
Dirichlet Process

Xiaobao Wu
Tsinghua University
xiaobao.wu1996@gmail.com

In general, Dirichlet Process is not as straightforward as Gaussian Process. The definition and construction are extremely opaque even after perusing. In this note, we conclude some main points of Dirichlet Process to explain it in a simple and direct way.

1 Motivation

Dirichlet Process is a method of Nonparametric Bayesian. One simple motivation of Dirichlet Process is to determine the number of clusters in mixture models like GMM. If the parameters of each data point are drawn from a continuous distribution, the probability is zero that two data points are from a same distribution and the number of clusters is the size of data points. Therefore, we need a discrete function to handle this issue.

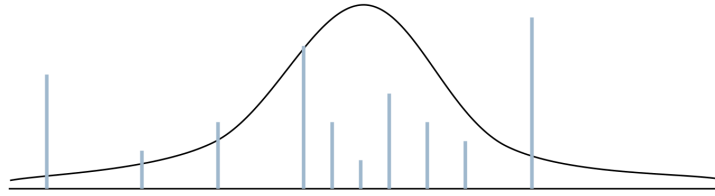


Figure 1: Illustration of G and H in Dirichlet Process.

2 Definition of Dirichlet Process

$$G \sim \text{DP}(\alpha, H) \quad (1)$$

$$\Leftrightarrow (G(a_1), G(a_2), \dots, G(a_K)) \sim \text{Dir}(\alpha H(a_1), \alpha H(a_2), \dots, \alpha H(a_K)) \quad (2)$$

H is called the base function. As shown in Figure 1, G can be considered as the discrete version of H where α determines the degree. If $\alpha = 0$, only one probability measure equal to 1 exists. If $\alpha = \infty$, G is H . G is also called the random probability measure. a_1, \dots, a_K are the partitions of H . $G(a_i)$ means the measure sum of G in partition a_i and $H(a_i)$ means the measure sum of H in partition a_i . The expectation and variance are

$$E[G(a_i)] = H(a_i) \quad (3)$$

$$\text{Var}[G(a_i)] = \frac{H(a_i)(1 - H(a_i))}{\alpha + 1} \quad (4)$$

3 Stick Breaking Construction

We consider the variable π_i where

$$\theta_i \sim H \quad (5)$$

$$\beta_i \sim \text{Beta}(1, \alpha) \quad (6)$$

$$\pi_1 = \beta_1 \quad (7)$$

$$\pi_i = \beta_i \prod_{k=1}^{i-1} (1 - \beta_k) \quad (8)$$

π_i is the probability measure of point θ_i . Here we construct the distribution G as

$$G(\cdot) = \sum_{i=1}^{\infty} \pi_i \delta_{\theta_i}(\cdot) \quad (9)$$

δ is the Dirac delta function, also called point mass function where $\delta_{\theta_i}(\theta) = 1$ if $\theta = \theta_i$ and 0 otherwise.

4 Posterior Dirichlet Process

$$G \sim DP(\alpha, H) \quad (10)$$

$$\theta_1, \dots, \theta_N \stackrel{i.i.d.}{\sim} G \quad (11)$$

We need to calculate the posterior as

$$p(G|\theta_1, \dots, \theta_N) \propto p(\theta_1, \dots, \theta_N|G)p(G)$$

where $G = p(\theta_1, \dots, \theta_N|G)$. However, it is inconvenient to directly use the Dirichlet Process as $p(G)$. In another way, we can turn to the conjugate property of Multinomial and Dirichlet distribution. For any partition (a_1, \dots, a_K) , n_i is the number of $\theta_j, j = 1, \dots, N$ in a_i . We have

$$p(G(a_1), \dots, G(a_K)|n_1, \dots, n_K) \quad (12)$$

$$\propto \text{Mult}(n_1, \dots, n_K|G(a_1), \dots, G(a_K)) \text{Dir}(\alpha H(a_1), \dots, \alpha H(a_K)) \quad (13)$$

$$= \text{Dir}(\alpha H(a_1) + n_1, \dots, \alpha H(a_2) + n_2) \quad (14)$$

$$= DP(\alpha + N, \frac{\alpha H + \sum_{i=1}^N \delta_{\theta_i}}{\alpha + N}) \quad (15)$$

Thus, the posterior of a Dirichlet Process is also a Dirichlet Process.

5 Chinese Restaurant Process

$n_{l,-i}$ means the number of data points in class l except i .

$$p(z_i = m|\mathbf{z}_{-i}) \quad (16)$$

$$= \frac{p(z_i = m, \mathbf{z}_{-i})}{p(\mathbf{z}_{-i})} \quad (17)$$

$$= \frac{\int p(z_i = m, \mathbf{z}_{-i}|p_1, \dots, p_K) \text{Dir}(p_1, \dots, p_K | \frac{\alpha}{K}, \dots, \frac{\alpha}{K}) d(p_1, \dots, p_K)}{\int p(\mathbf{z}_{-i}|p_1, \dots, p_K) \text{Dir}(p_1, \dots, p_K | \frac{\alpha}{K}, \dots, \frac{\alpha}{K}) d(p_1, \dots, p_K)} \quad (18)$$

$$= \frac{\Gamma(n_{m,-i} + \frac{\alpha}{K} + 1) \prod_{l=1, l \neq m}^K \Gamma(n_{l,-i} + \frac{\alpha}{K})}{\Gamma(\alpha + N)} \times \frac{\Gamma(\alpha + N - 1)}{\prod_{l=1}^K \Gamma(n_{l,-i} + \frac{\alpha}{K})} \quad (19)$$

$$= \frac{n_{m,-i} + \frac{\alpha}{K}}{N + \alpha - 1} \quad (20)$$

$$= \frac{n_{m,-i}}{N + \alpha - 1} \quad (\text{If } K \text{ is infinity}) \quad (21)$$

In CRP, if a new customer comes, he will choose an old table with the probability $\frac{\sum_{l=1}^K n_{l,-i}}{N + \alpha - 1} = \frac{N-1}{N + \alpha - 1}$ and a new table with $\frac{\alpha}{N + \alpha - 1}$.