

# Exchangeable random measures and stick-breaking priors

Ramsés H. Mena  
O'Bayes, 2022

IIMAS-UNAM, México

(work with María F. Gil-Leyva, Theodoros Nicolieris, Antonio Lijoi and Igor Prünster)

September 9, 2022

# Talk plan

- 1 Stick-breaking priors
- 2 Exchangeable SB prior
- 3 Markov stick-breaking processes
- 4 ESB-Mixture model

# The basic setup

Random phenomena encoded in  $\{X_i\}_{i=1}^{\infty}$  r.v.'s

- Statistical learning requires stochastic dependence !

▷ Logical/physical independence  $\not\Rightarrow$  stochastic independence

so  $\mathbb{P}(X_{n+1} \in B \mid X_1, \dots, X_n) = \mathbb{P}(X_{n+1} \in B)$  not always a good idea!

▷ Under physical independence of obs. all we can assume is certain stochastic symmetry among  $\{X_i\}$

# The basic setup

Random phenomena encoded in  $\{X_i\}_{i=1}^{\infty}$  r.v.'s

- **Statistical learning requires stochastic dependence !**

- ▷ Logical/physical independence  $\not\Rightarrow$  stochastic independence

so  $\mathbb{P}(X_{n+1} \in B \mid X_1, \dots, X_n) = \mathbb{P}(X_{n+1} \in B)$  not always a good idea!

- ▷ Under physical independence of obs. all we can assume is certain stochastic symmetry among  $\{X_i\}$

- **Exchangeability**

$$(X_1, \dots, X_n) \stackrel{d}{=} (X_{\pi(1)}, \dots, X_{\pi(n)}), \quad \forall n \geq 1$$

and for any permutation  $\pi$  of  $\{1, \dots, n\}$ .

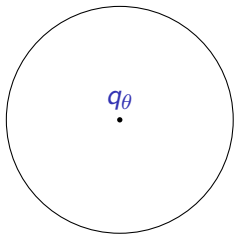
$\approx$  Distributional invariance under sampling order

# Exchangeability

$\mathbb{X}$ -valued  $\{X_i\}_{i=1}^{\infty}$  exchangeable sequence driven by  $P \sim Q$

- $Q(\cdot) = \delta_{q_\theta}(\cdot) \Rightarrow X_i$ 's are iid

$$\mathbb{P}(X_1 \in A_1, \dots, X_n \in A_n) = \int_{\mathcal{P}_{\mathbb{X}}} \prod_{i=1}^n P(A_i) \delta_{q_\theta}(dP) = \prod_{i=1}^n q_\theta(A_i)$$



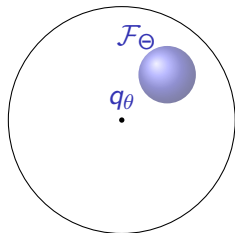
$\mathcal{P}_{\mathbb{X}}$  : Space of all distributions on  $\mathbb{X}$

## Exchangeability

$\mathbb{X}$ -valued  $\{X_i\}_{i=1}^{\infty}$  exchangeable sequence driven by  $P \sim Q$

- $Q(\mathcal{F}_{\Theta}) = 1 \Rightarrow$  Parametric family

$$\mathbb{P}(X_1 \in A_1, \dots, X_n \in A_n) = \int_{\mathcal{F}_{\Theta}} \prod_{i=1}^n \underbrace{F_{\theta}(A_i)}_{\text{Random uncertainty via param. model}} \overbrace{\pi_{\theta}(d\theta)}^{\text{Epistemic uncertainty}}$$



$\mathcal{P}_{\mathbb{X}}$  : Space of all distributions on  $\mathbb{X}$

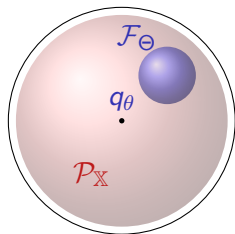
# Exchangeability

$\mathbb{X}$ -valued  $\{X_i\}_{i=1}^{\infty}$  exchangeable sequence driven by  $P \sim Q$

- $Q(P : d(P, \eta) < \varepsilon) > 0, \forall \eta \in \mathcal{P}_{\mathbb{X}} \text{ y } \varepsilon > 0 \Rightarrow \text{BNP}$

$$\mathbb{P}(X_1 \in A_1, \dots, X_n \in A_n) = \int_{\mathcal{P}_{\mathbb{X}}} \prod_{i=1}^n \underbrace{P(A_i)}_{Q(dP)}$$

Random and  
epistemic uncertainties  
in one stroke!

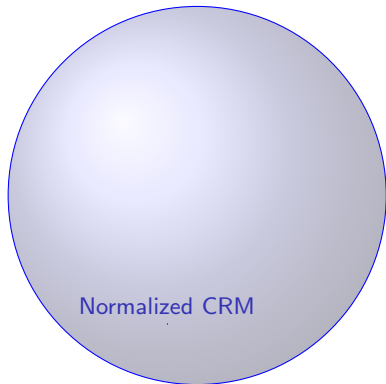


$\mathcal{P}_{\mathbb{X}}$  : Space of all distributions on  $\mathbb{X}$   
 ... or other infinite dimensional  
 sub-spaces of interest,  $\mathcal{P}_{\mathbb{X}}^d, \mathcal{P}_{\mathbb{X}}^c$ , etc.

How to construct suitable models for  $Q$  (nonparametric priors!)?

# Popular constructions of discrete BNP priors

Given a CRM  $\mu$  satisfying  $0 < \mu(\mathbb{X}) < \infty \Rightarrow P(\cdot) = \frac{\mu(\cdot)}{\mu(\mathbb{X})}$

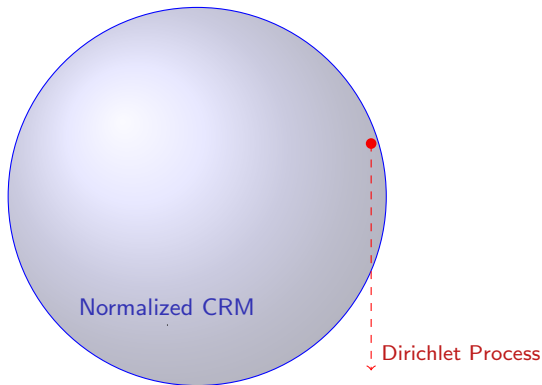




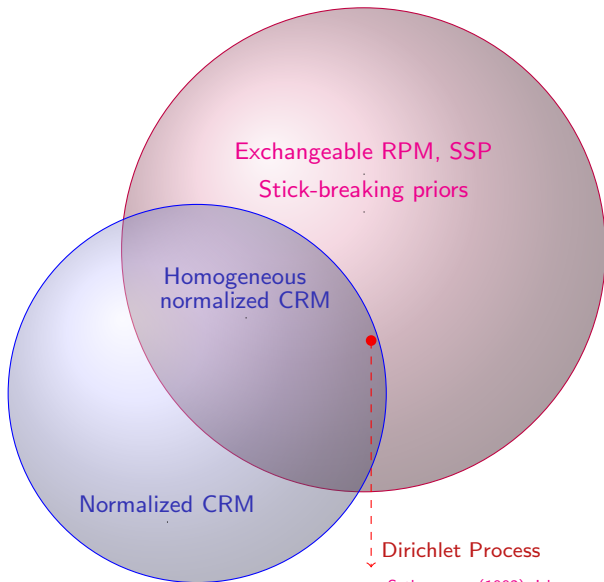
# Popular constructions of discrete BNP priors

Given a CRM  $\mu$  satisfying  $0 < \mu(\mathbb{X}) < \infty \Rightarrow P(\cdot) = \frac{\mu(\cdot)}{\mu(\mathbb{X})}$

If  $\mathbb{E}[\mu] = \nu(ds, dx) = s^{-1}e^{-s}\theta P_0(dx) \Rightarrow P \sim \mathcal{D}(\theta P_0)$



# Popular constructions of discrete BNP priors



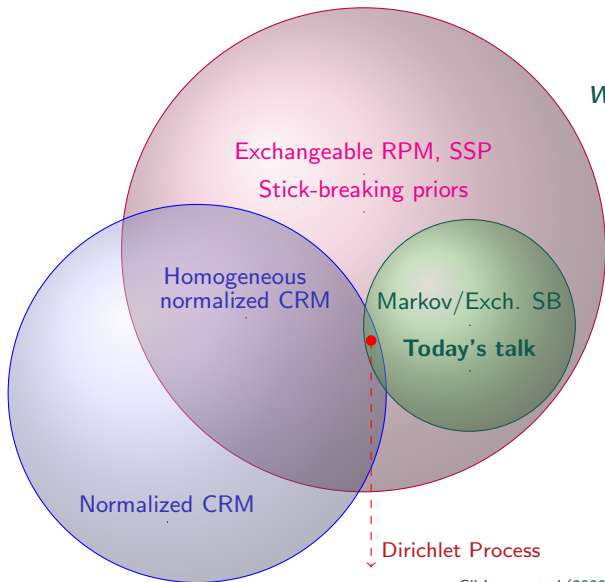
$$P = \sum_{i \geq 1} w_i \delta_{\xi_i}$$

$$\sum_{i \geq 1} w_i = 1$$

$$\xi_i \stackrel{\text{iid}}{\sim} P_0$$

$$(w_i)_{i \geq 1} \perp (\xi_i)_{i \geq 1}$$

# Popular constructions of discrete BNP priors



$$w_i = v_i \prod_{j=1}^{i-1} (1 - v_j)$$

Dependent  
 ↓  
 Length variables

## Stick breaking weights

$$P(B) = \sum_{i=1}^{\infty} w_i \delta_{\xi_i}(B), \quad B \in \mathcal{X}, \quad \sum_i w_i = 1$$

with  $w_i = v_i \prod_{j=1}^{i-1} (1 - v_j)$  and  $\xi_i \stackrel{\text{iid}}{\sim} P_0$

- **Full support if:** For every  $\varepsilon > 0$  there exist  $m \in \mathbb{N}$  such that

$$\mathbb{P}[v_1 < \varepsilon, \dots, v_m < \varepsilon] > 0$$

## Stick breaking weights

$$P(B) = \sum_{i=1}^{\infty} w_i \delta_{\xi_i}(B), \quad B \in \mathcal{X}, \quad \sum_i w_i = 1$$

with  $w_i = v_i \prod_{j=1}^{i-1} (1 - v_j)$  and  $\xi_i \stackrel{\text{iid}}{\sim} P_0$

- **Full support if:** For every  $\varepsilon > 0$  there exist  $m \in \mathbb{N}$  such that

$$\mathbb{P}[v_1 < \varepsilon, \dots, v_m < \varepsilon] > 0$$

- **Weights add up to one if**

$$\sum_{j \geq 1} w_j = 1 \quad \Leftrightarrow \quad \prod_{i=1}^j (1 - v_i) \xrightarrow{\text{a.s.}} 0 \quad \Leftrightarrow \quad \mathbb{E} \left[ \prod_{i=1}^j (1 - v_i) \right] \rightarrow 0$$

## Some SB representations of well-known RPMs

- Dirichlet process:  $v_i \stackrel{\text{iid}}{\sim} \text{Be}(1, \beta)$
- Two parameter Poisson-Dirichlet:  $v_i \stackrel{\text{ind}}{\sim} \text{Be}(1 - \sigma, \beta + i\sigma)$
- $\sigma$ -stable Poisson-Kingman: dependent  $(v_i)_{i \geq 1}$  with

$$g(v_j \mid t, v_1, \dots, v_{j-1}) = \frac{\sigma(tz_j)^{-\sigma}}{\Gamma(1 - \sigma)f_\sigma(tz_j)} v_j^{-\sigma} f_\sigma(tz_j(1 - v_j))$$

where  $f_\sigma$  denotes the positive  $\sigma$  stable density function and  $z_j := \prod_{i=1}^{j-1} (1 - v_i)$  with  $z_1 = 1$ .

⋮

- Homogeneous NRMs... also dependent  $(v_i)_{i \geq 1}$  with more involved conditional distributions for  $v_1$  and  $v_j \mid v_{i-1}, \dots, v_1$

## BNP clustering structure

- Relies on “analytical” expressions of the **EPPF**, i.e.  $\pi$  s.t.

$$\mathbb{P}(\Pi(\mathbf{x}_{1:n}) = A) = \pi(n_1, \dots, n_k) = \sum_{(j_1, \dots, j_k)} \mathbb{E} \left[ \prod_{i=1}^k w_{j_i}^{n_i} \right]$$

with  $A = \{A_1, \dots, A_k\}$  a partition of  $\mathbf{x}_{1:n}$  and  $n_j := |A_j|$

⇒ Similar inference can be achieved **via allocation variables**, i.e.,  
Given  $\{x_i\}_{i \geq 1}$  exch. driven by a SSP  $\mu = \sum w_j \delta_{\xi_j}$ ,  $d_i = j$  iff  $x_i = \xi_j$

$$\mathbb{P}(\mathbf{d}_1 = d_1, \dots, \mathbf{d}_n = d_n) = \mathbb{E} \left[ \prod_{j=1}^k v_j^{r_j} (1 - v_j)^{t_j} \right]$$

with  $k := \max\{d_1, \dots, d_n\}$ ,  $r_j := \sum_{i=1}^n \mathbf{1}_{(d_i=j)}$  and  $t_j := \sum_{i=1}^n \mathbf{1}_{(d_i>j)}$ .

## Dirichlet process &amp; Geometric process

Independent / Fully dependent random lengths

## Dirichlet process

- $v_i \stackrel{\text{iid}}{\sim} \text{Be}(1, \beta)$
- $w_j = v_j \prod_{i=1}^{j-1} (1 - v_i)$
- $\mathbb{E}[w_1] > \mathbb{E}[w_2] > \dots$
- ★ Size-biased random order

## Geometric process

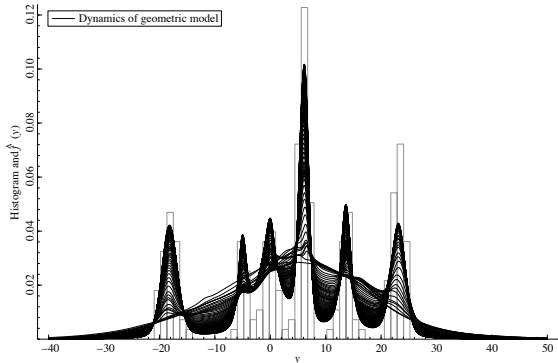
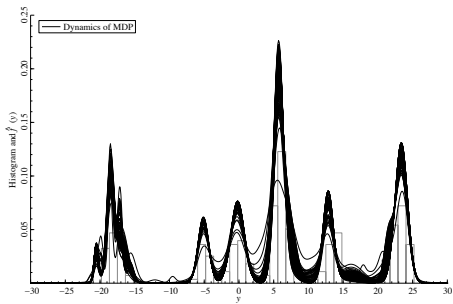
- $v_i = \lambda \sim \text{Be}(\alpha, \beta)$
- $w_j = \lambda(1 - \lambda)^{j-1}$
- $w_1 > w_2 > \dots$
- ★ Decreasing order

Both processes satisfy:

- $\sum_{j \geq 1} w_j = 1.$
- Have full support.
- Both are exchangeable!

★ We consider the general class of priors induced by exchangeable length variables





# Exchangeable stick-breaking processes

- $\text{ESB}(\nu, \mu_0)$

$$\mu = \sum_{j \geq 1} w_j \delta_{\xi_j} \xrightarrow{\text{atoms}} \xi_j \stackrel{\text{iid}}{\sim} \mu_0$$

with

$$w_j = \nu_j \prod_{i < j} (1 - \nu_i) \xrightarrow{\text{length var.}} (\nu_j) \text{ exchangeable seq. driven by } \nu$$

That is  $\nu_j | \nu \stackrel{\text{iid}}{\sim} \nu_0$  with  $\nu_0 := \mathbb{E}[\nu]$ .  $(\nu_j)$  are  $[0, 1]$ -valued.

## Theorem

- $\nu(\{0\}) < 1$  a.s iff  $\sum_{j \geq 1} w_j = 1$  a.s.
  - ▷ If  $\nu_0(\{0\}) = 0$  then  $\sum_{j \geq 1} w_j = 1$  a.s.
- If there exists  $\epsilon > 0$  such that  $(0, \epsilon)$  is contained in the support of  $\nu_0$ , then  $\mu$  has full support.

## Convergence to Dirichlet and Geometric processes

ESB( $\nu, \mu_0$ )

$$\begin{array}{ccc}
 \text{DP}(\theta, \mu_0) & \xleftarrow{\text{d}} & \mu & \xrightarrow{\text{d}} & \text{GP}(\nu_0, \mu_0) \\
 \text{Be}(1, \theta) & \xleftarrow{\text{d}} & \nu & \xrightarrow{\text{d}} & \delta_\nu, \text{ with } \nu \sim \nu_0
 \end{array}$$

If  $\nu = \sum_{j \geq 1} p_j \delta_{u_j}$  is a SSP with  $\rho := \mathbb{P}[v_1 = v_2] = \sum_{j \geq 1} \mathbb{E}[p_j^2]$

$$\begin{array}{ccc}
 \nu_0 & \xleftarrow{\text{d}} & \nu & \xrightarrow{\text{d}} & \delta_\nu, \text{ with } \nu \sim \nu_0 \\
 0 & \xleftarrow{\text{d}} & \rho & \xrightarrow{\text{d}} & 1
 \end{array}$$

Take  $\nu_0 = \text{Be}(1, \theta)$ . If  $\nu \sim \text{DP}(\beta, \nu_0)$ ,  $\rho = \frac{1}{1+\beta}$ .

# Convergence to Dirichlet and Geometric processes

## Ordering of the weights

$$\mu = \sum_{j \geq 1} w_j \delta_{\xi_j} \stackrel{d}{=} \sum_{j \geq 1} w_{\rho(j)} \delta_{\xi_j}$$

- One usually work with the ordering of weights that is the most tractable.

### Size-biased DP weights

$$\mathbb{E}[\tilde{w}_1] > \mathbb{E}[\tilde{w}_2] > \dots$$

$$(\tilde{w}_j) \xleftarrow{d} (w_j)$$

$$\text{Be}(1, \theta) \xleftarrow{d} \nu$$

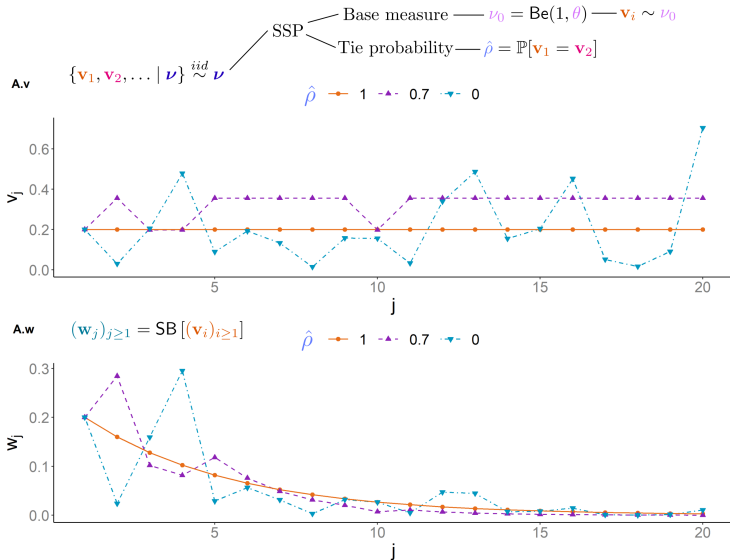
### Decreasing GP weights

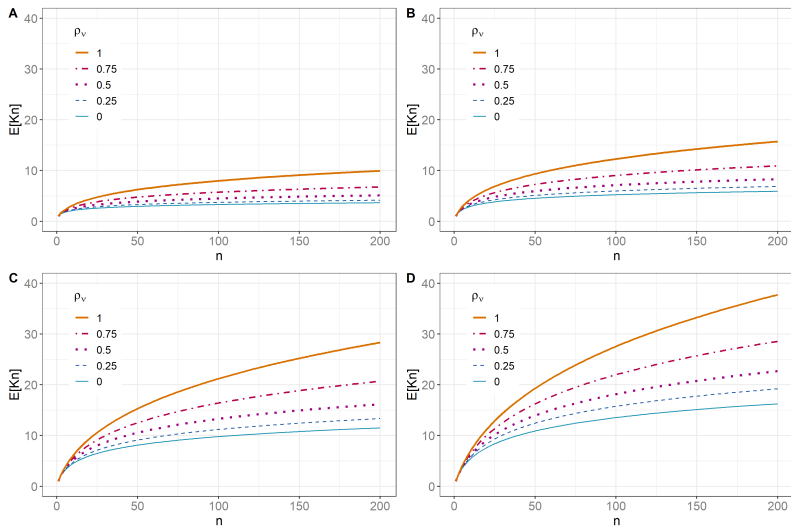
$$w_1^\downarrow > w_2^\downarrow > \dots$$

$$(w_j) \xrightarrow{d} (w_j^\downarrow)$$

$$\nu \xrightarrow{d} \delta_\nu, \text{ with } \nu \sim \nu_0$$

## Stick-breaking processes driven by SSP



Asymptotic behavior of  $E[K_n]$ 

$\{A, B, C, D\}$  corresponds to  $\theta = \{0.5, 1, 2.5, 4\}$

# Plan

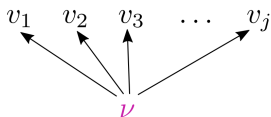
- 1 Stick-breaking priors
- 2 Exchangeable SB prior
- 3 Markov stick-breaking processes**
- 4 ESB-Mixture model

## Exchangeable vs Markov stick-breaking process

$$\mu = \sum_{j \geq 1} w_j \delta_{\xi_j} \dashrightarrow \xi_j \stackrel{iid}{\sim} \mu_0$$

$$w_j = v_j \prod_{i < j} (1 - v_i)$$

Exchangeable



$v$  takes values in  $\mathcal{P}([0, 1])$

$$\mathbb{E}[v] = \nu_0$$

$$v_j \sim \nu_0$$

$$\mu \sim \text{ESB}(v, \mu_0)$$

Markov (stationary)

$$v_1 \xrightarrow{\nu} v_2 \xrightarrow{\nu} v_3 \xrightarrow{\nu} \dots \xrightarrow{\nu} v_j$$

For each  $u \in [0, 1]$ ,  $\nu(u, \cdot) \in \mathcal{P}([0, 1])$

$$\mathbb{P}[v_{j+1} \in \cdot \mid v_j] = \nu(v_j, \cdot)$$

$$\nu_0(B) = \int \nu(u; B) \nu_0(du), \quad B \in \mathcal{B}([0, 1])$$

$$v_1 \sim \nu_0 \quad \Rightarrow \quad v_j \sim \nu_0$$

$$\mu \sim \text{MSB}(v, \mu_0)$$



# Exchangeable vs. Markov stick-breaking process

## Theorem (Exchangeable)

- $\nu(\{0\}) < 1$  a.s iff  $\sum_{j \geq 1} w_j = 1$  a.s.
  - \* If  $\nu_0(\{0\}) = 0$  then  $\sum_{j \geq 1} w_j = 1$  a.s.
- If there exists  $\epsilon > 0$  such that  $(0, \epsilon)$  is contained in the support of  $\nu_0$ , then  $\mu$  has full support.

## Theorem (Stationary Markov)

- $\nu_0 \neq \delta_0$  iff  $\sum_{j \geq 1} w_j = 1$  a.s.
- If there exists  $\epsilon > 0$  such that  $(0, \epsilon)$  is contained in the support of  $\nu_0$ , and for each  $u \in (0, \epsilon)$ ,  $(0, \epsilon)$  is contained in the support of  $\nu(u, \cdot)$ , then  $\mu$  has full support.

# Convergence to Dirichlet and Geometric processes

Size-biased DP weights

$$\mathbb{E}[\tilde{w}_1] > \mathbb{E}[\tilde{w}_2] > \dots$$

$$v_j \stackrel{iid}{\sim} \text{Be}(1, \theta)$$

$$(\tilde{w}_j)$$

$$\xleftarrow[\text{Be}(1, \theta) \leftarrow \nu]{d}$$

$$(w_j)$$

$$\xrightarrow[\nu \rightsquigarrow \delta_\bullet]{d} (w_j^\downarrow)$$

$$v_j \prod_{i < j} (1 - v_i)$$

$$v_j | \nu \stackrel{iid}{\sim} \nu$$

Decreasing GP weights

$$w_1^\downarrow > w_2^\downarrow > \dots$$

$$v_j = v \sim \nu_0$$

MSB( $\nu, \mu_0$ )

$$\text{DP}(\theta, \mu_0)$$

$$\xleftarrow[\text{Be}(1, \theta) \leftarrow \nu]{d}$$

$$\mu$$

$$\xrightarrow[\nu \rightsquigarrow \delta_\bullet]{d}$$

$$\text{GP}(\nu_0, \mu_0)$$

$$\nu_n \rightsquigarrow \varphi$$

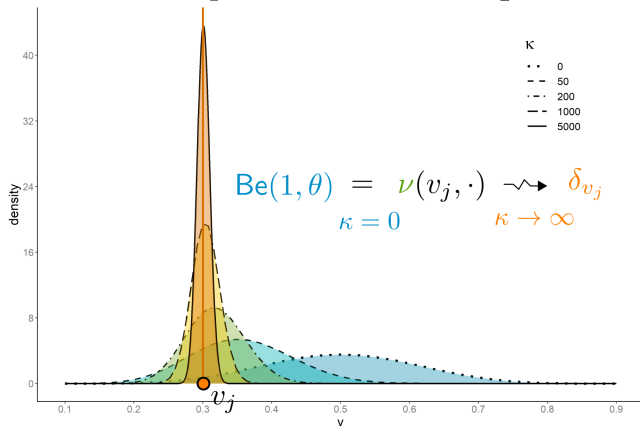
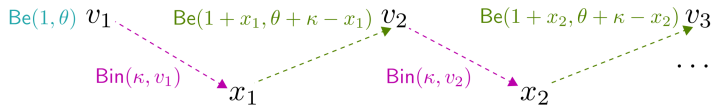


$$\nu_n(u_n, \cdot) \xrightarrow{w} \varphi(u, \cdot)$$

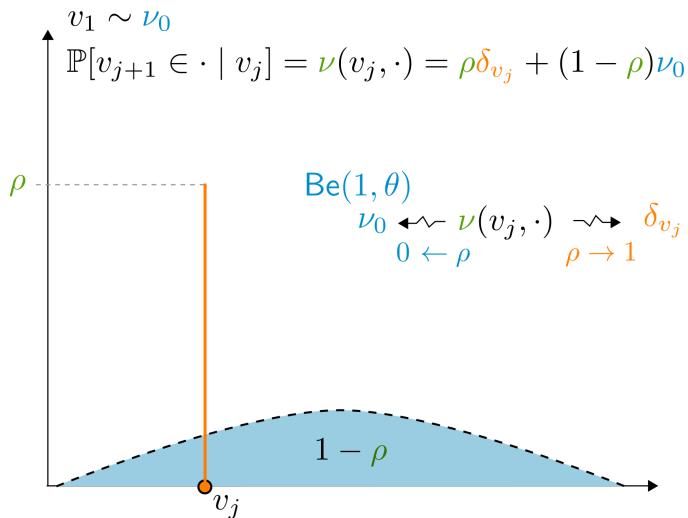
$$\forall u_n \rightarrow u \text{ in } [0, 1]$$

## Beta-Binomial transition

$$\nu(v_j, \cdot) = \sum_{x=0}^{\kappa} \text{Be}(\cdot \mid 1+x, \theta + \kappa - x) \text{Bin}(x \mid \kappa, v_j)$$



## Spike and slab transition



## Decreasing probability

If  $\nu$  is a spike and slab transition

$$\mathbb{P}[w_{j+1} \leq w_j] = \rho + (1 - \rho) \mathbb{E}[\vec{\nu}_0^{\rightarrow}(c(\nu))]$$

If  $\nu_0 = \text{Be}(1, \theta)$

$$\mathbb{P}[w_{j+1} \leq w_j] = 1 - \frac{{}_2F_1(1, 1; \theta + 2, 1/2)(1 - \rho)\theta}{2(\theta + 1)}$$

If  $\theta = 1$

$$\mathbb{P}[w_{j+1} \leq w_j] = \rho + (1 - \rho) \log(2)$$

# Plan

- 1 Stick-breaking priors
- 2 Exchangeable SB prior
- 3 Markov stick-breaking processes
- 4 ESB-Mixture model**

## ESB-Mixture model

$$\mu \sim \text{ESB}(\nu, \mu_0)$$

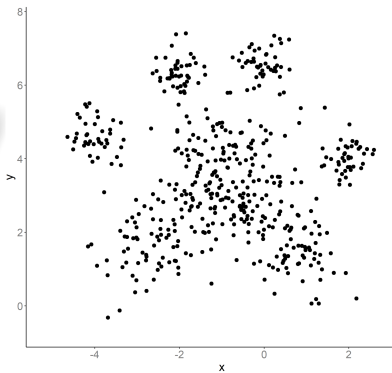
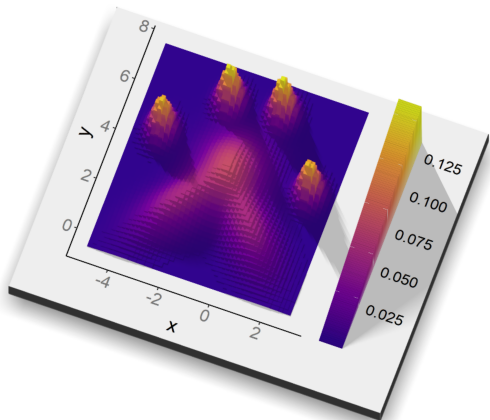
- Data modeled via  $y_i \mid \tilde{f} \stackrel{\text{iid}}{\sim} \tilde{f}$  for  $i = 1, 2, \dots$  with  $\tilde{f}$  a  $\mu$ -mixture, e.g.

$$\tilde{f}(y) = \int \text{N}(y \mid x) \mu(dx) = \sum_{j \geq 1} w_j \text{N}(y \mid \xi_j)$$

- ▷  $\mu_0 = \text{Normal} - \text{Inverse Wishart}$
- ▷  $\nu_0 = \text{Be}(1, \theta)$ ,  $\theta = 1$
- ▷  $\rho \sim \text{U}[0, 1]$

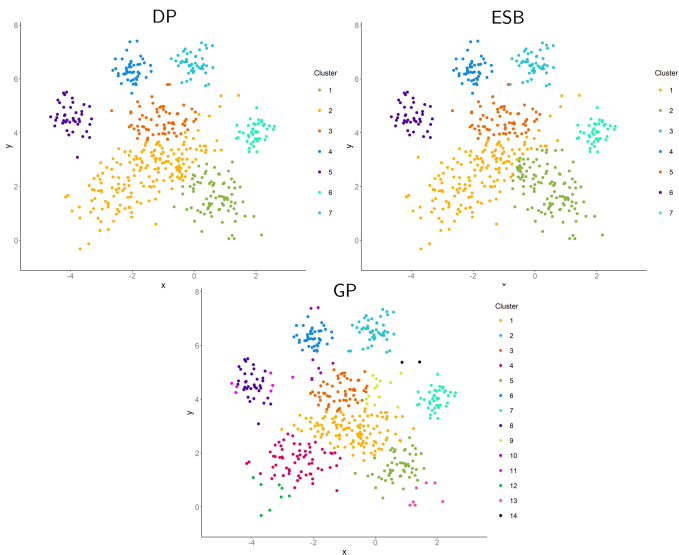
Gibbs sampler, e.g. Walker (2007), Kalli *et al.* (2011)

# Paw dataset

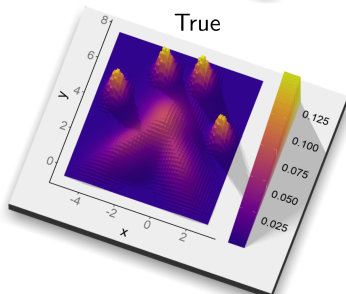
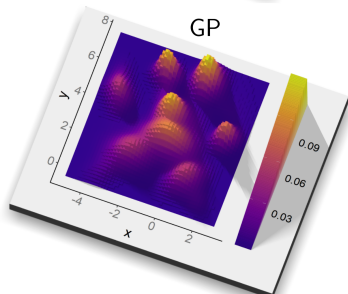
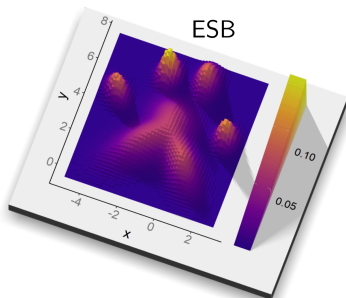
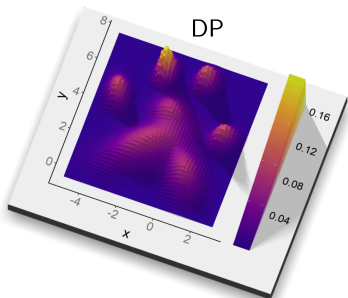




## Cluster estimation for Paw dataset



# Density estimation for Paw dataset



# References

- FAVARO, S., LIJOI, A., NAVA, C., NIPOTI, B., PRÜNSTER, I. AND TEH, Y. (2016). On the stick-breaking representation for homogeneous NRMIs, *Bayesian Analysis* 11: 697–724.
- FAVARO, S., LOMELI, M., NIPOTI, B. AND TEH, Y. (2014). On the stick-breaking representation of  $\sigma$ -stable Poisson-Kingman models, *Electronic Journal of Statistics* 8: 1063–1085.
- FERGUSON, T.S. (1973). A Bayesian analysis of some nonparametric problems. *Ann. Statist.* 1, 209–230.
- FUENTES-GARCÍA, R., MENA, R. H. AND WALKER, S. G. (2010). A new Bayesian nonparametric mixture model. *Communications in Statistics-Simulation and Computation*. 39, 669–682.
- GIL-LEYVA, M.F., MENA, R.H. AND NICOLERIS, T. (2020). Beta-Binomial stick-breaking non-parametric prior. *Electronic Journal of Statistics*. 14, 1479–1507.
- GIL-LEYVA, M.F. AND MENA, R.H. (2021). Stick-breaking processes with exchangeable length variables. *Journal of the American Statistical Association*. To appear.
- MENA, R. AND WALKER, S. G. (2009). On a construction of Markov models in continuous time, *METRON-International Journal of Statistics* LXVII: 303–323.
- PRÜNSTER, I. (2002). PhD in Mathematical Statistics. *Random probability measures derived from increasing additive processes and their application to Bayesian statistics*. Department of Mathematics, University of Pavia
- NIETO-BARAJAS, L. E. AND WALKER, S. G. (2002). Markov Beta and Gamma processes for modelling hazard rates, *Scandinavian Journal of Statistics* 29(3): 413–424.
- REGAZZINI, E., LIJOI, A. AND PRÜNSTER, I. (2003). Distributional results for means of random measures with independent increments. *Ann. Statist.*, 31, 560–585.
- SETHURAMAN, J. (1994). A constructive definition of Dirichlet priors. *Statist. Sinica* 4, 639–650.
- WALKER, S.G. (2007). Sampling the Dirichlet mixture model with slices. *Communications in Statistics*, 36, 45–54.



Gracias!