

ADVANCED BAYESIAN LEARNING BAYESIAN NONPARAMETRICS SPRING 2014

Mattias Villani

Division of Statistics
Department of Computer and Information Science
Linköping University

- ▶ Reminder: Multinomial data - Dirichlet prior
- ▶ Bayesian histograms
- ▶ Dirichlet process
- ▶ Beyond the Dirichlet process: Pitman-Yor and Probit stick-breaking
- ▶ The Dirichlet process mixtures
- ▶ MCMC for Dirichlet process mixtures

THE DIRICHLET DISTRIBUTION

- ▶ $\theta \sim \text{Dirichlet}(a_1, \dots, a_k)$ with density

$$p(\theta_1, \theta_2, \dots, \theta_k) \propto \prod_{j=1}^k \theta_j^{a_j-1}.$$
- ▶ Define $\alpha = \sum_{j=1}^k a_j$ and $\pi_0 = a/\alpha$.
- ▶ Expected value and variance of the *Dirichlet*(a_1, \dots, a_k) distribution

$$\mathbb{E}(\theta_j) = \frac{a_j}{\alpha} = \pi_{0j} \qquad \mathbb{V}(\theta_j) = \frac{\mathbb{E}(\theta_j) [1 - \mathbb{E}(\theta_j)]}{1 + \alpha}$$
- ▶ Note that α is a **precision** parameter (large α means low variance).

CONJUGATE ANALYSIS FOR MULTINOMIAL DATA

- ▶ **Data:** $y = (n_1, \dots, n_k)$, where n_j = number of items in category j .
- ▶ **Prior**

$$\theta \sim \text{Dirichlet}(a_1, \dots, a_k)$$
- ▶ **Likelihood**

$$p(n_1, n_2, \dots, n_k | \theta_1, \theta_2, \dots, \theta_k) \propto \prod_{j=1}^k \theta_j^{n_j}$$
- ▶ **Posterior**

$$\theta | n_1, \dots, n_k \sim \text{Dirichlet}(n_1 + a_1, \dots, n_k + a_k)$$
- ▶ **Posterior expected value**

$$\mathbb{E}(\theta_j | n_1, \dots, n_k) = \frac{n_j + a_j}{n + \alpha}$$

- ▶ Histogram partitions the data space $\tilde{\zeta}_0 < \tilde{\zeta}_1 < \dots < \tilde{\zeta}_k$ and records how many observations end up in each bin (B_h). Multinomial data.
- ▶ Probability model for **histograms**

$$f(y) = \sum_{h=1}^k 1_{\tilde{\zeta}_{h-1} < y \leq \tilde{\zeta}_h} \left(\frac{\pi_{\tau_h}}{\tilde{\zeta}_h - \tilde{\zeta}_{h-1}} \right)$$

- ▶ n_h = number of data points in partition (bin) h : $\tilde{\zeta}_{h-1} < y \leq \tilde{\zeta}_h$.
- ▶ **Prior** on $\pi = (\pi_1, \dots, \pi_k)$

$$\pi \sim \text{Dirichlet}(a_1, \dots, a_k)$$

- ▶ **Posterior**

$$\pi | n_1, \dots, n_k \sim \text{Dirichlet}(n_1 + a_1, \dots, n_k + a_k)$$

BAYESIAN HISTOGRAMS, CONT.

- ▶ **Posterior**

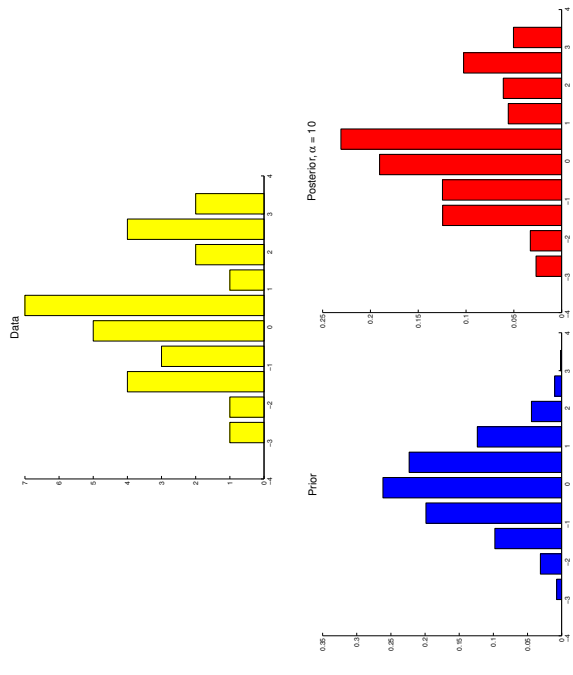
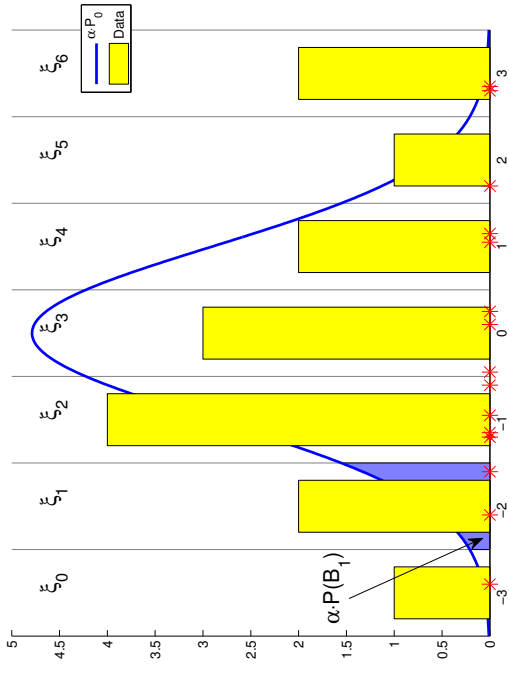
$$\pi | n_1, \dots, n_k \sim \text{Dirichlet}(n_1 + a_1, \dots, n_k + a_k)$$

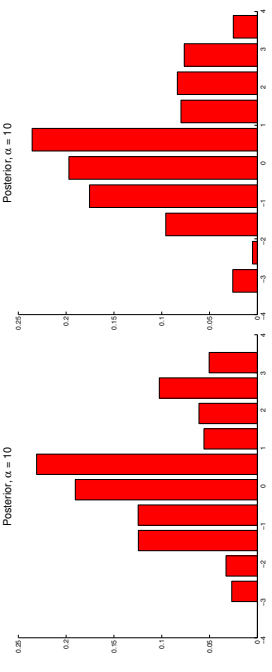
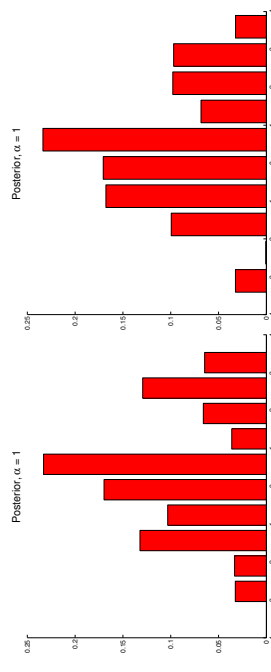
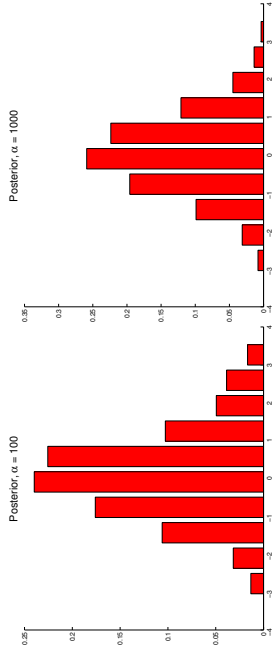
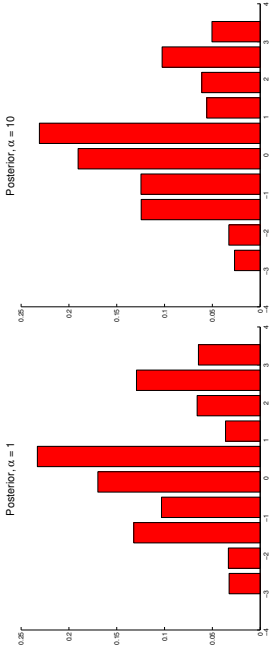
- ▶ Specify a_1, \dots, a_k through $\pi_0 = (\pi_{01}, \dots, \pi_{0k})$ and $\alpha = \sum_{j=1}^k a_j$.
- ▶ Specify π_0 from a **base distribution** P_0 . For the h th bin:

$$\pi_{0h} = P_0(B_h) = \Pr(\tilde{\zeta}_{h-1} < y \leq \tilde{\zeta}_h)$$

- ▶ The Dirichlet prior is a **computational dream**, and it is easy to specify the hyperparameters π_0 and α .
- ▶ But, the Dirichlet prior lacks **smoothness**: all pairs of bins have negative correlations, regardless of how near they are.
- ▶ Sensitive to the choice of bins.

ILLUSTRATION OF BAYESIAN HISTOGRAMS





THE DIRICHLET PROCESS

- ▶ Let B_1, B_2, \dots, B_k be a partition of the outcome space Ω .
- ▶ Let $P(B_1), \dots, P(B_k)$ denote the distribution over the partition.
- ▶ Dirichlet distribution is a **distribution over a space of distributions**:

$$(P(B_1), \dots, P(B_k)) \sim \text{Dirichlet}(\alpha P_0(B_1), \dots, \alpha P_0(B_k))$$

- where P_0 is a fixed probability measure (e.g. the $N(0, 1)$ density).
- ▶ Dirichlet distribution is closed under summation or splitting of bins.
- ▶ Can be used to define a **stochastic process** in a consistent way. Compare with GPs.
- ▶ A random probability measure P follows a **Dirichlet process** $P \sim DP(\alpha \cdot P_0)$ with base measure P_0 and precision parameter α iff

$$(P(B_1), \dots, P(B_k)) \sim \text{Dirichlet}(\alpha P_0(B_1), \dots, \alpha P_0(B_k))$$

for any finite (measurable) partition B_1, \dots, B_k .

THE DIRICHLET PROCESS - PROPERTIES

- ▶ If $P \sim DP(\alpha P_0)$ then

$$P(B) \sim \text{Beta}[\alpha P_0(B), \alpha(1 - P_0(B))], \text{ for any } B \in \mathcal{B}$$

$$E[P(B)] = P_0(B)$$

$$\text{Var}[P(B)] = P_0(B)[1 - P_0(B)] / (1 + \alpha)$$

- ▶ **Model**

$$y_i | P \stackrel{iid}{\sim} P, \text{ for } i = 1, \dots, n$$

- ▶ **Prior**

$$P \sim DP(\alpha P_0)$$

- ▶ **Posterior for a finite partition**, $P(B_1), \dots, P(B_k) | \mathbf{y}$ is

$$\text{Dir}\left(\alpha P_0(B_1) + \sum_{i=1}^n 1_{y_i \in B_1}, \dots, \alpha P_0(B_k) + \sum_{i=1}^n 1_{y_i \in B_k}\right)$$

- **Posterior** for the unknown probability distribution P

$$P|y_1, \dots, y_n \sim DP\left(\alpha P_0 + \sum_{i=1}^n \delta_{y_i}\right)$$

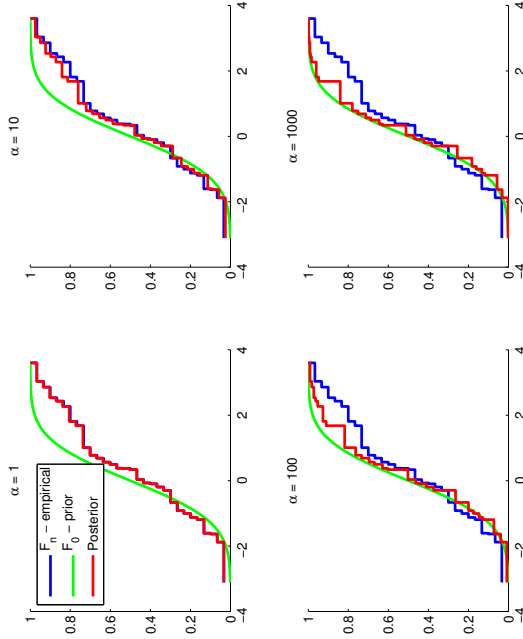
- Since

$$P(B) \sim \text{Beta}\left(\alpha P_0(B) + \sum_{i=1}^n \mathbf{1}_{y_i \in B}, \alpha(1 - P_0(B)) + \sum_{i=1}^n \mathbf{1}_{y_i \in B^c}\right)$$

so

$$E(P(B)|y_1, \dots, y_n) = \left(\frac{\alpha}{\alpha + n}\right) P_0(B) + \left(\frac{n}{\alpha + n}\right) \sum_{i=1}^n \frac{1}{n} \mathbf{1}_{y_i \in B^c}$$

ESTIMATING A D.F. WITH A DP PRIOR



- If $B = (-\infty, y]$ then

$$E(F(y)|y_1, \dots, y_n) = \left(\frac{\alpha}{\alpha + n}\right) F_0(y) + \left(\frac{n}{\alpha + n}\right) F_n(y)$$

- where

- $F(y)$ is the unknown d.f.
- $F_0(y)$ is the d.f. from P_0
- $F_n(y) = \frac{1}{n} \sum \mathbf{1}_{y_i \leq y}$ is the empirical d.f.

- Note: under the DP posterior, $F(\cdot)$ is discrete with probability one. Not great for continuous data ...

- This is true in general: **realisations from a DP are discrete with probability one.**

STICK-BREAKING CHARACTERIZATION OF THE DP

- $P \sim DP(\alpha P_0)$ is equivalent to an infinite mixture of point masses

$$P(\cdot) = \sum_{h=1}^{\infty} \pi_h \delta_{\theta_h}$$

$$\pi_h = V_h \prod_{\ell < h} (1 - V_\ell)$$

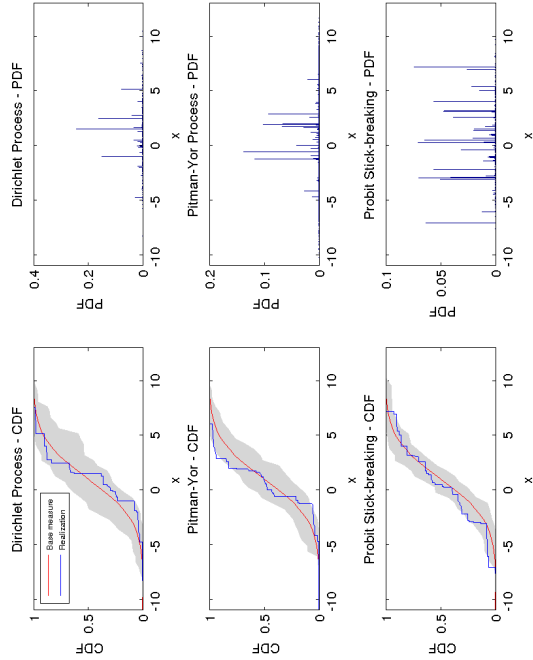
$$V_h \stackrel{iid}{\sim} \text{Beta}(1, \alpha)$$

$$\theta_h \stackrel{iid}{\sim} P_0$$

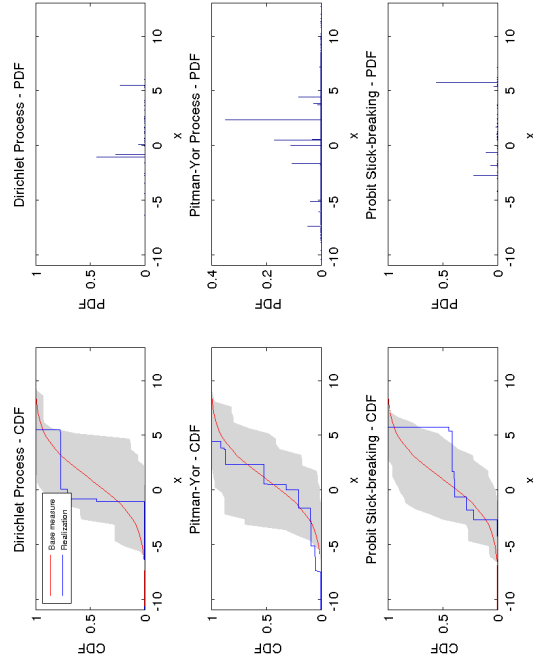
- Stick picture
- Alternative notation for $P \sim DP(\alpha P_0)$:

$$\pi = (\pi_1, \pi_2, \dots) \sim \text{Stick}(\alpha) \text{ and } \theta_h \sim P_0$$

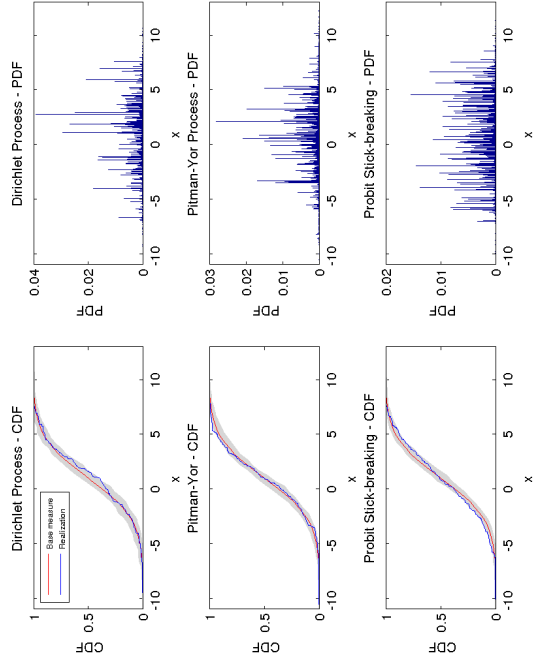
SIMULATING STICK-BREAKING PRIORS $\alpha = 10$



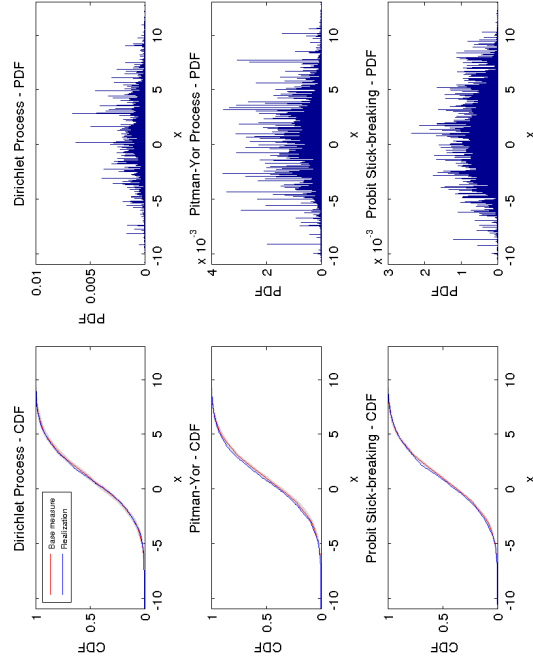
SIMULATING STICK-BREAKING PRIORS $\alpha = 1$



SIMULATING STICK-BREAKING PRIORS $\alpha = 100$



SIMULATING STICK-BREAKING PRIORS $\alpha = 1000$



- Pitman-Yor process with parameters P_0 , $0 \leq a < 1$ and $b > -a$:

$$P(\cdot) = \sum_{h=1}^{\infty} \pi_h \delta_{\theta_i} \quad \theta_h \stackrel{iid}{\sim} P_0$$

$$\pi_h = V_h \prod_{\ell < h} (1 - V_\ell)$$

$$V_h \stackrel{iid}{\sim} \text{Beta}(1 - a, b + ha)$$

- Probit stick-breaking with parameters μ and σ :

$$P(\cdot) = \sum_{h=1}^{\infty} \pi_h \delta_{\theta_i} \quad \theta_h \stackrel{iid}{\sim} P_0$$

$$\pi_h = V_h \prod_{\ell < h} (1 - V_\ell)$$

$$V_h = \Phi(x_h), \quad \text{where } x_h \stackrel{iid}{\sim} N(\mu, \sigma^2)$$

INFINITE MIXTURE MODELS - DP MIXTURES

- General mixture formulation

$$f(y|P) = \int \mathcal{K}(y|\theta) dP(\theta)$$

where $\mathcal{K}(y|\theta)$ is a kernel and $P(\theta)$ is a **mixing measure**.

- Example 1: **Student-t**, $t_\nu(\mu, \sigma^2)$. $\mathcal{K}(y|\theta) = \phi(y|\mu, \lambda)$ where μ is fixed, $\theta = \lambda$ and $P(\theta)$ is the *Inv- χ^2* (ν, σ^2) distribution.
- Example 2: **Finite mixture of normals**. $\phi(y|\mu, \sigma^2)$, $\theta = (\mu, \sigma^2)$. $P(\theta)$ is a discrete distr. with $\Pr[\theta = (\mu_j, \sigma_j^2)] = \pi_j$, for $j = 1, \dots, k$.
- Example 3: $P \sim DP(\alpha P_0)$ yields the **infinite mixture**

$$f(y) = \sum_{h=1}^{\infty} \pi_h \mathcal{K}(y|\theta_h^*), \quad \pi \sim \text{Stick}(\alpha).$$

FINITE MIXTURE MODELS

- Mixture of normals

$$p(y) = \sum_{j=1}^k \pi_j \cdot \phi(y; \mu_j, \sigma_j^2)$$

- Use **allocation variables**: $l_i = j$ if y_i comes from $\phi(y; \mu_j, \sigma_j^2)$.
- Let $l = (l_1, \dots, l_n)'$ and $n_j = \sum_{i=1}^n (l_i = j)$.
- **Gibbs sampling** algorithm:
 - $\pi_1, \dots, \pi_k \mid l, y \sim \text{Dirichlet}(a_1 + n_1, a_2 + n_2, \dots, a_k + n_k)$
 - $\sigma_j^2 \mid l, y \sim \text{Inv-}\chi^2$ and $\mu_j \mid l, \sigma_j^2, y \sim N$ for $j = 1, \dots, k$.
 - $l_i \mid \pi, \mu, \sigma^2, y \sim \text{Multinomial}(\omega_{i,1}, \dots, \omega_{i,k}), i = 1, \dots, n$, where

$$\omega_{i,j} = \frac{\pi_j \cdot \phi(y_i; \mu_j, \sigma_j^2)}{\sum_{q=1}^k \pi_q \cdot \phi(y_i; \mu_q, \sigma_q^2)}.$$

DP MIXTURE IS LIKE A FINITE MIXTURE WITH LARGE k

- In infinite mixtures every observation has its own parameter θ_i

$$y_i \sim \mathcal{K}(\theta_i)$$
- DP is almost surely discrete \Rightarrow **ties**: some of the θ_i will have exactly the same values. **DP leads to clustering** of the θ_i .
- Each observation has **potentially** its own parameter θ_i , but that **parameter may be shared by other observations**.
- In finite mixture models each observation also has its “own” parameter

$$y_i \mid l_i \sim \mathcal{K}(\theta_{l_i})$$

$$l_i \mid \pi \sim \text{Multinomial}(\pi_1, \dots, \pi_k)$$

$$\theta_i \sim P_0$$

$$\pi \sim \text{Dirichlet}(\alpha/k, \dots, \alpha/k)$$

- Neal (2000) shows that this finite mixture model approaches the DP mixture when $k \rightarrow \infty$.

- Hierarchical representation of DP mixtures

$$y_i \sim \mathcal{K}(\theta_i), \quad \theta_i \sim P \quad P \sim DP(\alpha P_0)$$

- We can actually marginalize out P to obtain the Polya scheme

$$p(\theta_i | \theta_1, \dots, \theta_{i-1}) \sim \left(\frac{\alpha}{\alpha + i - 1} \right) P_0(\theta_i) + \left(\frac{1}{\alpha + i - 1} \right) \sum_{j=1}^{i-1} \delta_{\theta_j}$$

- So $p(\theta_i | \theta_1, \dots, \theta_{i-1})$ is a mixture of the base measure P_0 and point masses at the previously “drawn” θ -values.
- Way to think about the scary ‘Marginalizing out P ’: integrate out π in the finite mixture model and let $k \rightarrow \infty$. [Neal, 2000].

GIBBS SAMPLING DP MIXTURES - MARGINALIZING P

- Similar to Gibbs sampling for finite mixtures. Data augmentation with mixture component indicators l_i .

1. **Update component allocation** for i th observation y_i by sampling from multinomial

$$\Pr(l_i = j | \cdot) \propto \begin{cases} n_j^{(-i)} \mathcal{K}(y_i | \theta_j^*) & \text{for } j = 1, \dots, k^{(-i)} \\ \alpha \int \mathcal{K}(y_i | \theta) dP_0(\theta) & \text{for } j = k^{(-i)} + 1 \end{cases}.$$

2. **Update the unique parameter values** θ^* by sampling from

$$p(\theta_j^* | \cdot) \propto P_0(\theta_c^*) \prod_{i: l_i = j} \mathcal{K}(y_i | \theta_j^*)$$

- Note that, unlike finite mixtures, the l_i are **not independent** conditional on θ^* . This because we have marginalized out P . They have to be sampled **sequentially**.

- The so called **Polya scheme**:

$$p(\theta_i | \theta_1, \dots, \theta_{i-1}) \sim \left(\frac{\alpha}{\alpha + i - 1} \right) P_0(\theta_i) + \left(\frac{1}{\alpha + i - 1} \right) \sum_{j=1}^{i-1} \delta_{\theta_j}$$

- **Chinese restaurant process**:

- first customer sits at empty table and obtains the dish θ_1^* from P_0 .
- second customer
 - sits at first customer’s table with probability $\frac{1}{1+\alpha}$ and has dish θ_1^* or
 - sits at a new table with probability $\frac{\alpha}{1+\alpha}$ and has dish $\theta_2^* \sim P_0$.
- \vdots
- the i th customer
 - sits at table with dish θ_j^* with a probability proportional to n_j , the number of customers sitting at table j or
 - sits at a new table with probability proportional to α .

GIBBS SAMPLING FOR TRUNCATED DP MIXTURES

- Set upper bound N for the number of components. Approximate DP mixture with $\pi_h = 0$ for $h = N + 1, \dots$
- Posterior sampling for infinite mixtures is now very similar to finite mixture. The l_i can be sampled independently.

1. Update component allocation for i th observation y_i by sampling from multinomial

$$\Pr(l_i = j | \cdot) \propto \pi_j \mathcal{K}(y_i | \theta_j^*) \quad \text{for } j = 1, 2, \dots, N.$$

2. Update the stick-breaking weights [recall: $\pi_h = V_h \prod_{\ell < h} (1 - V_\ell)$]

$$V_j | \cdot \sim \text{Beta} \left(1 + n_j, \alpha + \sum_{q=j+1}^N n_q \right) \quad \text{for } j = 1, \dots, N - 1.$$

3. Update the unique parameter values $\theta_1^*, \dots, \theta_N^*$ by sampling just like in the finite mixture model. Sample θ^* from prior $P_0(\theta)$ for empty clusters.

MCMC FOR DP MIXTURES

- ▶ Let's look at the updating step:

$$\Pr(l_i = j | \cdot) \propto \begin{cases} n_j^{(-i)} \mathcal{K}(y_i | \theta_c^*) & \text{for } j = 1, \dots, k^{(-i)} \\ \alpha \int \mathcal{K}(y_i | \theta) dP_0(\theta) & \text{for } j = k^{(-i)} + 1 \end{cases}.$$

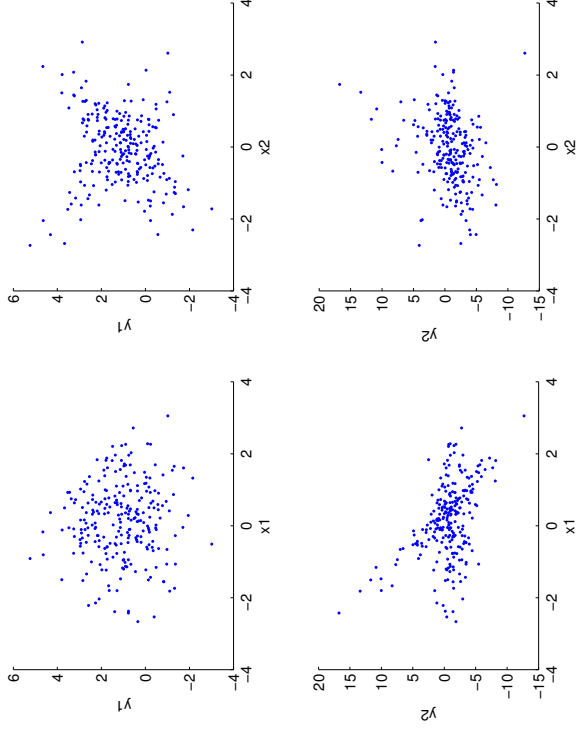
- ▶ **A customer chooses table based on:**

- ▶ the number of existing customers at the tables (with imaginary α customers at a new table)
- ▶ how compatible the taste of the customer (y_i) is to the different dishes served at occupied tables (θ_c^*)
- ▶ how compatible the taste of the customer (y_i) is to the different dishes that *may* be served at a new table.
- ▶ A $P_0(\theta)$ with large variance is equivalent to an very experimental cook. You never know what you get ...
- ▶ Hyperparameter α clearly matters for the number of clusters (tables), but so does P_0 .
- ▶ Hyperparameter α can be learned from data. Just add updating step.
- ▶ P_0 may contain hyperparameters (e.g. $P_0 = N(\mu, \sigma^2)$). Just add

MATTIAS VILLANI (STATISTICS, LIU) ADV BAYESIAN LEARNING

29 / 33

MIXTURE OF MULTIVARIATE REGRESSIONS - DATA



MATTIAS VILLANI (STATISTICS, LIU)

ADV BAYESIAN LEARNING

31 / 33

MIXTURE OF MULTIVARIATE REGRESSIONS - MODEL

- ▶ The response vector \mathbf{y} is p -dim. Covariates \mathbf{x} is q -dim.
- ▶ The model is of the form

$$p(\mathbf{y} | \mathbf{x}) = \sum_{j=1}^{\infty} \pi_j \cdot N(\mathbf{y}_j | \mathbf{B}_j \mathbf{x}_j; \Sigma_j)$$
- ▶ Each component in the mixture is a Gaussian multivariate regression with its own regression coefficient and covariance matrix:

$$\mathbf{y}_j = \mathbf{B}_j \mathbf{x}_j + \varepsilon_j, \quad \varepsilon_j \stackrel{iid}{\sim} N(0, \Sigma_j)$$

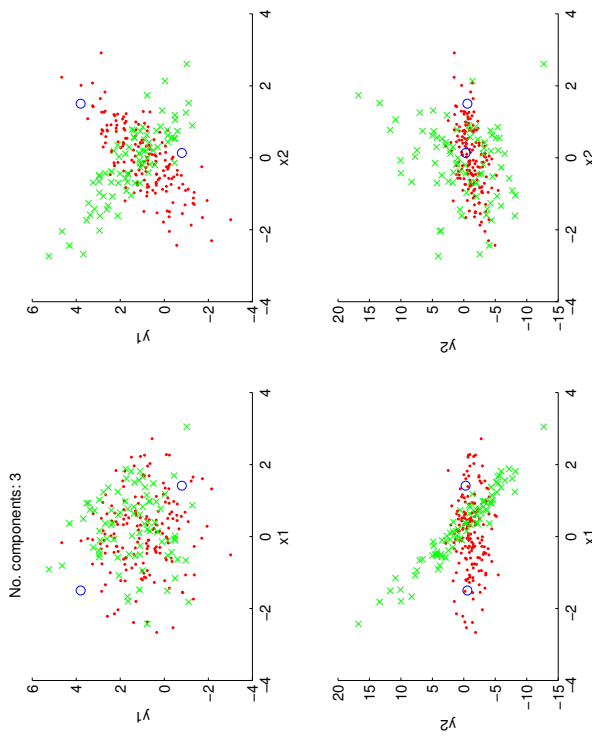
$$p \times 1 \quad p \times q \times 1 \quad p \times 1$$
- ▶ The mixture weights follow a DP stick prior $\pi \sim \text{Stick}(\alpha)$.

MATTIAS VILLANI (STATISTICS, LIU)

ADV BAYESIAN LEARNING

30 / 33

MIXTURE OF MULTIVARIATE REGRESSIONS - DPM



MATTIAS VILLANI (STATISTICS, LIU)

ADV BAYESIAN LEARNING

32 / 33

MIXTURE OF MULTIVARIATE REGRESSIONS - DPM

