
$$\delta_{x_{ij}}(x_{hj}) = \begin{cases} 1 \text{ if } x_{ij} = x_{hj} \\ 0 \text{ if } x_{ij} \neq x_{hj} \end{cases}$$

$$x_1 = [\clubsuit, \heartsuit, \bigstar]$$
 $x_2 = [\blacksquare, \heartsuit, \spadesuit]$
 $d(x_1, x_2) = 2$

Hamming distribution

- **center** parameter $\boldsymbol{c} = (c_1, \dots, c_p)^{\top} \in \Omega_p$
- scale parameter $\sigma = (\sigma_1, \dots, \sigma_p)^{\top}$, with $\sigma_j > 0, j = 1, \dots, p$

Hamming distribution

- **center** parameter $\boldsymbol{c} = (c_1, \dots, c_p)^{\top} \in \Omega_p$
- scale parameter $\sigma = (\sigma_1, \dots, \sigma_p)^{\top}$, with $\sigma_j > 0, j = 1, \dots, p$

Proposition 1

The function

$$p(\mathbf{x} \mid \mathbf{c}, \boldsymbol{\sigma}) = \prod_{j=1}^{p} \left(1 + \frac{m_j - 1}{\exp(1/\sigma_j)} \right)^{-1} \exp\left\{ -\frac{1 - \delta_{c_j}(x_j)}{\sigma_j} \right\}$$

is a probability mass function (p.m.f.) on Ω_p , i.e., $\sum_{x \in \Omega_n} p(x \mid c, \sigma) = 1$

A random vector $X = (X_1, \dots, X_p)$ with support Ω_p follows an **Hamming distribution** with center c and scale σ if its p.m.f. for $x \in \Omega_p$ is given by $p(x \mid c, \sigma)$ and we write

$$X \mid c, \sigma \sim \operatorname{Hamming}(c, \sigma)$$

Hamming distribution

- **center** parameter $\boldsymbol{c} = (c_1, \dots, c_p)^{\top} \in \Omega_p$
- scale parameter $\sigma = (\sigma_1, \dots, \sigma_p)^{\top}$, with $\sigma_j > 0, j = 1, \dots, p$

Proposition 1

The function

$$p(\mathbf{x} \mid \mathbf{c}, \boldsymbol{\sigma}) = \prod_{j=1}^{p} \left(1 + \frac{m_j - 1}{\exp(1/\sigma_j)} \right)^{-1} \exp\left\{ -\frac{1 - \delta_{c_j}(x_j)}{\sigma_j} \right\}$$

is a probability mass function (p.m.f.) on Ω_p , i.e., $\sum_{x \in \Omega_n} p(x \mid c, \sigma) = 1$

A random vector $X = (X_1, \dots, X_p)$ with support Ω_p follows an **Hamming distribution** with center c and scale σ if its p.m.f. for $x \in \Omega_p$ is given by $p(x \mid c, \sigma)$ and we write

$$X \mid c, \sigma \sim \operatorname{Hamming}(c, \sigma)$$

• When $\sigma_j = \sigma > 0$, $\forall j$, $p(\mathbf{x} \mid \mathbf{c}, \sigma) \propto \exp\left\{\frac{-d(\mathbf{c}, \mathbf{x})}{\sigma}\right\}$

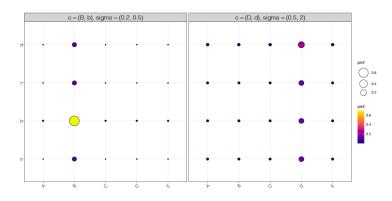


Figure: P.m.f. of p=2 categorical variables with different values of \mathbf{c} and $\boldsymbol{\sigma}$.

- The center represents the unique mode of the distribution
- The scale regulates the heterogeneity of the distribution (e.g., link with Gini's and Shannon's indexes)

Bayesian inference

- Sampling model: $X_i \mid c, \sigma \stackrel{iid}{\sim} \operatorname{Hamming}(c, \sigma), \qquad i = 1, \dots, n$
- Inference on c
 - ▶ Prior: $c_j \stackrel{iid}{\sim} U\{1, m_j\}$ $j = 1, \dots, p$
 - ► Full conditional probabilities: $p(c_j|\text{rest}) \propto \exp\left\{-\frac{n-\sum_{i=1}^n \delta_{c_j}(x_{ij})}{\sigma_j}\right\}$
- Inference on σ
 - ▶ Prior (Hypergeometric Inverse Gamma): $\sigma_j \mid u, v \stackrel{iid}{\sim} HIG(u, v)$ j = 1, ..., p, where

$$f(\sigma_{j} \mid u, v) = \frac{m_{j}^{(u+v)}(v+1)}{{}_{2}F_{1}(1, u+v; v+2; (m_{j}-1)/m_{j})} \left(1 + \frac{m_{j}-1}{\exp(1/\sigma_{j})}\right)^{-(u+v)} \exp\left(-\frac{v+1}{\sigma_{j}}\right) \frac{1}{\sigma_{j}^{2}}$$

where ${}_2F_1(\cdot,\cdot;\cdot;\cdot)$ is the hypergeometric function.

- Full conditional: $\sigma_j \mid \text{rest} \sim \text{HIG}(u^*, v^*)$, where $u^* = u + \sum_{i=1}^n \delta_{c_j}(x_{ij})$ and $v^* = v + n \sum_{i=1}^n \delta_{c_j}(x_{ij})$
- The marginal likelihood of the data is available in a closed-analytical form.

Hamming Mixture Model

HMM

$$h(\mathbf{x}_i \mid \mathbf{c}_1, \dots, \mathbf{c}_L, \boldsymbol{\sigma}_1, \dots, \boldsymbol{\sigma}_L, \boldsymbol{\pi}, L) = \sum_{l=1}^L \pi_l \, p(\mathbf{x}_i \mid \mathbf{c}_l, \boldsymbol{\sigma}_l)$$

- L: number of components
- π_l: mixing weight probability that observation i belongs to component l satisfies ∑^L_{l=1} π_l = 1 and 0 < π_l < 1
- $p(x_i | c_l, \sigma_l)$: mixture kernel Hamming p.m.f. of component l
- $c_l = (c_{1l}, \dots, c_{pl})^{\top}$: component-specific center
- $\sigma_l = (\sigma_{1l}, \dots \sigma_{pl})^{\top}$: component-specific scale

How do we choose L?

L is fixed: fit the model with different values of L and select the "best" one by using
information criteria, e.g., BIC, DIC, ICL (Biernacki et al., 2010; Celeux and Govaert, 2015)

How do we choose L?

- L is fixed: fit the model with different values of L and select the "best" one by using
 information criteria, e.g., BIC, DIC, ICL (Biernacki et al., 2010; Celeux and Govaert, 2015)
- L is random: a transdimensional sampler is needed for posterior inference
 - ▶ popular algorithm (but challenging) is the reversible jump MCMC (Green, 1995)
 - ▶ recent (painless!) alternatives mixtures of finite mixtures exploiting the link between finite and infinite mixture: Chinese restaurant process sampler (Miller and Harrison, 2018), telescoping sampler (Frühwirth-Schnatter et al., 2021), blocked Gibbs sampler (Argiento and De Iorio, 2022)

Prior on mixing weights

- Following Argiento and De Iorio (2022), the HMM can be framed in a BNP setting:
 - * infinite number of latent components
 - * only a finite number is used to generate the observed data
- We assign a prior distribution on the mixing weights by normalization:

$$\pi_1 = \frac{W_1}{T}, \dots, \pi_L = \frac{W_L}{T}, \qquad T = \sum_{l=1}^L W_l$$

$$W_1, \ldots, W_L, \mid L, \gamma \stackrel{iid}{\sim} \text{Gamma}(\gamma, 1)$$

- \Rightarrow Equivalent to: $\pi_1, \ldots, \pi_L \mid L, \gamma \sim \text{Dirichlet}_L(\gamma, \ldots, \gamma)$
- We assume a random number of components $L \sim q_L$

Clustering assignment

- Latent allocation variables $z_1 \dots, z_n$, where $z_i \in \{1, \dots, L\}$
- Observation x_i belongs to component $l \Leftrightarrow z_i = l$
- $z_i \mid W_1, \dots, W_L \stackrel{iid}{\sim} \text{Multinomial}(W_1/T, \dots, W_L/T)$

Clustering assignment

- Latent allocation variables $z_1 \dots, z_n$, where $z_i \in \{1, \dots, L\}$
- Observation x_i belongs to component $l \Leftrightarrow z_i = l$
- $z_i \mid W_1, \dots, W_L \stackrel{iid}{\sim} \text{Multinomial}(W_1/T, \dots, W_L/T)$
- z_1^*, \dots, z_K^* $K \leq L$, unique values of allocations
- $\rho := \{C_1, \dots, C_K\}$: partition of K clusters induced by z_1^*, \dots, z_K^* , where $C_k = \{i : z_i = z_k^*\}$ and $n_k = |C_k|$, for $k = 1, \dots, K$
- Prior of ρ : exchangeable partition probability function (eppf; Pitman 1995)

$$p(\rho) = p(n_1, \dots, n_K) \propto \prod_{k=1}^K \frac{\Gamma(\gamma + n_k)}{\Gamma(\gamma)}$$

• The prior on K is also available in a closed analytical form

HMM wrap up

$$egin{aligned} oldsymbol{x}_i \mid z_i, oldsymbol{c}_1, \dots, oldsymbol{c}_L, oldsymbol{\sigma}_L, \dots, oldsymbol{c}_L, oldsymbol{c}_L, \dots, oldsymbol{c$$

HMM wrap up

$$egin{aligned} x_i \mid z_i, c_1, \dots, c_L, \sigma_1, \dots, \sigma_L, L & \stackrel{ind}{\sim} & \operatorname{Hamming}(x_i \mid c_{z_i}, \sigma_{z_i}) & i = 1, \dots, n \\ & z_i \mid W_1, \dots, W_L, L & \stackrel{iid}{\sim} & \operatorname{Multinomial}(W_1/T, \dots, W_L/T) & i = 1, \dots, n \\ & W_l \mid L, \gamma & \stackrel{iid}{\sim} & \operatorname{Gamma}(\gamma, 1) & l = 1, \dots, L \\ & c_l \mid L & \stackrel{iid}{\sim} & \operatorname{U}\{1, m_j\} & l = 1, \dots, L \\ & \sigma_l \mid L & \stackrel{iid}{\sim} & \operatorname{HIG}(u, v) & l = 1, \dots, L \\ & L \mid \Lambda & \sim \operatorname{Poi_0}(\Lambda) & & & & \\ & \Lambda & \sim \operatorname{Gamma}(a, b) & & & & \end{aligned}$$

- Connections to Latent Class Models: the HMM is a novel parametrization of the parsimonious LCM (Celeux and Govaert, 2015) with two main benefits:
 - ✓ more straightforward interpretation of the parameters
 - \checkmark random L, so full posterior inference on the number of clusters and their structure

HMM wrap up

$$egin{aligned} x_i \mid z_i, c_1, \dots, c_L, \sigma_1, \dots, \sigma_L, L & \stackrel{ind}{\sim} & \operatorname{Hamming}(x_i \mid c_{z_i}, \sigma_{z_i}) & i = 1, \dots, n \\ & z_i \mid W_1, \dots, W_L, L & \stackrel{iid}{\sim} & \operatorname{Multinomial}(W_1/T, \dots, W_L/T) & i = 1, \dots, n \\ & W_l \mid L, \gamma & \stackrel{iid}{\sim} & \operatorname{Gamma}(\gamma, 1) & l = 1, \dots, L \\ & c_l \mid L & \stackrel{iid}{\sim} & \operatorname{U}\{1, m_j\} & l = 1, \dots, L \\ & \sigma_l \mid L & \stackrel{iid}{\sim} & \operatorname{HIG}(u, v) & l = 1, \dots, L \\ & L \mid \Lambda & \sim \operatorname{Poi}_0(\Lambda) & \\ & \Lambda & \sim \operatorname{Gamma}(a, b) & \\ & & & & & & & \\ & & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & \\ & & & & \\ & & & \\ & & & \\ & & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & \\ & & & \\ & &$$

- Connections to Latent Class Models: the HMM is a novel parametrization of the parsimonious LCM (Celeux and Govaert, 2015) with two main benefits:
 - ✓ more straightforward interpretation of the parameters
 - \checkmark random L, so full posterior inference on the number of clusters and their structure
- Posterior sampling
 - Blocked Gibbs sampler, a conditional algorithm of Argiento and De Iorio (2022)
 - Transdimensional moves which are automatic and implied by the prior process
 - Separate sampling of the weights and parameters corresponding to the allocated vs non-allocated components

USPS data analysis

- Handwritten digits from the US postal services (a subset of the USPS data from UCI machine learning repository)
- 1,756 images of the digits 3, 5 and 8 which are the most difficult digits to discriminate
- Each digit is a 16×16 image of m = 6 levels of gray, i.e., $X_i \in \{ \Box \ , \ \Box \ \}$, represented as a p = 256 dimensional vector

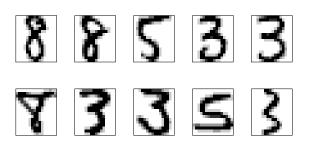


Figure: A sample of handwritten digits.

X Latent class model with fixed K: no way to select the number of clusters

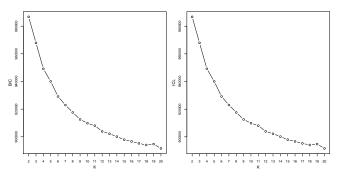


Figure: Values of the information criteria for fixed values of K; results obtained from the R package Rmixmod.

 \checkmark Hamming mixture model with random number of components L

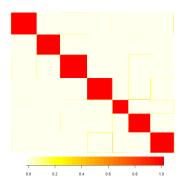


Figure: Posterior similarity matrix.

 \checkmark Hamming mixture model with random number of components L

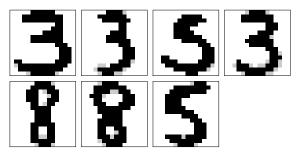
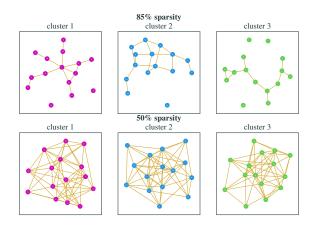
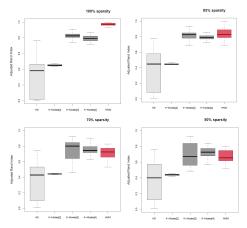


Figure: Posterior mode of the center parameters for each cluster.

- $n_k = 75, k = 1, ..., K = 3$, and $m_j = 4$ with j = 1, ..., p = 15
- Generating categorical data under a graphical modeling through Gaussian latent variables (Castelletti et al., 2024)
- Within each cluster, variables dependence structure based on a network with decreasing levels of sparsity, i.e., increasing strength of association among the variables



- $n_k = 75, k = 1, ..., K = 3$, and $m_j = 4$ with j = 1, ..., p = 15
- Generating categorical data under a graphical modeling through Gaussian latent variables (Castelletti et al., 2024)
- Within each cluster, variables dependence structure based on a network with decreasing levels of sparsity, i.e., increasing strength of association among the variables



- $n_k = 75, k = 1, \dots, K = 3$, and $m_j = 4$ with $j = 1, \dots, p = 15$
- Generating categorical data under a graphical modeling through Gaussian latent variables (Castelletti et al., 2024)
- Within each cluster, variables dependence structure based on a network with decreasing levels of sparsity, i.e., increasing strength of association among the variables

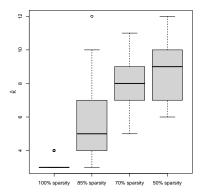


Figure: Estimated *K* under the HMM.

• Accounting for the dependence among the *p* categorical variables also within the clusters

- Accounting for the dependence among the p categorical variables also within the clusters
- Consider a mixture of HMMs: enriched HMM
- Enriched priors (Consonni and Veronese, 2001) for BNP: enriched Dirichlet (Wade et al., 2011), enriched Pitman—Yor process (Rigon et al., 2025), enriched Norm-IFPP (Franzolini et al., 2023)
- Connections to mixtures of LCMs (Malsiner-Walli et al., 2025)

Enriched HMM

$$h(\mathbf{x}_i \mid \mathbf{c}, \boldsymbol{\sigma}, \boldsymbol{\pi}, \boldsymbol{\eta}_1, \dots, \boldsymbol{\eta}_L, L, S_1, \dots, S_L) = \sum_{l=1}^{L} \pi_l \sum_{s=1}^{S_l} \eta_{ls} \, p\left(\mathbf{x}_i \mid \mathbf{c}_{ls}, \boldsymbol{\sigma}_{ls}\right)$$

- L: number of outer components
- π_l : outer mixing weight
- S_l : number of inner components in outer component l
- η_{ls} : inner mixing weight
- $p(x_i \mid c_{ls}, \sigma_{ls})$: mixture kernel Hamming p.m.f.
- c_{ls} = (c_{1ls},..., c_{pls})[⊤]: component-specific center parameter of variable j in outer component l and inner component s
- σ_{ls} = (σ_{1ls},...σ_{pls})^T: component-specific scale parameter of variable j in outer component l and inner component s
- Two-level clustering of observations (outer and inner clusters)

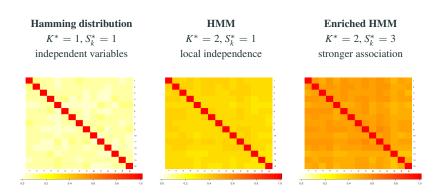


Figure: Cramer's V matrix of p=15 variables computed over n=450 data points with $m_j=4$ categories. Data simulated assuming $\sigma=0.4$.

- $n_k = 75, k = 1, ..., K = 3$, and $m_j = 4$ with j = 1, ..., p = 15
- Generating categorical data under a graphical modeling through Gaussian latent variables (Castelletti et al., 2024)
- Within each cluster, variables dependence structure based on a network with decreasing levels of sparsity, i.e., increasing strength of association among the variables

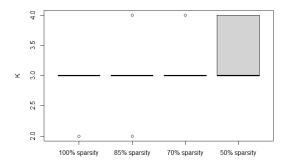


Figure: Estimated *K* under the Enriched HMM.

Summary

- A new probability mass function, based on the Hamming distance, to describe random vectors with support on a categorical space
- Conjugate Bayesian inference on the parameters of the Hamming distribution
- Hamming mixture model for clustering categorical data: finite mixture model with a random number of components
- Gibbs sampling strategy to provide full posterior inference of the cluster structure and the group-specific parameters and multiple imputation of missing values
- Theoretical results on model **identifiability** and **consistency** of the number of components
- Empirical analysis showed good accuracy in recovering the underlying clustering
- Enriched Hamming mixture model to overcome the local independence assumption
- Ongoing: study the properties of the enriched HMM (e.g., identifiability) and its link with a multivariate extension of the Hamming distribution

Thank you for your attention









References I

- Argiento, R. and De Iorio, M. (2022). Is infinity that far? A Bayesian nonparametric perspective of finite mixture models. Annals of statistics, 50(5):2641–2663.
- Biernacki, C., Celeux, G., and Govaert, G. (2010). Exact and Monte Carlo calculations of integrated likelihoods for the latent class model. *Journal of Statistical Planning and Inference*, 140(11):2991–3002.
- Castelletti, F., Consonni, G., and Vedova, M. L. D. (2024). Joint structure learning and causal effect estimation for categorical graphical models. *Biometrics*, 80.
- Celeux, G. and Govaert, G. (2015). Latent class models for categorical data. In Henning, C., Melia, M., Murtagh, F., and Rocci, R., editors, *Handbook of cluster analysis*. Chapman & Hall/CRC.
- Consonni, G. and Veronese, P. (2001). Conditionally reducible natural exponential families and enriched conjugate priors. Scandinavian Journal of Statistics, 28(3):377–406.
- Franzolini, B., Cremaschi, A., van den Boom, W., and Iorio, M. D. (2023). Bayesian clustering of multiple zero-inflated outcomes. *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences*, 381(2247):20220145.
- Frühwirth-Schnatter, S., Malsiner-Walli, G., and Grün, B. (2021). Generalized mixtures of finite mixtures and telescoping sampling. *Bayesian Analysis*, 16(4):1279–1307.
- Goodman, L. A. (1974). Exploratory latent structure analysis using both identifiable and unidentifiable models. *Biometrika*, 61(2):215–231.
- Green, P. J. (1995). Reversible jump Markov chain Monte Carlo computation and Bayesian model determination. *Biometrika*, 82(4):711–732.

References II

- Huang, Z. (1998). Extensions to the K-means algorithm for clustering large data sets with categorical values. Data Mining and Knowledge Discovery, 2(3):283–304.
- Malsiner-Walli, G., Grün, B., and Frühwirth-Schnatter, S. (2025). Without pain clustering categorical data using a Bayesian mixture of finite mixtures of latent class analysis models. Advances in Data Analysis and Classification.
- Miller, J. W. and Harrison, M. T. (2018). Mixture models with a prior on the number of components. *Journal of the American Statistical Association*, 113(521):340–356.
- Pitman, J. (1995). Exchangeable and partially exchangeable random partitions. Probability Theory and Related Fields, 102(2):145–158.
- Rigon, T., Petrone, S., and Scarpa, B. (2025). Enriched Pitman–Yor processes. Scandinavian Journal of Statistics. 52(2):631–657.
- Wade, S., Mongelluzzo, S., and Petrone, S. (2011). An enriched conjugate prior for Bayesian nonparametric inference. *Bayesian Analysis*, 6(3):359–385.
- Zhang, P., Wang, X., and Song, P. X.-K. (2006). Clustering categorical data based on distance vectors. Journal of the American Statistical Association, 101(473):355–367.