

# A review of Bayesian nonparametric regression through mixture models

Sonia Petrone

Bocconi University, Milano, Italy

with Sara Wade (University of Warwick, UK)  
and Michele Peruzzi (Bocconi University)

12-14 October 2016

Latent Variables, University of South Carolina, Columbia  
Conference in honor of Professor Jayaram Sethuraman

Increasing availability of large datasets,  
and developments in BNP

→ explosion of methods for Bayesian nonparametric regression.

However:

- still, less theory than for density estimation
- rich but fragmented literature
- little “ready-to-go” software

→ aim: an overview,  
having those three issues in mind

- focus on conditional density estimation and density regression (how  $f(y | x)$  varies with  $x$ ),
- through mixture models, based on Sethuraman's constructive definition of Dirichlet priors.
- Fragmented literature, how comparing?
  - frequentist properties (mostly asymptotics)
  - Bayesian? finite sample predictive properties.
    - Through a few examples, we underline how “details” of the prior can be overlooked, but do matter for prediction, and can be improved.
- software? R-packages?

- 1 Preliminaries
- 2 Random design: DP mixtures for  $f(x, y)$
- 3 Example 1 (improving prediction by restricted partition models)
- 4 Example 2 (Improving prediction by Enriched DP mixtures)
- 5 Fixed design: Dependent stick-breaking mixture models
- 6 discussion

- 1 Preliminaries
- 2 Random design: DP mixtures for  $f(x, y)$
- 3 Example 1 (improving prediction by restricted partition models)
- 4 Example 2 (Improving prediction by Enriched DP mixtures)
- 5 Fixed design: Dependent stick-breaking mixture models
- 6 discussion

# Density regression

$X$  predictor,  $p$ -dimensional;  $Y$  response.

\* Bayesian nonparametric (mean) regression:

$$: x \rightarrow m(x) = E(Y | x)$$

flexible model/prior on  $m(x)$  (basis expansions, Gaussian processes and splines (Wahba (1990), Denison, Holmes, Mallick, Smith (2002)), wavelets (Vidakovic, 2009), neural networks (Neal, 1996),...

\* Yet, the mean may be a too poor summary of the relationship between  $x$  and  $y$ .

$\Rightarrow$  median, quantile regression,.. density regression

$$: x \rightarrow f(y | x)$$

Limited literature on optimal estimators of  $f(y | x)$ . (Efromovich, Ann. Stat. 2007). How about Bayesian methods?

# Random or fixed design

- **regression: random design**

$(X_i, Y_i), i = 1, \dots, n$  are a random sample from  $f(x, y)$ .

Then, estimate the joint density  $f(x, y)$  and from this the conditional  $f(y | x) = \frac{f(x, y)}{f_x(x)}$ .

- **fixed design**

$x_1, \dots, x_n$  is a deterministic sequence. If predictor is  $x$ ,  $Y \sim f(y | x)$ . in fact,  $f_x(y)$ .

- **replicates of  $y$  values at a given  $x$**  ( $x$  typically categorical or ordinal), (e.g. ANOVA)
- **no replicates** ( $x$  typically continuous)  
still interest in  $f_x(y)$ : borrowing strength through smoothness conditions along  $x$ .

- 1 Preliminaries
- 2 Random design: DP mixtures for  $f(x, y)$
- 3 Example 1 (improving prediction by restricted partition models)
- 4 Example 2 (Improving prediction by Enriched DP mixtures)
- 5 Fixed design: Dependent stick-breaking mixture models
- 6 discussion



## Random design: DP mixtures for $f(x, y)$

$$(X_i, Y_i) \mid f \stackrel{iid}{\sim} f(x, y), \quad f \sim \text{prior prob. law}$$

## Random design: DP mixtures for $f(x, y)$

$$(X_i, Y_i) \mid f \stackrel{iid}{\sim} f(x, y), \quad f \sim \text{prior prob. law}$$

\* Model  $f$  as a mixture of kernels:

$$\begin{aligned} (X_i, Y_i) \mid G &\stackrel{iid}{\sim} f_G(x, y) = \int K(x, y \mid \theta) dG(\theta) \\ G &\sim DP(\alpha G_0) \end{aligned}$$

Usually,  $K(x, y \mid \theta) = N_{p+1}(x, y \mid \mu, \Sigma)$ , with  $\theta = (\mu, \Sigma)$ , and  $G_0(\theta)$  conjugate prior.

## Random design: DP mixtures for $f(x, y)$

$$(X_i, Y_i) | f \stackrel{iid}{\sim} f(x, y), \quad f \sim \text{prior prob. law}$$

\* Model  $f$  as a mixture of kernels:

$$(X_i, Y_i) | G \stackrel{iid}{\sim} f_G(x, y) = \int K(x, y | \theta) dG(\theta)$$
$$G \sim DP(\alpha G_0)$$

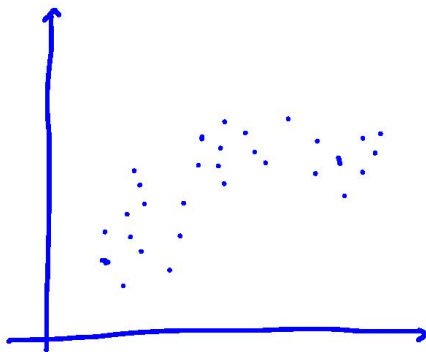
Usually,  $K(x, y | \theta) = N_{p+1}(x, y | \mu, \Sigma)$ , with  $\theta = (\mu, \Sigma)$ , and  $G_0(\theta)$  conjugate prior.

\* By Sethuraman construction, a.s.,  $G = \sum_{j=1}^{\infty} p_j \delta_{\theta_j^*}$ , where  $(p_j) \sim \text{stick-breaking}(\alpha)$  independent on  $\theta_j^* \stackrel{iid}{\sim} G_0$ . Then

$$f_G(x, y) = \sum_{j=1}^{\infty} w_j K(x, y | \theta_j^*).$$

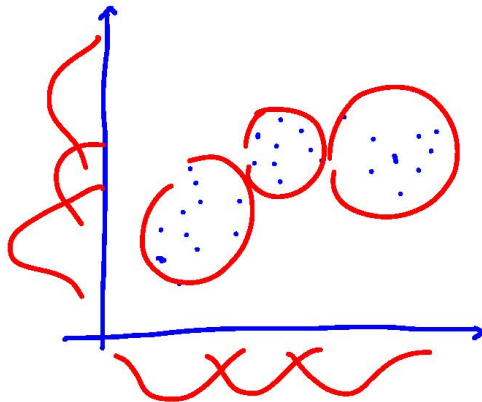
# mixture of Gaussian kernels

in picture



# mixture of Gaussian kernels

in picture



# Latent variable representation

The DP-mixture model is equivalently expressed as

$$\begin{aligned}(X_i, Y_i) | \theta_i &\stackrel{ind}{\sim} K(x, y | \theta_i) \\ \theta_i | G &\stackrel{iid}{\sim} G \\ G &\sim DP(\alpha G_0)\end{aligned}$$

Integrating the  $\theta_i$  out, one has back the countable mixture model

$$(X_i, Y_i) | G \stackrel{iid}{\sim} f_G(x, y) = \sum_{j=1}^{\infty} p_j K(x, y | \theta_j^*).$$

# Latent variable representation

The DP-mixture model is equivalently expressed as

$$\begin{aligned}(X_i, Y_i) | \theta_i &\stackrel{ind}{\sim} K(x, y | \theta_i) \\ \theta_i | G &\stackrel{iid}{\sim} G \\ G &\sim DP(\alpha G_0)\end{aligned}$$

Integrating the  $\theta_i$  out, one has back the countable mixture model

$$(X_i, Y_i) | G \stackrel{iid}{\sim} f_G(x, y) = \sum_{j=1}^{\infty} p_j K(x, y | \theta_j^*).$$

Then the conditional density  $f(y | x)$  is obtained as

$$f_G(y | x) = \frac{\sum_j p_j K(x, y | \theta_j^*)}{\sum_j p_j K(x | \theta_j^*)} = \sum_j p_j(x) K(y | x, \theta_j^*)$$

where  $p_j(x) = p_j K(x | \theta_j^*) / (\sum_{j'} p_{j'} K(x | \theta_{j'}^*))$ .

# Random partition

Since  $G$  is a.s. discrete, ties in a sample  $(\theta_1, \dots, \theta_n)$  from  $G$  have positive probability, so that

$(\theta_1, \dots, \theta_n)$  described by  $(\rho_n; \theta_1^*, \dots, \theta_{k(\rho_n)}^*)$

- a random partition  $\rho_n = (s_1, \dots, s_n)$
- the cluster-specific parameters  $\theta_j^*$

**Ex:** for  $n = 5$ ,  $\rho_n = (1, 1, 2, 2, 1)$  gives  $(\theta_1, \dots, \theta_n) = (\theta_1^*, \theta_1^*, \theta_2^*, \theta_2^*, \theta_1^*)$ ,  $k_n = 2$  two clusters of size  $n_1 = 3$ ,  $n_2 = 2$  resp., with cluster-specific parameters  $\theta_1^*, \theta_2^*$ .

- The DP induces a probability law of the random partition

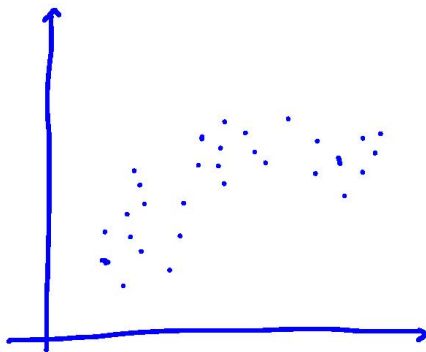
$$p(\rho_n) = \frac{\Gamma(\alpha)}{\Gamma(\alpha + n)} \alpha^{k_n} \prod_{j=1}^{k_n} \Gamma(n_j)$$

- Given the partition  $\rho_n$ , the cluster specific parameters  $\theta_j^*$  are i.i.d.  $\sim G_0$ .



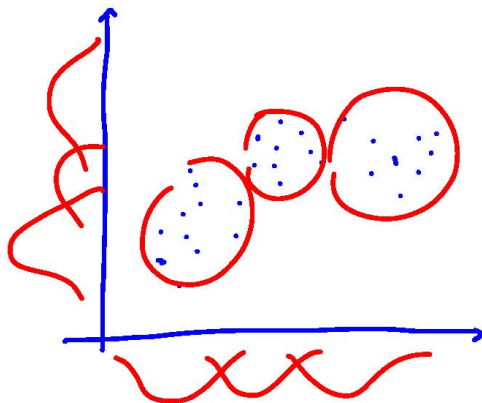
# mixture of Gaussian kernels

in picture

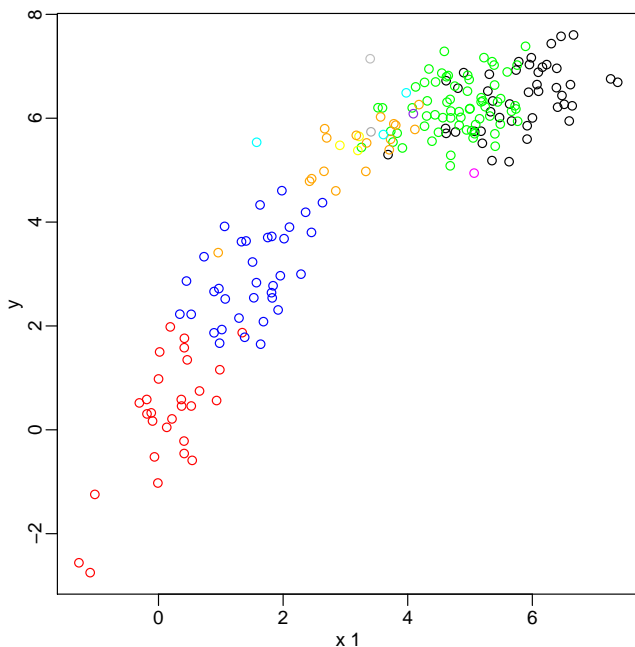


# mixture of Gaussian kernels

in picture



# mixture of Gaussian kernels



\* **posterior on  $\rho_n$**

$$\begin{aligned}
 p(\rho_n \mid x_{1:n}, y_{1:n}) &\propto p(\rho_n) \prod_{j=1}^{k_n} \int \prod_{i:(x_i, y_i) \in S_j} K(x_i, y_i \mid \theta_j^*) dP_0(\theta_j^*) \\
 &\propto p(\rho_n) \prod_{j=1}^{k_n} m(\{(x_i, y_i) \in C_j\} \mid \rho_n)
 \end{aligned}$$

prior times the independent marginal likelihoods in each clusters.

\* **posterior on  $\rho_n$**

$$\begin{aligned} p(\rho_n \mid x_{1:n}, y_{1:n}) &\propto p(\rho_n) \prod_{j=1}^{k_n} \int \prod_{i:(x_i, y_i) \in S_j} K(x_i, y_i \mid \theta_j^*) dP_0(\theta_j^*) \\ &\propto p(\rho_n) \prod_{j=1}^{k_n} m(\{(x_i, y_i) \in C_j\} \mid \rho_n) \end{aligned}$$

prior times the independent marginal likelihoods in each clusters.

\* **posterior on  $(\theta_j^*)$**  Given the partition, clusters are independent, and inference on  $\theta_j^*$  is based only on obs. in group  $C_j$

$$p(\theta_1^*, \dots, \theta_{k_n}^* \mid x_{1:n}, y_{1:n}, \rho_n) = \prod_{j=1}^{k_n} p(\theta_j^* \mid \rho_n, \{(x_i, y_i) \in C_j\})$$

# Predictive distribution (density estimate)

With quadratic loss, the Bayesian density estimate is the predictive density

$$(X_{n+1}, Y_{n+1}) \mid x_{1:n}, y_{1:n} \sim \hat{f}(x, y) = \mathbb{E}(f(x, y) \mid x_{1:n}, y_{1:n}).$$

\* Recalling that  $G \mid \theta_{1:n}, x_{1:n}, y_{1:n} \sim DP(\alpha G_0 + \sum_{i=1}^n \delta_{\theta_i})$ ,

$$\begin{aligned} \hat{f}(x, y) &= \mathbb{E}(\mathbb{E}(f_G(x, y) \mid \theta_{1:n}) \mid x_{1:n}, y_{1:n}) \\ &= \frac{\alpha}{\alpha + n} f_{G_0}(x, y) + \frac{n}{\alpha + n} \mathbb{E} \left( \sum_{i=1}^n \frac{K(x, y \mid \theta_i)}{n} \mid x_{1:n}, y_{1:n} \right) \end{aligned}$$

average of prior guess  $f_{G_0}$  and expectation of a **kernel estimate** with kernels centered at the  $\theta_j$ .

\* Recalling that  $(\theta_1, \dots, \theta_n) \leftrightarrow (\rho_n, \theta_1^*, \dots, \theta_k^*)$ , the joint density estimate is

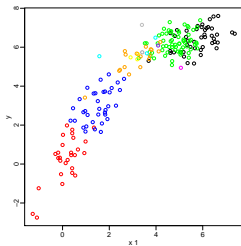
$$\hat{f}(x, y) = \frac{\alpha}{\alpha + n} f_{G_0}(x, y) + \sum_{\rho_n} \left( \sum_{j=1}^{k(\rho_n)} \frac{n_j(\rho_n)}{\alpha + n} f(x, y \mid (x_i, y_i) \in C_j(\rho_n)) \right) p(\rho_n \mid x_{1:n}, y_{1:n})$$

average of the prior guess  $f_{G_0}$ , and given the partition  $\rho_n$ , of the predictive densities in clusters  $C_j(\rho_n)$ .

From  $\hat{f}(x, y)$ , one can find an estimate of  $f(y \mid x)$ .

## 'clustering' and density estimate

The partition is not of main interest (mixture components just play the role of kernels), but  $p(\rho_n \mid x_{1:n}, y_{1:n})$  plays a crucial role.



Such role of the prior and posterior distribution of the random partition is often overlooked.



# Dirichlet process mixture models

A DPM is defined as

$$\begin{aligned} Y_i | x_i, \beta_i, \sigma_i^2 &\stackrel{\text{ind}}{\sim} N(\beta_i' x_i, \sigma_i^2), \\ \beta_i, \sigma_i^2 | G &\stackrel{\text{iid}}{\sim} G, \\ G &\sim DP(\alpha G_0). \end{aligned}$$

base measure  $G_0$ : conjugate Normal-Inverse Gamma prior  $(\beta_0, C^{-1}, a, b)$ .

$G$  is a.s. discrete, and integrating the parameters out

$$Y_i | x_i, G \stackrel{\text{ind}}{\sim} \sum_{j=1}^{\infty} w_j N(\beta_j^* x_i, \sigma_j^2(x_i)),$$

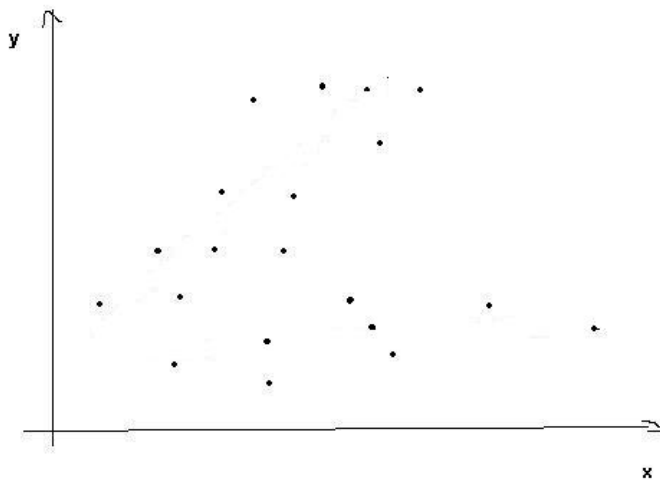
thus

$$m(x) = E[Y | x, w, G] = \sum_{j=1}^{\infty} w_j \beta_j^* x,$$

a mixture of cluster-specific linear regressions

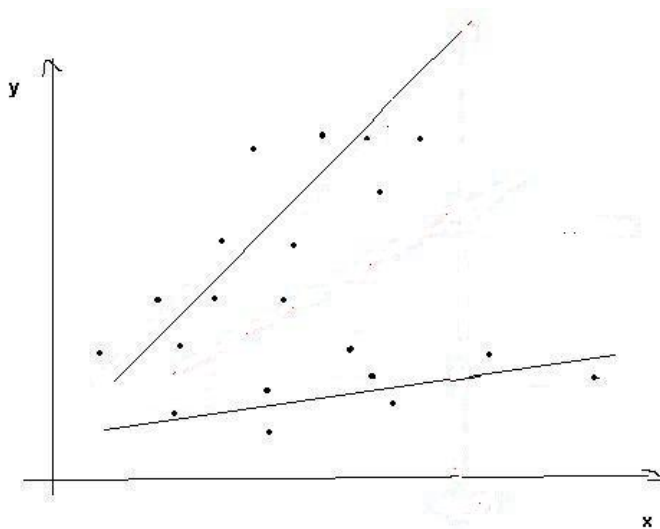
# DP mixture of linear regression

more an exploratory tool



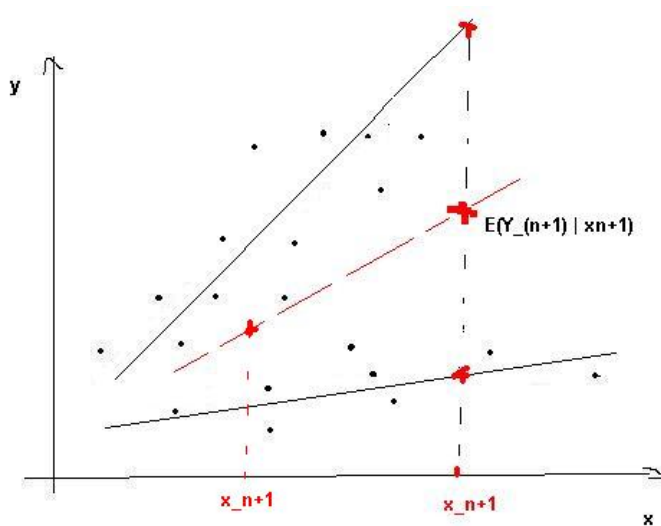
# DP mixture of linear regression

more an exploratory tool



# DP mixture of linear regression

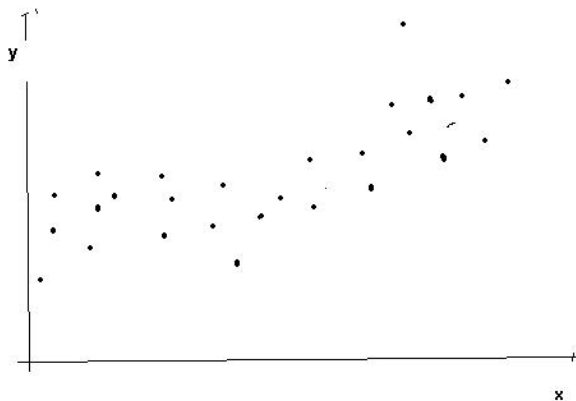
but prediction generally uninformative



- 1 Preliminaries
- 2 Random design: DP mixtures for  $f(x, y)$
- 3 Example 1 (improving prediction by restricted partition models)**
- 4 Example 2 (Improving prediction by Enriched DP mixtures)
- 5 Fixed design: Dependent stick-breaking mixture models
- 6 discussion

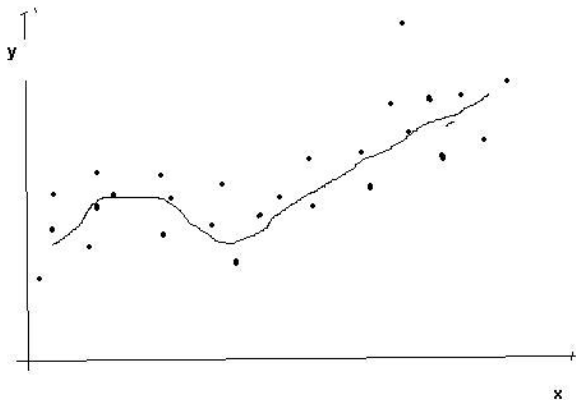
# Example

Consider a simple example,  $x$  and  $y$  univariate



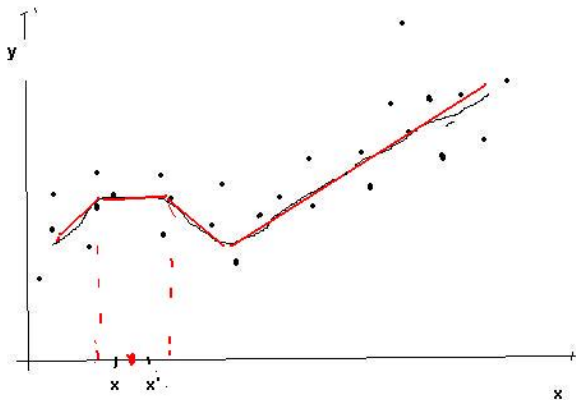
# BNP curve fitting

DP and 'clustering' is used for flexible regression



# BNP curve fitting

DP and 'clustering' is used for flexible regression





Here, the partition should depend on  $x$ , and if  $x \approx x'$ , we expect that they belong to the same cluster.

Here, the partition should depend on  $x$ , and if  $x \approx x'$ , we expect that they belong to the same cluster.

*joint DP and DDP mixture models* go in this direction: the random partition depends on covariates  $x$ , and the cluster allocation of  $Y_{n+1}$  depends on  $x_{n+1}$ .

However, we still unnecessarily give positive probability mass to 'bad' partitions, which still affect prediction.

# Restricted DP mixtures

- We want to incorporate the information that clusters should be based on **proximity** of  $x$ .
- However, because the **total number of partitions** is so **large**, placing higher prior mass on desirable partitions (joint DPM) is not enough to ensure:
  - prominence of these partitions in the posterior,
  - sufficiently small posterior mass for undesirable partitions.
- The only way to ensure this is to place **zero** mass on undesirable partitions in the prior (Wade, Walker, P., 2014).

# the space of partitions is huge!

- Focus on a particular case:  $x$  is one-dimensional and continuous.
- Partitions should be based on the natural ordering of  $x$ .
- Number of ways to partition the  $n$  subjects:

$$B_n = \sum_{k=1}^n S_{n,k}, \text{ a Bell number}$$

$$S_{n,k} = \frac{1}{k!} \sum_{j=0}^k (-1)^j \binom{k}{j} (k-j)^n \text{ a Stirling number of the second kind.}$$

- Under a ordering constraint, number of ways to partition the  $n$  subjects:

$$\sum_{k=1}^n \binom{n-1}{k-1} = 2^{n-1}.$$

- Example: if  $n=10$ ,  $B_{10} = 115,975$  and  $2^{n-1} = 512$ , 0.44% of the total partitions, and if  $n = 100$ , the percentage of partitions under this constraint is less than  $10^{-83}\%$  of the total partitions.

- We define a covariate-dependent random partition model that both **removes undesirable partitions** and **retains certain properties of the random partition model induced by the DP**.
- **Computations** for  $p(\rho_n | y, x)$  use reversible jump MCMC (Fuentes-Garcia et al., 2010). The smaller parameter space  $\rightarrow$  **faster computations and better mixing!**
- By removing bad partition, **we greatly improve prediction**.

(Wade, Walker, Petrone, 2014)

# can we formalize the gain of information?

**Theory.** Results on frequentist asymptotic properties for multivariate density estimation, and some results for regression and conditional density estimation (Wu & Ghosal (2008; 2010); Tokdar (2011), Norets & Pelenis (2012), Shen, Tokdar & Ghosal (2013), Canale & Dunson (2015),...)

But, how about (Bayesian) **finite sample properties and predictive performance?**

- 1 Preliminaries
- 2 Random design: DP mixtures for  $f(x, y)$
- 3 Example 1 (improving prediction by restricted partition models)
- 4 Example 2 (Improving prediction by Enriched DP mixtures)**
- 5 Fixed design: Dependent stick-breaking mixture models
- 6 discussion

## Problem: anisotropic case

$X$  continuous predictors; Gaussian kernels.

The DP mixture of multivariate Gaussian distributions uses joint clusters to fit the density  $f(x, y)$ .



## Problem: anisotropic case

$X$  continuous predictors; Gaussian kernels.

The DP mixture of multivariate Gaussian distributions uses joint clusters to fit the density  $f(x, y)$ .

**BUT**, the conditional density  $f_{y|x}$  and the marginal density  $f_x$  might have different smoothness; in regression, typically  $f_{y|x}$  is smoother than  $f_x$ .

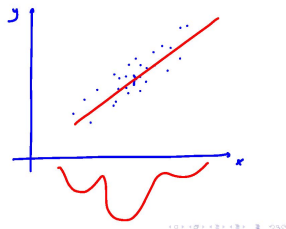
# Problem: anisotropic case

$X$  continuous predictors; Gaussian kernels.

The DP mixture of multivariate Gaussian distributions uses joint clusters to fit the density  $f(x, y)$ .

**BUT**, the conditional density  $f_{y|x}$  and the marginal density  $f_x$  might have different smoothness; in regression, typically  $f_{y|x}$  is smoother than  $f_x$ .

in picture



here, many small clusters (kernels) are needed to fit the  $f_x$  density, while much fewer kernels would suffice for  $f_{y|x}$ .

If the dimension of  $x$  is large, the likelihood is dominated by the  $x$  component and many small clusters are suggested by the posterior on  $\rho_n$ . This impoverishes the performance of the model.

This undesirable behavior does not seem to vanish with increasing sample size.

If  $f_x(x)$  requires many clusters, the unappealing behaviour of the random partition could be reflected in worse convergence rates. Efromovich [2007] shows that if the conditional density is smoother than the joint, it can be estimated at a faster rate.

Thus, improving inference on the random partition to take into account the different degree of smoothness of  $f_x$  and  $f_{y|x}$  is crucial.

Consider the Dirichlet mixture of Gaussian kernels

$$(X_i, Y_i) | G \sim \int N_{p+1}(\mu, \Sigma) dG(\mu, \Sigma), \quad G \sim DP(\alpha G_0).$$

The base measure of the DP,  $G_0(\mu, \Sigma)$ , is usually Normal-Inv Wishart. **BUT, this conjugate prior is restrictive if  $p$  is large**

# Improving prediction

- Write the kernels as

$$N_{p+1}(x, y | \mu, \Sigma) = N_p(x | \mu_x, \Sigma_x) N(y | x' \beta, \sigma_{y|x}^2)$$

and use simple spherical  $x$ -kernels (Shahbaba and Neal, 2009, Hannah et al., 2011). Thus,

$$f_x(x | P) = \sum_{j=1}^{\infty} w_j N_p(x | \mu_{xj}^*, \begin{pmatrix} \sigma_{x1,j}^{2*} & 0 & \cdots & 0 \\ 0 & \sigma_{x2,j}^{2*} & 0 & 0 \\ 0 & \cdots & \cdots & \cdots \\ 0 & 0 & 0 & \sigma_{xp,j}^{2*} \end{pmatrix})$$

$$f(y | x, P) = \sum_{j=1}^{\infty} w_j(x) N(y | x' \beta_j^*, \sigma_{y|x,j}^{2*}).$$

# Improving prediction

- Write the kernels as

$$N_{p+1}(x, y | \mu, \Sigma) = N_p(x | \mu_x, \Sigma_x) N(y | x' \beta, \sigma_{y|x}^2)$$

and use simple spherical  $x$ -kernels (Shahbaba and Neal, 2009, Hannah et al., 2011). Thus,

$$f_x(x | P) = \sum_{j=1}^{\infty} w_j N_p(x | \mu_{xj}^*, \begin{pmatrix} \sigma_{x1,j}^{2*} & 0 & \cdots & 0 \\ 0 & \sigma_{x2,j}^{2*} & 0 & 0 \\ 0 & \cdots & \cdots & \cdots \\ 0 & 0 & 0 & \sigma_{xp,j}^{2*} \end{pmatrix})$$

$$f(y | x, P) = \sum_{j=1}^{\infty} w_j(x) N(y | x' \beta_j^*, \sigma_{y|x,j}^{2*}).$$

- Denote the  $x$ -parameters  $\phi = ((\mu_{x1}, \dots, \mu_{xp}), (\sigma_{x1}^2, \dots, \sigma_{xp}^2))$  and the  $y|x$ -parameters  $\theta = (\beta, \sigma_{y|x}^2)$ . Assign independent conjugate priors  $G_{0,\phi}(\phi)$  and  $G_{0,\theta}(\theta)$  (that leads to an enriched conjugate prior for  $(\mu, \Sigma)$ )

# Improving prediction

- Write the kernels as

$$N_{p+1}(x, y \mid \mu, \Sigma) = N_p(x \mid \mu_x, \Sigma_x) N(y \mid x' \beta, \sigma_{y|x}^2)$$

and use simple spherical  $x$ -kernels (Shahbaba and Neal, 2009, Hannah et al., 2011). Thus,

$$f_x(x \mid P) = \sum_{j=1}^{\infty} w_j N_p(x \mid \mu_{xj}^*, \begin{pmatrix} \sigma_{x1,j}^{2*} & 0 & \cdots & 0 \\ 0 & \sigma_{x2,j}^{2*} & 0 & 0 \\ 0 & \cdots & \cdots & \cdots \\ 0 & 0 & 0 & \sigma_{xp,j}^{2*} \end{pmatrix})$$

$$f(y \mid x, P) = \sum_{j=1}^{\infty} w_j(x) N(y \mid x' \beta_j^*, \sigma_{y|x,j}^{2*}).$$

- Denote the  $x$ -parameters  $\phi = ((\mu_{x1}, \dots, \mu_{xp}), (\sigma_{x1}^2, \dots, \sigma_{xp}^2))$  and the  $y|x$ -parameters  $\theta = (\beta, \sigma_{y|x}^2)$ . Assign independent conjugate priors  $G_{0,\phi}(\phi)$  and  $G_{0,\theta}(\theta)$  (that leads to an enriched conjugate prior for  $(\mu, \Sigma)$ )
- In the DP mixture model, assume individual-specific  $(\phi_i, \theta_i)$  as a random sample from  $G \sim DP(\alpha G_{0,\phi} G_{0,\theta})$

# Joint DP mixture model

$$\begin{aligned} Y_i | x_i, \beta_i, \sigma_{y,i}^2 &\stackrel{\text{ind}}{\sim} N(\beta_i' x_i, \sigma_{y,i}^2), \quad \theta_i = (\beta_i, \sigma_{y,i}^2), \\ X_i | \mu_i, \sigma_{x,i}^2 &\stackrel{\text{ind}}{\sim} \prod_{h=1}^p N(\mu_{x,h,i}, \sigma_{x,h,i}^2), \quad \phi_i = (\mu_{x,i}, \sigma_{x,i}^2) \\ (\theta_i, \phi_i) | G &\stackrel{\text{i.i.d.}}{\sim} G, \\ \mathbf{G} &\sim DP(\alpha G_{0\theta} \times G_{0\phi}). \end{aligned}$$

with  $G_{0\theta}$  and  $G_{0\phi}$  independent conjugate Normal-Inverse Gamma priors.



The model is flexible, and MCMC computations are standard.

BUT, if  $p$  is large, many (independent) kernels will be typically needed to describe the (dependent) marginal  $f_x$ , while the relationship  $Y | x$  can be smoother.

However, the DP only allows joint clusters of  $(\phi_i, \theta_i)$ ,  $i = 1, \dots, n$ .

Given its crucial role, difficulties in the random partition have relevant consequences on prediction

We would like to use a prior on  $P$  that allows many  $\phi_i$  clusters, to fit  $f_x$ , but fewer  $\theta_j$  clusters, and it is still conjugate, so that computations remain simple.

The model is flexible, and MCMC computations are standard.

BUT, if  $p$  is large, many (independent) kernels will be typically needed to describe the (dependent) marginal  $f_x$ , while the relationship  $Y | x$  can be smoother.

However, the DP only allows joint clusters of  $(\phi_i, \theta_i)$ ,  $i = 1, \dots, n$ .

Given its crucial role, difficulties in the random partition have relevant consequences on prediction

We would like to use a prior on  $P$  that allows many  $\phi_i$  clusters, to fit  $f_x$ , but fewer  $\theta_j$  clusters, and it is still conjugate, so that computations remain simple.

⇒ Enriched Dirichlet Process (Wade, Mongelluzzo, P., 2011)

# Enriched Dirichlet Process

Extending the idea that leads from the Dirichlet distribution (conjugate to the multinomial) to the enriched (or generalized) Dirichlet (Connor & Mosiman (1969); and from the DP to Doksum (1974) neutral-to-the-right priors for random probability measures on the real line.

EDP: Assign a prior for the random prob. measure  $G(\theta, \phi)$  by assuming:

- $P_\theta \sim DP(\alpha_\theta P_{0,\theta})$
- for any  $\theta$ ,  $P_{\phi|\theta} \sim DP(\alpha_\phi(\theta) P_{0,\phi|\theta})$

all independent.

The EDP gives a **nested random partition**:  $\rho_n = (\rho_{n,\theta}, \rho_{n,\phi})$ , that **allows many  $\phi$ -clusters inside each  $\theta$ -cluster**.

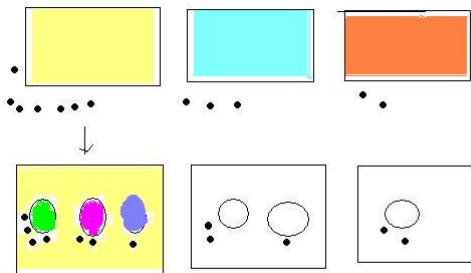
## EDP: nested random partition

$\rho_n = (\rho_{n,\theta}, \rho_{n,\phi})$  : many  $\phi$ -clusters inside each  $\theta$ -cluster.

# EDP: nested random partition

$\rho_n = (\rho_{n,\theta}, \rho_{n,\phi})$  : many  $\phi$ -clusters inside each  $\theta$ -cluster.

- $P_\theta \sim DP(\alpha P_{0\theta})$  gives a Chinese restaurant process: customers choose *restaurants*, and restaurants are colored with colors  $\theta_h^* \stackrel{iid}{\sim} P_{0\theta}$  (nonatomic):
- $P_{\phi|\theta} \sim DP(\alpha_\phi(\theta) P_{0\phi|\theta})$  gives a *nested* CRP: within each restaurant, customers sits at tables as in the CRP. Tables in restaurant  $\theta_h^*$  are colored with colors  $\phi_{j|h}^* \stackrel{iid}{\sim} P_{0,\phi|\theta}(\phi | \theta)$ .



nested Chinese restaurant

- Model (replace DP with EDP):

$$Y_i | x_i, \beta_i, \sigma_{y,i}^2 \stackrel{\text{ind}}{\sim} N(\beta_i' x_i, \sigma_{y,i}^2), \quad \theta_i = (\beta_i, \sigma_{y,i}^2),$$

$$X_i | \mu_i, \sigma_{x,i}^2 \stackrel{\text{ind}}{\sim} \prod_{h=1}^P N(\mu_{x,h,i}, \sigma_{x,i,h}^2), \quad \phi_i = (\mu_{x,i}, \sigma_{x,i}^2)$$

$$(\theta_i, \phi_i) | \mathbf{G} \stackrel{\text{i.i.d.}}{\sim} \mathbf{G},$$

$$\mathbf{G} \sim \text{EDP}(\alpha_\theta, \alpha_\phi(\cdot), G_{0,\theta} \times G_{0,\phi|\theta}).$$

- Computations remain simple, as the EDP is a conjugate prior;
- Inference on a cluster-specific  $\theta_j^* = (\beta_j; \sigma_{y|x^*})$ , (thus, ultimately, on the conditional density  $f(y | x, \theta)$ ), exploits the information from the observations in all the  $\phi_h^*$ -clusters that share the same  $\theta_j^* \Rightarrow$  **much improved inference and prediction**

(Wade, Dunson, Petrone, Trippa, JMLR (2014))

# Simulation study

Toy example to demonstrate two key advantages of the EDP model

- it can recover the true coarser  $\theta$ -partition
- improved prediction and smaller credible intervals result.

Data: 200 obs  $(x_i, y_i)$  where

The covariates  $X_i$  are sampled from a  $p$ -variate normal,

$$X_i = (X_{i,1}, \dots, X_{i,p})' \stackrel{iid}{\sim} N(\mu, \Sigma),$$

centered at  $\mu = (4, \dots, 4)'$  with  $\Sigma_{h,h} = 4$  for  $h = 1, \dots, p$ , and covariances that model two groups of covariates with different correlation structure.

# Simulation study

Toy example to demonstrate two key advantages of the EDP model

- it can **recover the true coarser  $\theta$ -partition**
- **improved prediction** and **smaller credible intervals** result.

Data: 200 obs  $(x_i, y_i)$  where

The covariates  $X_i$  are sampled from a  $p$ -variate normal,

$$X_i = (X_{i,1}, \dots, X_{i,p})' \stackrel{iid}{\sim} N(\mu, \Sigma),$$

centered at  $\mu = (4, \dots, 4)'$  with  $\Sigma_{h,h} = 4$  for  $h = 1, \dots, p$ , and covariances that model two groups of covariates with different correlation structure.

The true regression only depends on the first covariate; it is a nonlinear regression obtained as a mixture

$$Y_i | x_i \stackrel{ind}{\sim} p(x_{i,1})N(y_i | \beta_{1,0} + \beta_{1,1}x_{i,1}, \sigma_1^2) + (1 - p(x_{i,1}))N(y_i | \beta_{2,0} + \beta_{2,1}x_{i,1}, \sigma_2^2)$$



Table: Prediction error for both models as  $p$  increases.

	$p = 1$		$p = 5$		$p = 10$		$p = 15$	
	$\hat{l}_1$	$\hat{l}_2$	$\hat{l}_1$	$\hat{l}_2$	$\hat{l}_1$	$\hat{l}_2$	$\hat{l}_1$	$\hat{l}_2$
DP	<b>0.03</b>	<b>0.05</b>	0.16	0.2	0.25	0.34	0.26	0.34
EDP	0.04	<b>0.05</b>	<b>0.06</b>	<b>0.1</b>	<b>0.09</b>	<b>0.16</b>	<b>0.12</b>	<b>0.21</b>

# Real data: Alzheimer disease diagnosis

Data were obtained from the [Alzheimer's Disease Neuroimaging Initiative \(ADNI\)](#) database, which is publicly accessible at UCLA's Laboratory of Neuroimaging.

**covariates:** summaries of  $p = 15$  brain structures computed from structural MRI obtained at the first visit for 377 patients, of which 159 have been diagnosed with AD and 218 are cognitively normal (CN).

**response:**  $Y_i = 1$  (cognitively normal subject); or  $= 0$  (diagnosed with AD).

**Aim:** Prediction of AD status

# Extension of the model to binary response

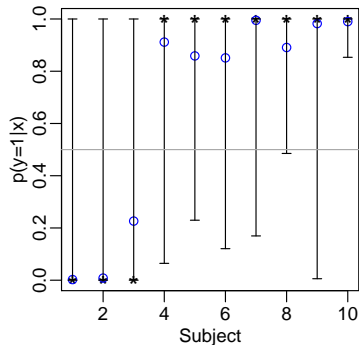
The model is extended to a local probit model:

$$Y_i | x_i, \beta_i \stackrel{iid}{\sim} \text{Bern}(\Phi(\underline{x}_i \beta_i)), \quad X_i | \mu_i, \sigma_i^2 \stackrel{iid}{\sim} \prod_{h=1}^p \text{N}(\mu_{i,h}, \sigma_{i,h}^2),$$
$$(\beta_i, \mu_i, \sigma_i^2) | G \stackrel{iid}{\sim} G, \quad \mathbf{G} \sim Q.$$

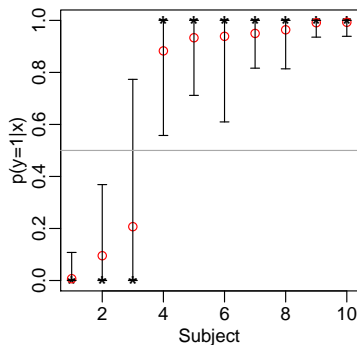
First, **DP prior for  $G$** :  $G \sim DP(\alpha, G_{0\beta} \times G_{0\psi})$ , with  $G_{0\beta} = \text{N}(0_p, C^{-1})$  and  $G_{0\psi}$  product of  $p$  normal-inverse gamma. We let  $\alpha \sim \text{Gamma}(1, 1)$ .

**EDP prior for  $G$** : correlation between the measurements of the brain structures and non-normal univariate histograms of the covariates suggest that many Gaussian kernels with local independence will be needed to approximate the density of  $X$ . The conditional density of the response, on the other hand, may not be so complicated. This motivates the choice of an EDP prior.

# Prediction for 10 new subjects



(a) DP



(b) EDP

**Figure:** Predicted probability of being healthy against subject index for 10 new subjects and represented with circles (DP in blue and EDP in red) with the true outcome as black stars. The bars about the prediction depict the 95% credible intervals.

- 1 Preliminaries
- 2 Random design: DP mixtures for  $f(x, y)$
- 3 Example 1 (improving prediction by restricted partition models)
- 4 Example 2 (Improving prediction by Enriched DP mixtures)
- 5 Fixed design: Dependent stick-breaking mixture models**
- 6 discussion

## Fixed design: Conditional models

Sample  $(x_i, Y_{i,\nu}), \nu = 1, \dots, n_i$ ,  $x$  fixed input.

Interest in  $f_x(y)$  (no longer a conditional  $f(y | x)$ ).

## Fixed design: Conditional models

Sample  $(x_i, Y_{i,\nu}), \nu = 1, \dots, n_i$ ,  $x$  fixed input.

Interest in  $f_x(y)$  (no longer a conditional  $f(y | x)$ ).

Joint DP mixture models are used in this context, too. Yet, they unnecessarily require to model the marginal density  $f_x(x)$ .

## Fixed design: Conditional models

Sample  $(x_i, Y_{i,\nu}), \nu = 1, \dots, n_i$ ,  $x$  fixed input.

Interest in  $f_x(y)$  (no longer a conditional  $f(y | x)$ ).

Joint DP mixture models are used in this context, too. Yet, they unnecessarily require to model the marginal density  $f_x(x)$ .

In a Bayesian approach, one wants to assign a prior on the set of random prob. measures  $\{F_x(\cdot), x \in \mathcal{X}\}$ .

The random measures  $F_x$  must be dependent, as we want some smoothness along  $x$ , for borrowing strength; and have a given marginal distribution, e.g.  $F_x \sim DP$ .

Early proposals: Cifarelli & Regazzini (1978), extending mixtures of DP (Antoniak, 1974).

Influential idea: dependent Dirichlet processes (DDP), McEachern (1999; 2000), based on dependent stick-breaking constructions.



## DDP mixture models

Use a mixture model for  $f_x(y)$

$$Y_{i,\nu} | x, G_x \stackrel{ind}{\sim} f_{G_x}(y) = \int K(y | x, \theta) dG_x(\theta)$$

where the mixing distribution  $G_x$  is indexed by  $x$ .

# DDP mixture models

Use a mixture model for  $f_x(y)$

$$Y_{i,\nu} | x, G_x \stackrel{\text{ind}}{\sim} f_{G_x}(y) = \int K(y | x, \theta) dG_x(\theta)$$

where the mixing distribution  $G_x$  is indexed by  $x$ .

A DDP prior on  $\{G_x, x \in \mathcal{X}\}$  assumes that  $G_x \sim DP(\alpha(x)G_{0,x})$ , and dependence is introduced through the dependent stick-breaking constructions:

$$G_x(\cdot) = \sum_{j=1}^{\infty} p_j(x) \delta_{\theta_j^*(x)}(\cdot)$$

where, for each  $j$ :

- $(w_j(x), x \in \mathcal{X})$  is a stochastic process, with stick-breaking construction  $w_1(x) = v_1(x)$ ;  $w_2(x) = v_2(x)(1 - v_1(x))$ , ... where  $(v_j(x), x \in \mathcal{X})$  is a stochastic process with marginals  $v_j(x) \sim \text{Beta}(1, \alpha(x))$ , and the  $v_j(\cdot)$  are independent across  $x$ ;
- $(\theta_j^*(x), x \in \mathcal{X})$  is a stochastic process with marginals  $G_{0,x}$ . The processes  $\theta_j(\cdot)$  are independent across  $j$ , and indep. of the  $v_j(\cdot)$ .

# Conditional approach: Single weights

The DDP allows both the mixing weights and the atoms to depend on  $x$ . But this is redundant, and one either considers 1) models with **single weights** and 2) models with **covariate dependent weights**.

**Single weights:** assume  $w_j(x) = w_j$  with flexible  $\theta_j(x)$ :

$$f_{G_x}(y|x) = \sum_{j=1}^{\infty} w_j K(y|\theta_j(x)).$$

- Ex.  $K(y|\theta_j(x)) = N(y|\mu_j(x), \sigma_j^2)$  with  $\mu_j \stackrel{iid}{\sim}$  GP.
- Popular because inference relies on established algorithms for BNP mixtures.

# Conditional approach: Covariate dependent weights

**Covariate dependent weights:** flexible  $w_j(x)$  with  $\theta_j(x) = \theta_j$ :

$$f_{P_x}(y|x) = \sum_{j=1}^{\infty} w_j(x) K(y|\theta_j, x),$$

for ex.  $K(y|\theta_j(x)) = N(y|\beta_j'x, \sigma_j^2)$ .

- Most techniques to define  $w_j(x)$  s.t.  $\sum_j w_j(x) = 1$  use a **stick-breaking approach**:

$$w_1(x) = v_1(x) \text{ and for } j > 1 \quad w_j(x) = v_j(x) \prod_{j' < j} (1 - v_{j'}(x)),$$

- Proposals for  $v_j(x)$  include Griffin and Steel (2006), Dunson and Park (2008), Ren et al. (2011), Rodriguez and Dunson (2011). A formally related, but differently motivated, approach is the **normalized weights model** by Antoniano, Wade, Walker (2014).

How can we compare?

The two examples show clear improvement in the predictive performance. We would like to **formally express the gain of information** that a model/prior can provide.

The two examples show clear improvement in the predictive performance. We would like to **formally express the gain of information** that a model/prior can provide.

So far

- Careful and detailed study of the information and assumption introduced through the prior/model; and the analytic implications on the predictive distribution
- compare on simulated and real data.

Some findings (Peruzzi, Petrone & Wade, 2016+)

- The 'single weights' mixture model needs flexible atoms, for giving reasonable prediction. But this has drawbacks, including identifiability.
- Covariate-dependent weights allow local selection of the clustering, and generally better prediction.

In particular, Peruzzi & Wade (2016+) compared the kernel stick-breaking model (Dunson & Park, 2008) and the normalized weights model (Antoniano, Wade, Walker, 2014). The latter appears to allow faster computations.

A substantial challenge for comparison is due to computations: find an algorithm that can give fast computations, and be fairly easily adapted to different models.

Peruzzi & Wade (2016+) suggest a clever algorithm, based on adaptive truncation and Metropolis-Hastings, moving from a proposal by Griffin (2016).

This is a necessary basis for providing R-packages for density regression and conditional density estimation.



- 1 Preliminaries
- 2 Random design: DP mixtures for  $f(x, y)$
- 3 Example 1 (improving prediction by restricted partition models)
- 4 Example 2 (Improving prediction by Enriched DP mixtures)
- 5 Fixed design: Dependent stick-breaking mixture models
- 6 discussion**

I tried to give an overview of proposals for Bayesian density regression, based on Sethuraman's construction of Dirichlet priors.

Careful study of the finite-sample properties of these models, for comparison

Together with flexible, easily-exportable computational tools, these are necessary steps for providing R-packages for density regression.

I limited the presentation to the Dirichlet process. But much richer classes of stick-breaking priors, with different properties and implied partition models, have originated by Sethuraman's influential work.