

Gibbs Sampling Methods for Stick-Breaking Priors

Hermant Ishwaran

Lancelot F. James

JASA 2001

Presented by Daniel Johnson, Geoffrey Sun

Overview

- General class of *stick-breaking priors*
- Truncation result
- *Polya urn Gibbs sampler*
- *Blocked Gibbs sampler*
- Comparison

Stick-Breaking Priors

- Discrete random probability measures

$$\mathcal{P}(\cdot) = \sum_{k=1}^N p_k \delta_{Z_k}(\cdot), \quad 0 \leq p_k \leq 1 \text{ and } \sum_{k=1}^N p_k = 1 \text{ almost surely}$$

- Random measure $\mathbf{P}_N(\mathbf{a}, \mathbf{b})$ is *stick-breaking random measure*

$$p_1 = V_1 \quad \text{and} \quad p_k = (1 - V_1)(1 - V_2) \cdots (1 - V_{k-1}) V_k, \quad k \geq 2,$$

- $V_k \sim \text{Beta}(a_k, b_k)$, independent
- N finite or infinite

The Case $N < \text{Infinity}$

$$p_1 = V_1 \quad \text{and} \quad p_k = (1 - V_1)(1 - V_2) \cdots (1 - V_{k-1}) V_k, \quad k = 2, \dots, N.$$

- $N-1$ degrees of freedom

Setting $V_N = 1$ guarantees that $\sum_{k=1}^N p_k = 1$ with probability 1, because

$$1 - \sum_{k=1}^{N-1} p_k = (1 - V_1) \cdots (1 - V_{N-1}). \quad (4)$$

The Case $N = \text{Infinity}$

$$p_1 = V_1 \quad \text{and} \quad p_k = (1 - V_1)(1 - V_2) \cdots (1 - V_{k-1}) V_k, \quad k \geq 2,$$

- Necessary and sufficient conditions

Lemma 1. For the random weights in the $\mathcal{P}_\infty(\mathbf{a}, \mathbf{b})$ random measure,

$$\sum_{k=1}^{\infty} p_k = 1 \quad \text{a.s. iff} \quad \sum_{k=1}^{\infty} E(\log(1 - V_k)) = -\infty. \quad (5)$$

Alternatively, it is sufficient to check that $\sum_{k=1}^{\infty} \log(1 + a_k/b_k) = +\infty$.

- Computation?

The Pitman-Yor Process $PY(a,b)$

- Special case of *stick-breaking prior* $P_N(a,b)$
- $a_k = 1 - a \quad 0 \leq a < 1$
- $b_k = b + ka \quad b > -a$

Notable Pitman-Yor processes:

- Dirichlet process $a = 0, b = \alpha$
- 'Stable law' process $a = \alpha, b = 0$

Generalized Polya Urn Characterization

For a P-Y process $\mathbf{PY}(a,b)$

$$\mathbb{P}\{Y_i \in \cdot | Y_1, \dots, Y_{i-1}\} = \frac{b + am_i}{b + i - 1} H(\cdot) + \sum_{j=1}^{m_i} \frac{n_{j,i}^* - a}{b + i - 1} \delta_{Y_{j,i}^*}(\cdot), \quad i = 2, 3, \dots, n,$$

ζ_1, \dots, ζ_n iid H

$$Y_1 = \zeta_1,$$

$$(Y_i | Y_1, \dots, Y_{i-1}) = \begin{cases} \zeta_i & \text{with probability} \\ & (b + am_i)/(b + i - 1), \\ Y_{j,i}^* & \text{with probability} \\ & (n_{j,i}^* - a)/(b + i - 1), \end{cases}$$

for $i = 2, 3, \dots, n$.

Finite Dimensional Dirichlet $DP_N(\alpha^*H)$

- Like a Pitman-Yor process with
 - $a = -\alpha/N$
 - $b = \alpha > 0$
 - $N \geq n$

$$\mathcal{P}(\cdot) = \sum_{k=1}^N \frac{G_{k,N}}{\sum_{k=1}^N G_{k,N}} \delta_{Z_k}(\cdot), \quad G_{k,N} \stackrel{\text{iid}}{\sim} \text{Gamma}\left(\frac{\alpha}{N}\right).$$

- $P_N(\mathbf{a}, \mathbf{b})$ measure, but not actually P-Y as $(a < 0)$

$$DP_N(\alpha H)(g) \xrightarrow{d} DP(\alpha H)(g)$$

Truncation of $\mathcal{P}_\infty(\mathbf{a}, \mathbf{b})$ Measures

- Truncations $\mathcal{P}_N(\mathbf{a}, \mathbf{b})$ are computationally tractable
- Produce virtually indistinguishable measures

Theorem 2. Let p_k denote the random weights from a given $\mathcal{P}_\infty(\mathbf{a}, \mathbf{b})$ measure. If $\|\cdot\|_1$ denotes the \mathcal{L}_1 distance, then

$$\|\mu_N - \mu_\infty\|_1 \leq 4 \left(1 - \mathbb{E} \left[\left(\sum_{k=1}^{N-1} p_k \right)^n \right] \right).$$

- Ex. Dirichlet Process

$$\|\mu_N - \mu_\infty\|_1 \sim 4n \exp(-(N-1)/\alpha)$$

Proof of Theorem 2

Proof sketch on board if time & interest.

Summary

- Stick-breaking prior
- Pitman-Yor
- Truncation

Two Gibbs Samplers

Polya Urn Gibbs Sampler

- extension to models from Escobar, West, MacEachern, Ferguson.
- integrates out P in the hierarchical model
- must have a known prediction rule ($Y_i | \mathbf{Y}_{-i}$)
- P can be infinite

Blocked Gibbs Sampler

- works when prediction rule is unknown
- directly involve the prior in the sampler
- P needs to be finite (but we can apply truncation for infinite P)

Polya Urn Gibbs Samplers

hierarchical model with stick-breaking priors:

$$\begin{aligned}(X_i|Y_i, \theta) &\stackrel{\text{ind}}{\sim} \pi(X_i|Y_i, \theta), & i = 1, \dots, n, \\(Y_i|P) &\stackrel{\text{iid}}{\sim} P, \\ \theta &\sim \pi(\theta), \\ P &\sim \mathcal{P},\end{aligned}\tag{11}$$

integrate out P :

$$\begin{aligned}(X_i|Y_i, \theta) &\stackrel{\text{ind}}{\sim} \pi(X_i|Y_i, \theta), & i = 1, \dots, n, \\(Y_1, \dots, Y_n) &\sim \pi(Y_1, \dots, Y_n), \\ \theta &\sim \pi(\theta),\end{aligned}\tag{12}$$



CRP

Polya Urn Gibbs Samplers (PG)

- Assume priors = $\mathbf{PY}(a,b)$ or $\mathbf{DP}_N(\text{alpha}^*H)$
- Know the prediction rule ($Y_i | \mathbf{Y}_{-i}$)
- Want the posterior $\square(\mathbf{Y}, \theta | \mathbf{X})$
- Iterate between the two steps:
 - (a) ($Y_i | \mathbf{Y}_{-i}, \theta, \mathbf{X}$):

$$\begin{aligned}
 & \mathbb{P}\{Y_i \in \cdot | \mathbf{Y}_{-i}, \theta, \mathbf{X}\} && \text{assigned clusters} \\
 & q_0^* \propto (b + am) \int_y f(X_i | Y, \theta) H(dY), && \downarrow \\
 & q_j^* \propto (n_j^* - a) f(X_i | Y_j^*, \theta) && = q_0^* \mathbb{P}\{Y_i \in \cdot | \theta, X_i\} + \sum_{j=1}^m q_j^* \delta_{Y_j^*}(\cdot), \\
 & && \uparrow \\
 & && \text{unseen clusters}
 \end{aligned}$$

(b) ($\theta | \mathbf{Y}, \mathbf{X}$):

$$f(\theta | \mathbf{Y}, \mathbf{X}) \propto \pi(d\theta) \prod_{i=1}^n f(X_i | Y_i, \theta)$$

Polya Urn Gibbs Samplers (PG_a)

- Problem: \mathbf{Y}^* (unique Y 's) get stuck if q_0^* is large
- Acceleration step: resample \mathbf{Y}^*
- Let $\mathbf{C}=(C_1, \dots, C_n)$, indexing into \mathbf{Y}^* as a look-up table

(c) $(Y_j^* | \mathbf{C}, \theta, \mathbf{X})$

$$f(Y_j^* | \mathbf{C}, \theta, \mathbf{X}) \propto H(dY_j^*) \prod_{\{i: C_i=j\}} f(X_i | Y_j^*, \theta).$$

Limitations of PG and PG_a

1. slow mixing: a single Y_i at a time
2. relies on conjugacy for q_0^*
3. prior P is not directly involved, only depend on \mathbf{Y}
4. *requires a known urn scheme / prediction rule*

Blocked Gibbs Sampler

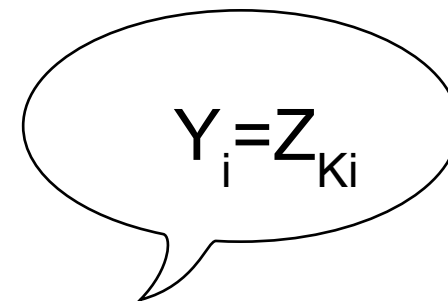
- need finite prior $P_N(\mathbf{a}, \mathbf{b})$ (use truncations as in section 3.2)
- update blocks of parameters (draw from multivariate distr)

$$(X_i | \mathbf{Z}, \mathbf{K}, \theta) \stackrel{\text{iid}}{\sim} \pi(X_i | Z_{K_i}, \theta), \quad i = 1, \dots, n,$$

$$(K_i | \mathbf{p}) \stackrel{\text{iid}}{\sim} \sum_{k=1}^N p_k \delta_k(\cdot),$$

$$(\mathbf{p}, \mathbf{Z}) \sim \pi(\mathbf{p}) \times H^N(\mathbf{Z}),$$

$$\theta \sim \pi(\theta),$$



where $\mathbf{K} = (K_1, \dots, K_n)$, $\mathbf{Z} = (Z_1, \dots, Z_N)$, $\mathbf{p} = (p_1, \dots, p_N) \sim \mathcal{GD}(\mathbf{a}, \mathbf{b})$, and Z_k are iid H .

Blocked Gibbs Sampler

direct posterior inference:

iterate

$$(\mathbf{Z}|\mathbf{K}, \theta, \mathbf{X}),$$

$$(\mathbf{K}|\mathbf{Z}, \mathbf{p}, \theta, \mathbf{X})$$

$$(\mathbf{p}|\mathbf{K}),$$

$$(\theta|\mathbf{Z}, \mathbf{K}, \mathbf{X}).$$

equilibrium
distribution

$$(\mathbf{Z}, \mathbf{K}, \mathbf{p}, \theta|\mathbf{X})$$

each draw gives
a random measure

$$P(\cdot) = \sum_{k=1}^N p_k \delta_{Z_k}(\cdot),$$

Blocked Gibbs Sampler

(a) Conditional for \mathbf{Z} : Simulate $Z_k \stackrel{\text{iid}}{\sim} H$ for each $k \in \mathbf{K} - \{K_1^*, \dots, K_m^*\}$. Also, draw $(Z_{K_j^*} | \mathbf{K}, \theta, \mathbf{X})$ from the density

$$f(Z_{K_j^*} | \mathbf{K}, \theta, \mathbf{X}) \propto H(dZ_{K_j^*}) \prod_{\{i: K_i = K_j^*\}} f(X_i | Z_{K_j^*}, \theta),$$

$$j = 1, \dots, m. \quad (18)$$

(b) Conditional for \mathbf{K} : Draw values

$$(K_i | \mathbf{Z}, \mathbf{p}, \theta, \mathbf{X}) \stackrel{\text{iid}}{\sim} \sum_{k=1}^N p_{k,i} \delta_k(\cdot), \quad i = 1, \dots, n,$$

where

$$(p_{1,i}, \dots, p_{N,i}) \propto (p_1 f(X_i | Z_1, \theta), \dots, p_N f(X_i | Z_N, \theta)).$$

(c) Conditional for \mathbf{p} : By the conjugacy of the generalized Dirichlet distribution to multinomial sampling, it follows that our draw is

$$p_1 = V_1^* \quad \text{and} \quad p_k = (1 - V_1^*)(1 - V_2^*) \cdots (1 - V_{k-1}^*) V_k^*,$$

$$k = 2, \dots, N - 1,$$

where

$$V_k^* \stackrel{\text{iid}}{\sim} \text{Beta}\left(a_k + M_k, b_k + \sum_{l=k+1}^N M_l\right),$$

for $k = 1, \dots, N - 1$,

(d) Conditional for θ : As before, draw θ from the density (remembering that $Y_i = Z_{K_i}$)

$$f(\theta | \mathbf{Z}, \mathbf{K}, \mathbf{X}) \propto \pi(d\theta) \prod_{i=1}^n f(X_i | Y_i, \theta).$$

(a) ~ acceleration step (c)

(d) = (b) in PG and PG_a

Evaluation

- $\{DP, DP_{50}, PY\} \times \{PG, PG_a, BG\}$
- experimental results (batch means, std, etc)

- PG is bad
- PG_a works well when prediction rule is known
- BG is more flexible

- complexity? linear in n (PG) vs multivariate draws (BG)
- easy to get stuck?
- which sampler for which prior?

Discussions

- PG : PG_a is the same as Algorithm 1 : 2 in [Neal 1999]
- what's semiparametric?
- what's "almost sure"?
- what's the graphical model like for (16) in Section 5?
- in the non-conjugate case, are we doing Metropolis-Hastings inside the Gibbs sampler? (Section 5.4)
- with known prediction rule, is BG preferred? How?
- sampling in equivalence class space vs label space?