
Variational Inference for Dirichlet Process to Stratify Cancer Patients Using DNA Methylation

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Variational Inference (VI) is an alternative strategy to Markov Chain Monte Carlo
2 (MCMC) which tends to be faster and easier to scale to larger datasets (Blei
3 et al., 2016). This is especially advantageous in applications that interact with
4 high dimensional data. Previous work has been done on a VI method for Dirichlet
5 Process Mixture Models (DPMMs) based on the well-known stick-breaking process
6 (Blei et al., 2016). In contrast, we propose a similar model, but based on the Chinese
7 Restaurant Process (CRP). We hypothesize that our model will perform faster
8 because it has fewer parameters to estimate. To test our hypothesis, we present
9 an experiment on a large-scale stratification problem using DNA methylation to
10 compare both implementations.

11 1 Introduction

12 1.1 Variational Inference

13 Variational inference is an alternative strategy to Markov Chain Monte Carlo (MCMC) sampling
14 which tends to be faster and easier to scale to larger datasets (Blei et al., 2016). The key characteristic
15 of variational inference is that it casts Bayesian inference as an optimization problem (Salimans
16 et al., 2015). Variational inference attempts to approximate the posterior with another distribution
17 $q_\theta(z|x)$ by choosing its parameters θ to optimize the evidence lower bound (ELBO) on the marginal
18 likelihood,

$$\begin{aligned}\log p(x) &\geq \log p(x) - D_{KL}(q_\theta(z|x)||p(z|x)) \\ &= E_{q_\theta(z|x)}(\log p(x, z) - \log q_\theta(z|x))\end{aligned}$$

19 In recent years, there have been many advances in the field of VI, which are aptly summarized in a
20 review by Zhang et al., (2019).

21 1.2 Dirichlet Process and Chinese Restaurant Process

22 The Dirichlet process is a stochastic process used in Bayesian nonparametrics; one speci, for con-
23 structing Dirichlet process mixture models (Neal, 2000). A Dirichlet process G is a distribution of
24 distributions consisting of a base distribution G_0 and a positive real number α , written as,

$$G \sim \text{DP}(G_0, \alpha)$$

where G_0 is a continuous distribution such that the probability of any two samples generated from this distribution being equal is zero, whereas G is a discrete distribution consisting of infinitely many number of point masses, so the probability of two samples colliding is non-zero. For any measurable finite k partitions $\{B_i\}_{i=1}^k$, if,

$$G(B_1), \dots, G(B_k) \sim \text{Dir}(\alpha G(B_1), \dots, \alpha G(B_k))$$

then,

$$\Pr[X_1, \dots, X_k] = \int \Pr[G] \prod_{i=1}^k \Pr[X_i | G] dG$$

which represents the dependencies among $\{X_i\}_{i=1}^k$ by marginalizing out G .

The Chinese restaurant process, stick-breaking process and Polya urn scheme are three common perspectives regarding the Dirichlet Process; we restrict our attention only to the Chinese restaurant process perspective in this project. Assume a restaurant has infinitely many tables, and let $\{X_i\}_{i=1}^k$ be the customers of the restaurant. The $\{X_i\}_{i=1}^k$ are partitioned based on which table each customer is seated at. Consider the behavior of a single customer X_n given $\{X_i\}_{i=1}^{n-1}$:

$$X_n | (X_1 = x_1, X_2 = x_2, \dots, X_{n-1} = x_{n-1}) = \begin{cases} x_n^* & \text{with probability } \frac{|\{j: x_j = x_n\}|}{n-1+\alpha} \\ \text{new draw from } G_0 & \text{with probability } \frac{\alpha}{n-1+\alpha} \end{cases}$$

where $|\{j: x_j = x_n\}|$ is the number of times the value x_n occurs in $\{x_1, x_2, \dots, x_n\}$.

The intuition is that customers are more likely to sit at tables with more customers, but will sit at a new table with a probability proportional to α .

1.3 Dirichlet Process Mixture Model

The Dirichlet Process Mixture Model (DPMM) is a hierarchical model for classifying data points into an unbounded number of mixture components. Given a sample $\{x_1, \dots, x_N\}$, the aim of a DPMM is to compute the posterior predictive distribution,

$$\Pr[X = \hat{x} | x_1, \dots, x_N, \alpha, G_0] = \int \Pr[\hat{x} | x] \Pr[x | x_1, \dots, x_N, \alpha, G_0] dx,$$

The posterior distribution $\Pr[x | x_1, \dots, x_N]$ does not have a closed form. Since there is an unbounded number of mixtures, sampling methods are commonly used to estimate the posterior. The most popular methods include MCMC and VI.

1.4 DNA Methylation

Cancers develop via the acquisition of genomic changes. These changes in turn alter the cells harbouring them, leading to changes in cell state (phenotype). One common effect of these mutations is to induce a series of modifications to the genome, which do not alter the encoded DNA but rather the ability of the DNA to be read and processed into protein. These heritable non-genetic changes are broadly referred to as epigenetic changes. According to Baylin and Jones (2011), "Epigenetic alterations are leading candidates for the development of specific markers for cancer detection, diagnosis and prognosis". One such epigenetic change is DNA methylation, which is unambiguously linked with transcriptional repression. When present in promoter regions, DNA methylation correlates negatively with gene expression; furthermore, gene promoter CpG islands acquire abnormal hypermethylation resulting in transcriptional silencing in cancer (Bock, 2012). DNA methylation can be used to stratify cancer patients into functionally similar groups. These groups can elucidate shared altered gene pathways to inform treatment strategies.

61 2 Related Work

62 Blei and Jordan (2006) introduced a VI algorithm for DPMMs based on the stick-breaking process.
 63 They compared the mean convergence time consumption of their VI algorithm to two MCMC
 64 sampling algorithms, i.e., Collapsed Gibbs and Blocked Gibbs. We now provide a brief description
 65 of their algorithm. Let,

$$\Pr[G_0 = x^* \mid \lambda] = h(x^*) \exp(\lambda_1^\top x^* + \lambda_2(-a(x^*)) - a(\lambda))$$

66 where λ 's are hyperparameters, and $a(x^*)$ is a cumulant function. Then the target function for
 67 prediction, which is based on the stick-breaking process, is

$$\Pr[x_n \mid z_n, x_1^*, x_2^*, \dots] = \prod_{i=1}^{\infty} [h(x_n) \exp(x_i^{*\top} x_n - a(x_i^*))]^{\mathbf{1}[z_n=i]}$$

68 This is an intractable distribution, so Blei and Jordan use VI, along with the mean-field variational
 69 approximations assuming the independence of latent variables and a derived coordinate ascent
 70 algorithm. They arrive at the following expression,

$$\Pr[x_{N+1} \mid x_1, x_2, \dots, x_{N-1}, \alpha, \lambda] \approx \sum_{t=1}^T \mathbb{E}[\pi_t(\mathbf{v})] \mathbb{E}[\Pr[x_{N+1} \mid x_t^*]]$$

71 where q is an approximation to the predictive distribution depending on $x_1, x_2, \dots, \alpha, \lambda$, and each
 72 component v_i in \mathbf{v} follows the beta distribution with parameters $(1, \alpha)$.

73 3 Methods

74 3.1 Model

75 This paper builds upon Blei and Jordan's work (2006). We propose a VI algorithm for DPMMs based
 76 on a CRP, to stratify cancer patients using DNA methylation. Our method differs in that we base our
 77 model on a CRP as opposed to a stick-breaking process. We will use the probabilistic programming
 78 language Pyro to implement both models (see Figure 1 for graphical models).

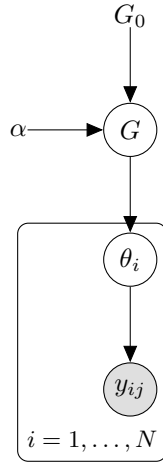


Figure 1: Graphical model for DPMM based on CRP

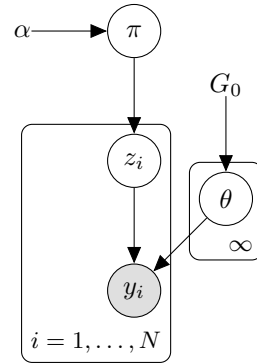


Figure 2: Graphical model for DPMM based on the stick-breaking process

79 3.2 Data Preprocessing

80 DNA methylation data will be obtained from the International Cancer Genome Consortium (ICGC)
 81 and processed in R. Preprocessing steps will include:

- 82 • Filtering / Cleaning
- 83 • Probe to gene mapping using FDb.InfiniumMethylation.hg19 (Triche, 2014)
- 84 • Normalization
- 85 • Dimensionality reduction

86 4 Expected Results

87 4.1 Clustering

88 We intend to use clustermaps where color bars on the left hand side indicates cancer type and purple
 89 squares represent the clusters. The more data points in the cluster, the larger the squares. The intensity
 90 of a square describes the proportion of iterations the data point of interest appeared in that cluster.
 91 We expect that our clustermaps will resemble Figure 3, as patients should cluster by cancer type.

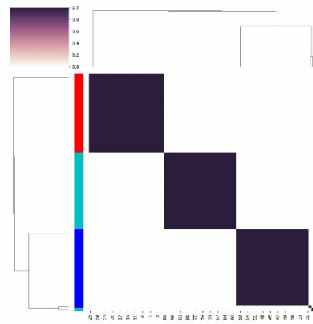


Figure 3: Example of a clustermap. Color labels indicate cancer types: Type I, Type II, Type III

92 4.2 Mean Convergence Time

93 We expect to observe that our model, based on a CRP will have a shorter mean convergence time
 94 compared to Blei and Jordan's model. Plots depicting the relationship between dimension, time, and
 95 algorithm type will be used to visualize and evaluate our hypothesis. We expect that our plot will
 96 look similar to Figure 4, and will show that our model performs faster than Blei and Jordan's.

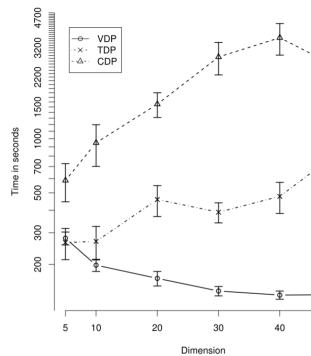


Figure 4: Mean convergence time and standard error across ten data sets per dimension for variational inference, TDP Gibbs sampling, and the collapsed Gibbs sampler (Blei and Jordan, 2006)

97 References

- 98 [1] Baylin, Stephen B., and Peter A. Jones. A Decade of Exploring the Cancer Epigenome – Biological and
 99 Translational Implications. *Nature Reviews Cancer*, vol 11, no. 10, 23 Sept. 2011, pp. 726-734, 10.1038/nrc3130.

- 100 [2] Blei, David M, et al. Variational Inference: A Review for Statisticians. *Journal of the Ameri-*
101 *can Statistical Association*, vol. 112, no. 518, 27 Feb. 2017, pp. 859-877, arxiv.org/abs/1601.00670,
102 [10.1080/01621459.2017.1285773](https://doi.org/10.1080/01621459.2017.1285773).
- 103 [3] Blei, David M., and Michael I. Jordan. Variational Inference for Dirichlet Process Mixtures. *Bayesian*
104 *Analysis*, vol.1, no.1, Mar.2006, pp. 121-143, [10.1214/06-ba104](https://doi.org/10.1214/06-ba104).
- 105 [4] Bock, Christoph, et al. Quantitative Comparison of Genome-Wide DNA Methylation Mapping Technologies.
106 *Nature Biotechnology*, vol. 28, no. 10, 19 Sept. 2010, pp. 1106-1114, [10.1038/nbt.1681](https://doi.org/10.1038/nbt.1681).
- 107 [5] Salimans, Tim, et al. Markov Chain Monte Carlo and Variational Inference: Bridging the Gap. ArXiv:
108 [1410.6460](https://arxiv.org/abs/1410.6460), 19 May 2015, arxiv.org/abs/1410.6460.
- 109 [6] Zhang, Cheng, et al. Advances in Variational Inference. *IEEE Transactions on Pattern Analysis and Machine*
110 *Intelligence*, vol. 41, no. 8, 1 Aug. 2019, pp. 2008-2026, [10.1109/tpami.2018.2889774](https://doi.org/10.1109/tpami.2018.2889774).
- 111 [7] Triche, Jr. T (2014). FDb.InfiniumMethylation.hg19: Annotation package for Illumina Infinium DNA
112 methylation probes. R package version 2.2.0.