

Slice sampling in nested IBP

Jinsan Yang, Jinseok Nam and Byoung-Tak Zhang
Biointelligence Laboratory
School of Computer Science and Engineering
Seoul National University
Seoul 151-742, Korea

November 30, 2011

Abstract

We develop a nonparametric Bayesian method that explores the infinite space of latent features and finds the best subset in the sense of posterior probability. When the data appear in several groups, there should be different measures reflecting the differences between the groups. We formalize this as a nested Indian buffet process (nIBP) by assuming different measures according to the specific types where corresponding set of groups belong to. For efficiently running the Gibbs sampling, slice sampling method is applied to our model with proper choice of stick components. Our contributions here are two fold. First we extended the Indian buffet process model by imposing nested structure above the groups of data sets and hence providing a hierarchically clustering method for objects having multiple features like images in nonparametric approach without imposing class numbers or feature numbers. We followed similar formulation as in nDP model except extending the IBP model instead of the DP model (Rodriguez, Dunson & Gelfand, 2008). Secondly in the computation of the stick breaking expressions, we applied the slice sampling method used in IBP model for properly expanding the stick components (Teh, Gorur & Ghahramani, 2007).

1 Introduction

In the analysis of the observed objects characterized by a certain set of latent features, the Indian buffet process (IBP) method is widely applied (Griffiths & Ghahramani, 2005). In IBP, the number of classes or features is assumed potentially unbounded and the observed objects reflect a subset of these features (Rasmusen & Ghahrahmani, 2001). Instead of using a fixed set of latent features, each observed object is modeled very generally using unbounded number of latent features selected based on the beta process (Thibaux & Jordan 2007). Hierarchical Dirichlet process (HBP) models assume multiple groups of data with observed objects in each group. In HDP, each group follows a Dirichlet

process with common measure and this common measure is again comes from a Dirichlet process with a base measure (Teh, Jordan, Beal & Blei, 2006). Therefore, different groups of data share clusters from a common set of clusters.

Recently, several approaches have been proposed to combine feature selections and data clusterings together. DP-IBP and IBP-IBP models (Doshi-Velez, & Ghahramani, 2009) assign different feature sets to each data cluster by assuming single (DP-IBP) or multiple (IBP-IBP) cluster memberships to each data observation. These models consider data clusterings and feature selections for each data together and provide models for the correlated features.

In the nested Dirichlet process (Rodriguez, Dunson & Gelfand, 2008), two stage clustering is considered, where clustering separates distributions of each data group as well as observed data in each nested group. The nested Dirichlet process looks similar to the nested Chinese restaurant process (Blei, Griffiths, Jordan & Tenenbaum, 2004) in selecting clusters hierarchically in each level but in nCRP, individual data point is provided without group memberships and the selection is repeated to make a path over the tree structure for each observed data.

We consider the case when there are multiple groups of data and classifying each group depends on the different set of latent features used to characterize each data group. When there are several groups of observed objects, we consider the clustering of the groups as well as the finding of latent features for each observed objects nested in each group. Therefore, the measure for modeling the set of observed objects are sampled from a Dirichlet process and the set of latent features are selected nested in that measure.

2 Nested IBP model

2.1 Nonparametric Bayesian Models

In nonparametric Bayesian models, the model complexity is unbounded and the prior over the underlying distribution is not limited by the parametric family but includes the space of all distributions. The Chinese restaurant process (CRP) and the Indian buffet process (IBP) are two of the most popular constructive processes in nonparametric Bayesian models relating to Dirichlet process and beta process, respectively (Teh, Jordan, Beal & Blei, 2006; Griffiths & Ghahramani, 2005).

$G \sim DP(\alpha G_0)$ is a Dirichlet process where α is a concentration parameter, G_0 is a base measure over some probability space Θ and G is a random measure over any partition A_1, A_2, \dots, A_K of Θ with distributional property, $(G(A_1), \dots, G(A_K)) \sim dir(\alpha G_0(A_1), \dots, \alpha G_0(A_K))$. Since G is a distribution over Θ , the parameters could be sampled from G .

The Indian buffet process is a constructive sampling process for generating sets of feature selections based on the beta process (Griffiths & Ghahramani, 2005; Thibaux & Jordan, 2007) formulated as: $Z \sim BeP(B)$, $B \sim BP(cB_0)$ where c is a concentration parameter, B_0 is a base measure over latent feature space

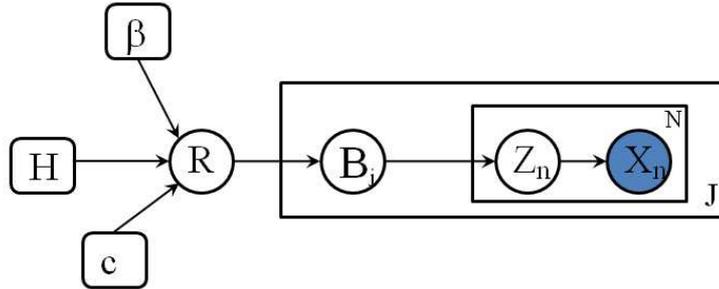


Figure 1: Graphical representation of the nested IBP model. Nodes represent the fixed hyperparameters (rounded cornered squares), variables (circles) and an observable (shaded circle). Rectangles denote the replications with the replication number on the corner. Arrows between the nodes represent the dependency relations.

Ω and Z is a binary row of feature selections sampled from a Bernoulli process with parameters depending on B . B is a random measure with distributional property, $B(A_k) \sim \text{beta}(cB_0(A_k), c(1 - B_0(A_k)))$ independently for any $A_k \subset \Omega$. Stick breaking construction (Sethuraman, 1994) provides an explicit method to construct a sampled measure: For G sampled from $DP(\alpha G_0)$, $G = \sum_{k=1}^{\infty} \pi_k \delta(\theta_k^*)$ where $\delta(\theta_k^*)$ is a point measure concentrated on a representing element θ_k^* for the k^{th} partition of Θ , $\beta_k \sim \text{Beta}(1, \alpha)$, $\pi_k = \beta_k \prod_{i=1}^{k-1} (1 - \beta_i)$ and $\theta_k^* \sim G_0$ for each k . Stick breaking construction is very useful in computing the likelihoods in the case of nested DP model due to the expressive representation of arbitrary distributions (Rodriguez, Dunson & Gelfand, 2008).

2.2 Nested Indian buffet process

When the data are observed in several groups, in addition to the finding of latent feature sets in each group, it is also important to categorize each group of objects according to their properties. For example, when there are sets of images and we want to know which image sets can be characterized based on the same sets of latent features in common. We can consider to categorize the groups of images according to their characteristics. Nested Indian buffet process is a model for clustering the general classes as well as finding the set of latent features for each observed object. Nested IBP model should not be confused with the DP-IBP model (Doshi-Velez, & Ghahramani, 2009) where assuming feature correlations, each object is associated with a finite number of categories and each category is associated with a finite number of features.

Denoting $Z_n^{(j)}$ for a set of latent features corresponding to the n^{th} object in the j^{th} group, we can formulate the nested IBP for $n = 1, \dots, n_j$ and $j = 1, \dots, J$ as follows:

$$Z_n^{(j)} | B_j \sim \text{BeP}(B_j) \quad (1)$$

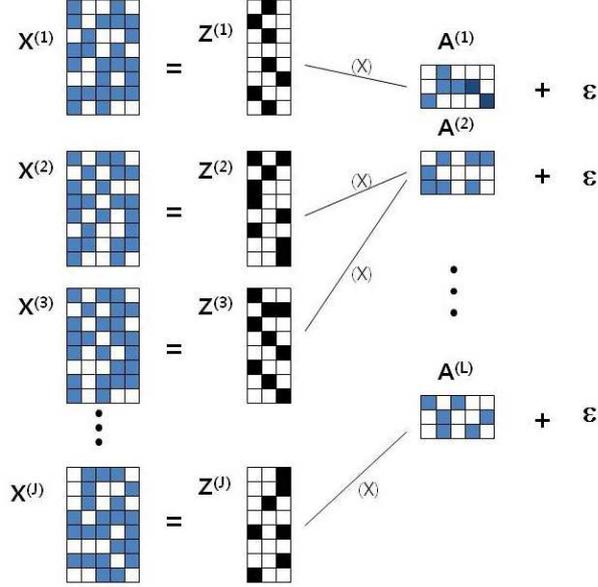


Figure 2: In nested IBP model, each data $X^{(j)}$ is the product of $Z^{(j)}$ and $A^{(l)}$ in one of the latent types with some noise.

$$B_j \sim R \sim DP(\beta\nu) \quad (\nu \equiv BP(cH)) \quad (2)$$

In (1), each element of binary vector $Z_n^{(j)}$ is sampled from a Bernoulli distribution with success probability determined by a measure B_j over parameter space. In (2), the measure B_j for generating features among the j^{th} group is distributed according to R which is also a measure over $\mathcal{C} = \{B\}$ (\mathcal{C} is a collection of measures like B_j). ν is a measure for generating R with positive concentration parameter β and set to be a beta process with base measure H with a positive concentration parameter c . If B_k^* is a representing member of the k^{th} category in \mathcal{C} , it can be sampled from ν in the stick breaking representation. We will show detailed representations in the next section.

3 Stick breaking construction in nested IBP by slice sampling

3.1 Stick breaking construction

Applications of stick breaking method in the defining relations (1) and (2) can provide expressions of R as well as B_k^* , a representing element in the k^{th} category

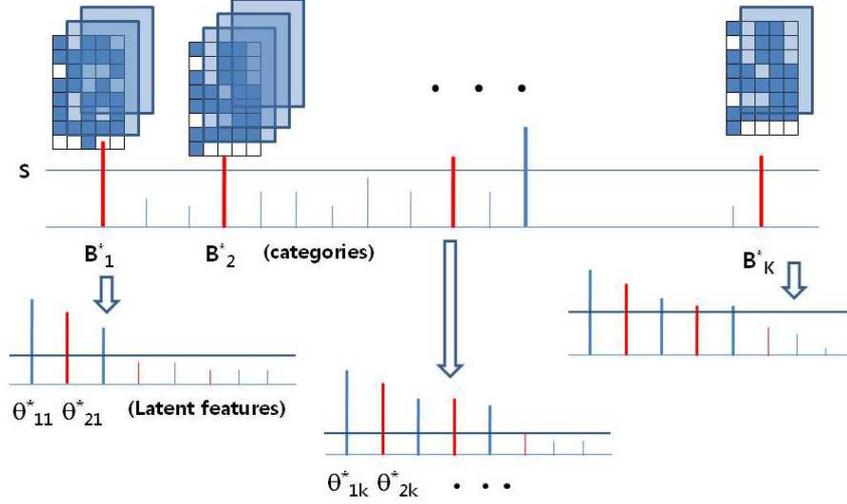


Figure 3: In nested IBP model, each data set is allocated to the constructed stick components taller than the slice value and also each data entry is taking latent features among the feature components with larger weights than the given slice values enabling one to consider only finitely many components out of infinitely many components.

of \mathcal{C} for $k = 1, 2, \dots$. Denoting $\delta(x)$ as a point mass concentrated to a point x :

$$R = \sum_{k=1}^{\infty} \pi_k \delta(B_k^*), \quad \pi_k = u_k \prod_{i=1}^{k-1} (1 - u_i), \quad u_i \sim \text{beta}(1, \beta), \quad B_k^* \sim \text{BP}(cH) \quad (3)$$

$$B_k^* = \sum_{l=1}^{\infty} w_{lk} \delta(\theta_{lk}^*), \quad w_{(l)k} = v_{(l)k} w_{(l-1)k} = \prod_{j=1}^l v_{(j)k}, \quad v_{(j)k} \sim \text{beta}(c, 1), \quad \theta_{lk}^* \sim H \quad (4)$$

where H is a discrete uniform base measure, $w_{(1)k} > w_{(2)k} > \dots > w_{(L)k}$ are rearrangements of weights $w_{1k}, w_{2k}, \dots, w_{Lk}$ in the k^{th} category and θ_{lk}^* is the l^{th} feature of the k^{th} category (Teh, Gorur & Ghahramani, 2007). Since there are infinite terms of weights in each representation, the exact procedure is infeasible and truncation is an easy implementing method using the approximation model instead of the full model. But by introducing the slice sampling (Kalli, Griffin, & Walker, 2011), it is possible to avoid infinite number of weights and still take advantages of the full model. In the following section, we will explain the use of slice sampling method in our model.

3.2 Slice sampling

Slice sampling (Neal, 2003) is an efficient way of sampling by introducing a latent (slice) variable that can regulate sampling process adaptively. For stick breaking construction, slice sampling can regulate the number of sticks adaptively without assuming alternative approximate models (Kalli, Griffin, & Walker, 2011; Teh, Gorur & Ghahramani, 2007). We apply this sampling method in our nested IBP model for efficient computation. From (3), if we use χ for arbitrary data set $X^{(j)}$, the distributional form of χ can be expressed as:

$$f_{\pi, \mathbf{B}^*}(\chi) = \sum_{j=1}^{\infty} \pi_j p(\chi | B_j^*) \quad (5)$$

Now, by introducing a latent variable s to the distributional expression (5) (Kalli, Griffin, & Walker, 2011):

$$f_{\pi, \mathbf{B}^*}(\chi, s) = \sum_{j=1}^{\infty} \mathbb{I}(s < \pi_j) p(\chi | B_j^*) \quad (6)$$

The conditional distribution of choosing sticks given a slice value becomes 0 for sticks smaller than the slice value:

$$f_{\pi, \mathbf{B}^*}(k | \chi, s) \propto \mathbb{I}(s < \pi_k) p(\chi | B_k^*) \quad (7)$$

We introduce another slice variable r for expanding sticks in each category of data groups,

$$f_{\mathbf{w}, \boldsymbol{\theta}^*}(z_{ik}, r) = \mathbb{I}(r < w_{(L^\dagger)}) w_{(k)} \quad \text{for } k = 1, 2, \dots \quad (8)$$

where L^\dagger is the index of the last expanded stick for each k . The conditional probability of choosing sticks becomes 0 for sticks smaller than the given slice value.

The conditional distributions of newly made sticks in each category are sampled from,

$$p(w_{(l)} | w_{(l-1)}, rest) \propto \exp\left(c \sum_{i=1}^N \frac{1}{i} (1 - w_{(l)})^i\right) w_{(l)}^{c-1} (1 - w_{(l)})^N \mathbb{I}(0 < w_{(l)} < w_{(l-1)}) \quad (9)$$

Existing stick weights are updated by

$$p(w_{(l)} | rest) \propto w_{(l)}^{m_{\cdot l}} (1 - w_{(l)})^{N - m_{\cdot l}} \mathbb{I}(w_{(l+1)} < w_{(l)} < w_{(l-1)}) \quad (10)$$

where $m_{\cdot l} = \sum_i^N z_{il}$.

For a given r , the conditional probability of z_{il} is,

$$p(z_{il} = 1 | rest) \propto \frac{w_{(l)}}{w^*} f(x_i | z_{il} = 1, rest) \quad (11)$$

{Slice sampling for the categories of the data group}

- For each data set $x_{:,j}$, sample s uniformly by the length of the lastly used stick.
- Expand the sticks until $s + \sum_{l=1}^{K^\dagger} \pi_l > 1$ (where K^\dagger is the index of the lastly selected stick)
- Assign $x_{:,j}$ to one of the components above the slice using (7) and denote its index by k_j
- Repeat above process for each of the $j = 1, \dots, J$ data groups.

{Slice sampling for the latent features of data in the k^{th} category}

- For each row Z_n in the k^{th} category, sample a slice r uniformly between 0 and the last active stick weight in the k^{th} category.
- Expand the sticks until $r > \mu_{(L^\dagger)}$ using (9)
- Sample parameters for the newly expanded features from prior
- For all the existing sticks, update z_{il} using (11)
- Repeat above process for each row in the k^{th} category
- Update all of the parameters corresponding to the expanded sticks
- Update all the stick weights using (10)

4 Conclusions and discussions

We proposed a model for analyzing large scales data sets using both clustering and feature allocation methods with effective sampling method. We will test our model with various image data for further properties and possibilities.

References

- [1] Antoniak, C. (1974). Mixtures of Dirichlet processes with applications to Bayesian nonparametric problems. *The Annals of Statistics*, 2, 1152-1174.
- [2] Blei, D. M., Griffiths, T. L., Jordan, M. I., and Tenenbaum, J. B. (2004). Hierarchical topic models and the nested Chinese restaurant process. In *Advances in Neural Information Processing Systems*, volume 16, Cambridge, MA. MIT Press.
- [3] Doshi-Velez, F. & Ghahramani, Z. (2009) Correlated non-parametric latent feature models In: *UAI 2009 (Conference on Uncertainty in Artificial Intelligence)*, 18-21 June 2009, Montreal, Quebec.
- [4] Escobar, M. D., & West, M. (1995). Bayesian density estimation and inference using mixtures. *Journal of the American Statistical Association*, 90, 577-588.
- [5] Ferguson, T. S. (1983). Bayesian density estimation by mixtures of normal distributions. In M. Rizvi, J. Rustagi, & D. Siegmund (Eds.), *Recent advances in statistics* (p. 287-302). New York: Academic Press.

- [6] Ishwaran, H., & James, L.F. (2001). Gibbs sampling methods for stick-breaking priors. *Journal of the American Statistical Association*, 96(453):161-173.
- [7] Kalli, M., Griffin, J.E., & Walker, S.G. (2011). Slice sampling mixture models. *Journal of Statistics and Computing*, Vol. 21, Issue 1.
- [8] Griffiths, T., & Ghahramani, Z. (2005). Infinite latent feature models and the Indian buffet process. TR 2005-001, Gatsby Comp. Neuroscience Unit.
- [9] McLachlan, G. J. and Peel, D. (2000). *Finite Mixture Models*, Wiley.
- [10] Neal, R. M. (2000). Markov chain sampling methods for Dirichlet process mixture models. *Journal of Computational and Graphical Statistics*, 9, 249-265.
- [11] Neal, R. M. (2003). Slice sampling, *Ann. Statist.* Volume 31, Number 3, 705-767.
- [12] Paisley, J., Zaas, A., Woods, C. W., Ginsburg, G. S., & Carlin, L. (2010). A Stick-Breaking Construction of the Beta Process. ICML2010, Haifa, Israel.
- [13] Rodriguez, A., Dunson, D. B., Gelfand, A. E. (2008). The nested Dirichlet process. *Journal of the American Statistical Association*, Vol. 103, No. 483, pp. 1131-1154.
- [14] Sethuraman, J. (1994). A constructive definition of Dirichlet priors. *Statistica Sinica*, 4:639-650.
- [15] Teh, Y. W., Jordan, M. I., Beal, M. J., and Blei, D. M. (2006). Hierarchical Dirichlet processes. *Journal of the American Statistical Association*, 101(476):1566-1581.
- [16] Teh, Y.W., Gorur, D., Ghahramani, Z. (2007). Stick-breaking Construction for the Indian Buffet Process. *Proceedings of the International Conference on Artificial Intelligence and Statistics*, Vol. 11.
- [17] Thibaux, R. & Jordan, M. I. (2007). Hierarchical beta processes and the Indian buffet process. In *Proceedings of the International Workshop on Artificial Intelligence and Statistics*, volume 11.