

# Dirichlet process note

Daeyoung Lim\*  
Department of Statistics  
Korea University

May 12, 2016

## 1 Dirichlet Distribution

As Dirichlet process is the infinite-dimensional generalization of the Dirichlet distribution, we should first take a look at what Dirichlet distribution is and what properties it possesses. Afterwards, it will be quite straightforward to make the stochastic process have those properties as we construct it. The Dirichlet distribution of dimension  $D$  is a continuous probability measure on  $\Delta_D$  having the density function

$$p(\boldsymbol{\pi} | \beta_1, \dots, \beta_D) = \frac{\Gamma(\sum_i \beta_i)}{\prod_i \Gamma(\beta_i)} \prod_{i=1}^D \pi_i^{\beta_i-1} \quad (1)$$

where the parameters  $\beta_i \geq 0, \forall i$ . In preparation of the generalization to Dirichlet process, let us reparameterize the distribution as

$$p(\boldsymbol{\pi} | \alpha g_{01}, \dots, \alpha g_{0D}) = \frac{\Gamma(\alpha)}{\prod_i \Gamma(\alpha g_{0i})} \prod_{i=1}^D \pi_i^{\alpha g_{0i}-1} \quad (2)$$

where  $\alpha = \sum_i \beta_i$  and  $g_{0i} = \beta_i / (\sum_i \beta_i)$ . We will hereafter denote the distribution by  $\boldsymbol{\pi} \sim \text{Dir}(\alpha g_0)$ ; the mean and variance of such a reparameterized Dirichlet random variable are

$$\mathbb{E}[\pi_i] = g_{0i}, \text{Var}[\pi_i] = \frac{g_{0i}(1-g_{0i})}{\alpha+1}. \quad (3)$$

## 2 How to interpret Dirichlet distribution

It is a well-known fact that the sum of all elements in a Dirichlet random vector is unity. Therefore, it is not difficult to admit that somehow a Dirichlet random vector is a realization of some sort of discrete distribution with a finite number of possible values. In other words, we also refer to a Dirichlet distribution as a distribution on a probability simplex, which is basically the same thing.

$$X_{n+1} | X_1, \dots, X_n \sim \sum_{i=1}^n \frac{1}{\alpha(X) + n} \delta_{X_i} + \frac{1}{\alpha(X) + n} \alpha \quad (4)$$

## 3 DPM: Rao-Blackwellized MCMC

$c_i$  is the cluster to which  $y_i$  belongs. Fast MCMC algorithm integrates out  $\theta$  and only constructs a chain for the categorical variable  $c$ .

- If  $c = c_j$  for some  $j \neq i$ :

$$P(c_i = c | c_{-i}, y_i) = b \frac{n_{-i,c}}{n-1+M} \int F(y_i, \theta) dH_{-i,c}(\theta). \quad (5)$$

---

\*Prof. Taeryon Choi

- Otherwise,

$$P(c_i \neq c_j \text{ for all } j \neq i \mid c_{-i}, y_i) = b \frac{M}{n-1+M} \int F(y_i, \theta) dG_0(\theta). \quad (6)$$

Now let  $G_0 \equiv N(\mu_0, \sigma_0^2)$  and  $F \equiv N(\theta, \sigma^2)$ .

$$p(\theta \mid \{y_j \mid j \neq i, c_j = c\}, G_0) \equiv N\left(\frac{\sigma_0^2 \sum_{c_j=c} y_j + \mu_0 \sigma^2}{n_c \sigma_0^2}, \frac{\sigma^2 \sigma_0^2}{n_c \sigma_0^2 + \sigma^2}\right) \quad (7)$$

$$\int F(y_i, \theta) dH_{-i,c}(\theta) \equiv N(y_i \mid \mu_\theta, \sigma_\theta^2 + \sigma_y^2) \quad (8)$$

where the mean of (7) is  $\mu_\theta$  in (8) and the same for  $\sigma_\theta^2$ .  $\sigma_y^2$  in (8) is simply  $\sigma^2$  but the subscript was attached to make the distinction clear.

## 4 DPM: Latent variable MCMC

We have marginalized out  $\theta$  in the previous section, which is also called *Rao-Blackwellization*. We examine another form of MCMC which also has a chain with respect to  $\theta$  along with the categorical variable  $c$ .

- If  $c = c_j$  for some  $j \neq i$ ,

$$P(c_i = c \mid c_{-i}, y_i, \theta) = b \frac{n_{-i,c}}{n-1+M} F(y_i, \theta_c). \quad (9)$$

- Otherwise,

$$P(c_i \neq c_j \text{ for all } j \neq i \mid c_{-i}, y_i, \theta) = b \frac{M}{n-1+M} \int F(y_i, \theta) dG_0(\theta). \quad (10)$$

In Neal's paper, the Gibbs sampler construction is summerized as follows:

- For  $i = 1, \dots, n$ : If the present value of  $c_i$  is associated with no other observation (i.e.,  $n_{-i,c_i} = 0$ ), remove  $\theta_{c_i}$  from the state. Draw a new value for  $c_i$  from  $c_i \mid c_{-i}, y_i, \theta$  as defined above. If the new  $c_i$  is not associated with any other observation, draw a value for  $\theta_{c_i}$  from  $H_i$  and add it to the state.  $H_i$  is the posterior distribution based on the prior  $G_0$  and the single observation  $y_i$ .
- For all  $c \in \{c_1, \dots, c_n\}$ : Draw a new value from  $\theta_c \mid$  all  $y_i$  for which  $c_i = c$ — that is, from the posterior distribution based on the prior  $G_0$  and all the data points currently associated with latent class  $c$ .