



Figure 1: Graphical Representation of NMVM Model

Table 1: NOTATIONS

$V$	size of the vocabulary
$D$	number of documents in the corpus
$L$	average length of the documents
$\vec{d}$	documents in the corpus
$\vec{z}$	cluster assignment of the documents
$z_d$	cluster assignment of document $d$
$e_d$	embedding of document $d$
$I$	number of iterations
$m_z$	number of documents in cluster $z$
$n_z$	number of words in cluster $z$
$e_z$	sum of the embedding of documents in cluster $z$
$n_z^w$	number of occurrences of word $w$ in cluster $z$
$N_d$	number of words in document $d$
$N_d^w$	number of occurrences of word $w$ in document $d$

## 1 COLLAPSED GIBBS SAMPLING FOR NMVM

### 1.1 The Probability to Choose an Existed Cluster

In this part, we will give a detailed derivation of the probability to choose an existed cluster using collapsed Gibbs Sampling.

$$\begin{aligned}
 & p(z_d = z | \vec{z}_{-d}, \vec{d}, \alpha, \beta, \lambda, \sigma_0^2, \mu_0, \sigma^2) \\
 &= p(z_d = z, s_d = 1 | \vec{z}_{-d}, \vec{d}, \alpha, \beta, \lambda, \sigma_0^2, \mu_0, \sigma^2) \\
 & \quad + p(z_d = z, s_d = 0 | \vec{z}_{-d}, \vec{d}, \alpha, \beta, \lambda, \sigma_0^2, \mu_0, \sigma^2)
 \end{aligned} \tag{1}$$

$$\begin{aligned}
 &= p(s_d = 1 | \lambda) p(z_d = z | \vec{z}_{-d}, \vec{d}, \alpha, \beta) \\
 & \quad + p(s_d = 0 | \lambda) p(z_d = z | \vec{z}_{-d}, \vec{d}, \sigma_0^2, \mu_0, \sigma^2)
 \end{aligned} \tag{2}$$

$$\begin{aligned}
 &= \lambda p(z_d = z | \vec{z}_{-d}, \vec{d}, \alpha, \beta) \\
 & \quad + (1 - \lambda) p(z_d = z | \vec{z}_{-d}, \vec{d}, \sigma_0^2, \mu_0, \sigma^2)
 \end{aligned} \tag{3}$$

Equation (1) uses the Sum Rule of Probability and Equation (2) uses the property of D-Separation. Now, we derive the first term of Equation (3), the generation of bag-of-words part, as follows:

$$\begin{aligned}
& p(z_d = z | \vec{z}_{-d}, \vec{d}, \alpha, \beta) \\
& \propto p(z_d = z | \vec{z}_{-d}, \vec{d}_{-d}, \alpha, \beta) p(d | z_d = z, \vec{z}_{-d}, \vec{d}_{-d}, \alpha, \beta) \quad (4) \\
& = p(z_d = z | \vec{z}_{-d}, \alpha) p(d | z_d = z, \vec{d}_{z, -d}, \beta) \quad (5)
\end{aligned}$$

Here, we use the Bayes Rule in Equation (4) and apply the property of D-Separation in Equation (5). The first term in Equation (5) is the probability of choosing cluster  $z$  for document  $d$  when the cluster assignments of the other documents have been known. The second term of Equation (5) indicates the predictive probability of document  $d$  when we know  $\vec{d}_{z, -d}$ , i.e., the other documents in cluster  $z$  currently.

Now, we will first derive the first term in Equation (5) as follows:

$$\begin{aligned}
& p(z_d = z | \vec{z}_{-d}, \alpha) \\
& = \int p(z_d = z, \theta | \vec{z}_{-d}, \alpha) d\theta \quad (6)
\end{aligned}$$

$$= \int p(\theta | \vec{z}_{-d}, \alpha) p(z_d = z | \vec{z}_{-d}, \theta, \alpha) d\theta \quad (7)$$

$$= \int p(\theta | \vec{z}_{-d}, \alpha) p(z_d = z | \theta) d\theta \quad (8)$$

Here, we use Sum Rule of Probability in Equation (6) and apply Product Rule of Probability in Equation (7). We finally employ the property of D-Separation in Equation (8). The first term of Equation (8) is the posterior of  $\theta$  and the second term in Equation (8) is a *Categorical* distribution:  $Mult(z_d = z | \theta) = \theta_z$

Now, the following is the derivation of the first term in Equation (8)

$$\begin{aligned}
& p(\theta | \vec{z}_{-d}, \alpha) \\
& = \frac{p(\theta | \alpha) p(\vec{z}_{-d} | \theta)}{\int p(\theta | \alpha) p(\vec{z}_{-d} | \theta) d\theta} \quad (9)
\end{aligned}$$

$$= \frac{\frac{1}{\Delta(\alpha)} \prod_{k=1}^K \theta_k^{\alpha/K-1} \prod_{k=1}^K \theta_k^{m_{k, -d}}}{\int \frac{1}{\Delta(\alpha)} \prod_{k=1}^K \theta_k^{\alpha/K-1} \prod_{k=1}^K \theta_k^{m_{k, -d}} d\theta} \quad (10)$$

$$= \frac{1}{\Delta(\vec{m}_{-d} + \alpha/K)} \prod_{k=1}^K \theta_k^{m_{k, -d} + \alpha/K - 1} \quad (11)$$

$$= Dir(\theta | \Delta(\vec{m}_{-d} + \alpha/K)) \quad (12)$$

Here, we use the Bayes Rule in Equation (9). Then, we can derive the first term in Equation (5) as follows:

$$\begin{aligned}
& p(z_d = z | \vec{z}_{-d}, \alpha) \\
& = \int Dir(\theta | \vec{m}_{-d} + \alpha/K) Mult(z_d = z | \theta) d\theta \quad (13)
\end{aligned}$$

$$= \int \frac{1}{\Delta(\vec{m}_{-d} + \alpha/K)} \theta_z \prod_{k=1}^K \theta_k^{m_{k, -d} + \alpha/K - 1} d\theta \quad (14)$$

$$= \frac{\Delta(\vec{m} + \alpha/K)}{\Delta(\vec{m}_{-(d)} + \alpha/K)} \quad (15)$$

$$= \frac{\prod_{k=1}^K \Gamma(m_k + \alpha/K)}{\Gamma(\sum_{k=1}^K (m_k + \alpha/K))} \frac{\Gamma(\sum_{k=1}^K (m_{k, -d} + \alpha/K))}{\prod_{k=1}^K \Gamma(m_{k, -d} + \alpha/K)} \quad (16)$$

$$= \frac{\Gamma(m_{z, -d} + \alpha/K + 1)}{\Gamma(m_{z, -d} + \alpha/K)} \frac{\Gamma(D - 1 + \alpha)}{\Gamma(D + \alpha)} \quad (17)$$

$$= \frac{m_{z, -d} + \alpha/K}{D - 1 + \alpha} \quad (18)$$

In Equation (15), we employ the  $\Delta$  function, which is defined as  $\Delta(\vec{\alpha}) = \frac{\prod_{k=1}^K \alpha_k}{\Gamma(\sum_{k=1}^K \alpha)}$ . Based on the property of  $\Gamma$  function:  $\Gamma(x + 1) = x\Gamma(x)$ , we can obtain Equation (18) from Equation (17). In Equation (18),  $m_{z, -d}$  denotes the number of documents in cluster  $z$  with document  $d$

removed, and  $D$  is the total number of documents in the dataset. Equation (18) means that, the document  $d$  is more likely to choose the cluster with more documents when we only consider the cluster assignments of other documents.

Now, we will derive the second term in Equation (5) as follows:

$$\begin{aligned} p(d|z_d = z, \vec{d}_{z,-d}, \beta) \\ = \int p(d, \Phi_z | z_d = z, \vec{d}_{z,-d}, \beta) d\Phi_z \end{aligned} \quad (19)$$

$$= \int p(\Phi_z | z_d = z, \vec{d}_{z,-d}, \beta) p(d | \Phi_z, z_d = z, \vec{d}_{z,-d}, \beta) d\Phi_z \quad (20)$$

$$= \int p(\Phi_z | \vec{d}_{z,-d}, \beta) p(d | \Phi_z, z_d = z) d\Phi_z \quad (21)$$

Here, we use the Sum Rule of Probability in Equation (19) and exploits the Product Rule of Probability in Equation (20). In Equation (21), we apply the property of D-Separation.

Next, we will derive the first term of Equation (21) as follows:

$$\begin{aligned} p(\Phi_z | \vec{d}_{z,-d}, \beta) \\ = \frac{p(\Phi_z | \beta) p(\vec{d}_{z,-d} | \Phi_z)}{\int p(\Phi_z | \beta) p(\vec{d}_{z,-d} | \Phi_z) d\Phi_z} \end{aligned} \quad (22)$$

$$= \frac{\frac{1}{\Delta(\beta)} \prod_{t=1}^V \Phi_{z,t}^{\beta-1} \prod_{t=1}^V \Phi_{k,t}^{n_{z,-d}^t}}{\int \frac{1}{\Delta(\beta)} \prod_{t=1}^V \Phi_{z,t}^{\beta-1} \prod_{t=1}^V \Phi_{k,t}^{n_{z,-d}^t} d\Phi_z} \quad (23)$$

$$= \frac{1}{\Delta(\vec{n}_z + \beta)} \prod_{t=1}^V \Phi_{z,t}^{n_{z,-d}^t + \beta - 1} \quad (24)$$

$$= \text{Dir}(\Phi_z | \vec{n}_{z,-d} + \beta) \quad (25)$$

Here, we use the Bayes Rule in Equation (22),

Next, we derive the second term of Equation (21) as follows:

$$\begin{aligned} p(d | z_d = z, \vec{d}_{z,-d}, \beta) \\ = \int \text{Dir}(\Phi_z | \vec{n}_{z,-d} + \beta) \prod_{w \in d} \text{Mult}(w | \Phi_z) d\Phi_z \end{aligned} \quad (26)$$

$$= \int \frac{1}{\Delta(\vec{n}_{z,-d})} \prod_{t=1}^V \Phi_{z,t}^{n_{z,-d}^t + \beta - 1} \prod_{w \in d} \Phi_{z,w}^{n_d^w} d\Phi_z \quad (27)$$

$$= \frac{\Delta(\vec{n}_z + \beta)}{\Delta(\vec{n}_{z,-d} + \beta)} \quad (28)$$

$$= \frac{\prod_{t=1}^V \Gamma(n_{z,-d}^t + \beta)}{\Gamma(\sum_{t=1}^V (n_{z,-d}^t + \beta))} \frac{\Gamma(\sum_{t=1}^V (n_{z,-d}^t + \beta))}{\prod_{t=1}^V \Gamma(n_{z,-d}^t + \beta)} \quad (29)$$

$$= \frac{\prod_{w \in d} \prod_{j=1}^{N_d^w} (n_{z,-d}^w + \beta + j - 1)}{\prod_{i=1}^{N_d} (n_{z,-d} + V\beta + i - 1)} \quad (30)$$

Here, we use the property of  $\Gamma$  function:  $\frac{\Gamma(x+m)}{\Gamma(x)} = \prod_{i=1}^m (x+i-1)$  when we obtain Equation (30) from Equation (29). In Equation (30),  $N_d^w$  and  $N_d$  denote the number of word  $w$  in document  $d$  and the total number of words in document  $d$  respectively, and  $N_d = \sum_{w \in d} N_d^w$ . In addition,  $n_{z,-d}^w$  and  $n_{z,-d}$  denotes the number of word  $w$  in cluster  $z$  and the total number of words in cluster  $z$  with document  $d$  removed, respectively, and  $n_{z,-d} = \sum_{w=1}^V n_{z,-d}^w$ . In fact, Equation (30) measures similarity between document  $d$  and cluster  $z$  to a certain extent. And, document  $d$  is more likely to choose a cluster whose documents share more words with it.

Eventually, we obtain the form of the first term of Equation (3) as follows:

$$p(z_d = z | \vec{z}_{-d}, \vec{d}, \alpha, \beta) \propto \frac{m_{z,-d} + \alpha/K}{D - 1 + \alpha} \frac{\prod_{w \in d} \prod_{j=1}^{N_d^w} (n_{z,-d}^w + \beta + j - 1)}{\prod_{i=1}^{N_d} (n_{z,-d} + V\beta + i - 1)} \quad (31)$$

When  $K \rightarrow \infty$ , the following is the final form of the first term of Equation (3)

$$p(z_d = z | \vec{z}_{-d}, \vec{d}, \alpha, \beta) \propto \frac{m_{z,-d}}{D - 1 + \alpha} \frac{\prod_{w \in d} \prod_{j=1}^{N_d^w} (n_{z,-d}^w + \beta + j - 1)}{\prod_{i=1}^{N_d} (n_{z,-d} + V\beta + i - 1)} \quad (32)$$

Now, we will derive the second term of Equation (3), the generation of text embedding part, as follows:

$$p(z_d = z | \vec{z}_{-d}, \sigma^2, \sigma_0^2, \mu_0, \vec{d}, \alpha) = p(z_d = z | \vec{z}_{-d}, \sigma^2, \sigma_0^2, \mu_0, \alpha, d, \vec{d}_{-d}) \quad (33)$$

$$\propto p(z_d = z | \vec{z}_{-d}, \sigma^2, \sigma_0^2, \mu_0, \alpha, \vec{d}_{-d}) \times p(d | z_d = z, \vec{z}_{-d}, \sigma^2, \sigma_0^2, \mu_0, \alpha, \vec{d}_{-d}) \quad (34)$$

$$= p(z_d = z | \vec{z}_{-d}, \alpha) p(d | \vec{d}_{z,-d}, \sigma^2, \sigma_0^2, \mu_0) \quad (35)$$

Here, we use the Bayes Rule of Probability in Equation (34) and apply the property of D-Separation in Equation (35). We can notice that the first term of Equation (35) is as same as the first term of Equation (5), which we have derived already and the result is as follows:

$$p(z_d = z | \vec{z}_{-d}, \alpha) = \frac{m_{z,-d} + \alpha/K}{D - 1 + \alpha} \quad (36)$$

Now, we will derive the second term of Equation (35) as follows:

$$p(d | \vec{d}_{z,-d}, \sigma^2, \sigma_0^2, \mu_0) = \int_{\mu_z} p(d, \mu_z | \vec{d}_{z,-d}, \sigma^2, \sigma_0^2, \mu_0) d\mu_z \quad (37)$$

$$= \int_{\mu_z} p(\mu_z | \vec{d}_{z,-d}, \sigma^2, \sigma_0^2, \mu_0) p(d | \mu_z, \vec{d}_{z,-d}, \sigma^2, \sigma_0^2, \mu_0) d\mu_z \quad (38)$$

$$= \int_{\mu_z} p(\mu_z | \vec{d}_{z,-d}, \sigma^2, \sigma_0^2, \mu_0) p(d | \mu_z, \sigma^2) d\mu_z \quad (39)$$

Here, we use the Sum Rule of Probability in Equation (37) and apply Product Rule of Probability in Equation (38). Equation (39) exploits the property of D-Separation. The first term of Equation (39) denotes the posterior distribution of  $\mu_z$  and the second term of Equation (39) is a gaussian distribution as follows:

$$p(d | \mu_z, \sigma^2) = \mathcal{N}(d | \mu_z, \sigma^2) \quad (40)$$

Now, we will derive the first term of Equation (39) as follows:

$$p(\mu_z | \vec{d}_{z,-d}, \sigma^2, \sigma_0^2, \mu_0) = \frac{p(\mu_z | \sigma^2, \sigma_0^2, \mu_0) p(\vec{d}_{z,-d} | \mu_z, \sigma^2, \sigma_0^2, \mu_0)}{\int_{\mu_z} p(\mu_z | \sigma^2, \sigma_0^2, \mu_0) p(\vec{d}_{z,-d} | \mu_z, \sigma^2, \sigma_0^2, \mu_0) d\mu_z} \quad (41)$$

$$= \frac{p(\mu_z | \sigma_0^2, \mu_0) p(\vec{d}_{z,-d} | \mu_z, \sigma^2)}{\int_{\mu_z} p(\mu_z | \sigma_0^2, \mu_0) p(\vec{d}_{z,-d} | \mu_z, \sigma^2) d\mu_z} \quad (42)$$

$$= \mathcal{N}(\mu_z | \mu_n, \sigma_n^2) \quad (43)$$

Here, we use the Bayes Rule in Equation (41) and Equation (43) use the conjugacy for Gaussian distribution as follows:

$$p(\mu|\mu_0, \sigma_0^2) = \mathcal{N}(\mu|\mu_0, \sigma_0^2) \quad (44)$$

$$p(x|\mu, \sigma^2) = \mathcal{N}(x|\mu, \sigma^2) \quad (45)$$

$$p(\mu|x_1, x_2, \dots, x_n) = \mathcal{N}(\mu|\frac{\sigma_0^2 \sum_i x_i + \sigma^2 \mu_0}{n\sigma_0^2 + \sigma^2}, \frac{\sigma^2 \sigma_0^2}{n\sigma_0^2 + \sigma^2}) \quad (46)$$

Here, Equation (44) and Equation (45) is the prior distribution of  $\mu$  and likelihood of  $x$ , respectively. Then, we can obtain the posterior distribution of  $\mu$  as Equation (46). In our method,  $\vec{d}_{z,-d}$ , the documents in cluster  $z$  with document  $d$  removed, is  $x_1, x_2, \dots, x_n$  in Equation (46) and therefore,  $m_{z,-d}$ , the number of documents in cluster  $z$  without considering document  $d$ , is  $n$  in Equation (46). Given above comparison, we can obtain Equation (43) and  $\mu_n, \sigma_n^2$  are as follows:

$$\mu_n = \frac{1}{m_{z,-d}\sigma_0^2 + \sigma^2} (\sigma^2 \mu_0 + \sigma_0^2 \sum_{i \in z, i \neq d} \vec{e}_i) \quad (47)$$

$$\sigma_n^2 = (\frac{1}{\sigma_0^2} + \frac{m_{z,-d}}{\sigma^2})^{-1} \quad (48)$$

Here,  $e_i$  is text embedding representation of document  $i$ . Specially, in order to connect  $\sigma$  and  $\sigma_0$ , we make a trick:  $\sigma_0^2 = \sigma^2/\kappa$ . Then, the above two equations become the following form:

$$\mu_n = \frac{\kappa \mu_0 + \sum_{i \in z, i \neq d} \vec{e}_i}{\kappa + m_{z,-d}} \quad (49)$$

$$\sigma_n^2 = \frac{\sigma^2}{\kappa + m_{z,-d}} \quad (50)$$

Then, using the property of gaussian distribution, we can obtain the result of Equation (39) as follows:

$$\int \mathcal{N}(d|\mu_z, \sigma^2) \mathcal{N}(\mu_z|\mu_n, \sigma_n^2) d\mu_z = \mathcal{N}(d|\mu_n, \sigma_n^2 + \sigma^2) \quad (51)$$

Next, we get the probability choosing an existed cluster for text embedding part as follows:

$$\begin{aligned} p(z_d = z|\vec{z}_{-d}, \sigma^2, \sigma_0^2, \mu_0, \vec{d}, \alpha) \\ \propto \frac{m_{z,-d} + \alpha/K}{D - 1 + \alpha} \mathcal{N}(d|\mu_n, \sigma_n^2 + \sigma^2) \end{aligned} \quad (52)$$

when  $K \rightarrow \infty$ , we obtain the second term of (3) as follows:

$$\begin{aligned} p(z_d = z|\vec{z}_{-d}, \sigma^2, \sigma_0^2, \mu_0, \vec{d}, \alpha) \\ \propto \frac{m_{z,-d}}{D - 1 + \alpha} \mathcal{N}(d|\mu_n, \sigma_n^2 + \sigma^2) \end{aligned} \quad (53)$$

given above deduction, we can obtain the probability choosing an existed cluster as follows:

$$\begin{aligned} \lambda \frac{m_{z,-d}}{D - 1 + \alpha} \frac{\prod_{w \in d} \prod_{j=1}^{N_d^w} (n_{z,-d}^w + \beta + j - 1)}{\prod_{i=1}^{N_d} (n_{z,-d} + V\beta + i - 1)} \\ + (1 - \lambda) \frac{m_{z,-d}}{D - 1 + \alpha} \mathcal{N}(d|\mu_n, \sigma_n^2 + \sigma^2) \end{aligned} \quad (54)$$

## 1.2 the probability to choose a new cluster

In this part, we will give a detailed derivation of the probability to choose a new cluster using collapsed Gibbs Sampling.

$$\begin{aligned}
& p(z_d = K + 1 | \vec{z}_{-d}, \vec{d}, \alpha, \beta, \lambda, \sigma_0^2, \mu_0, \sigma^2) \\
&= p(z_d = K + 1, s_d = 1 | \vec{z}_{-d}, \vec{d}, \alpha, \beta, \lambda, \sigma_0^2, \mu_0, \sigma^2) \\
&\quad + p(z_d = K + 1, s_d = 0 | \vec{z}_{-d}, \vec{d}, \alpha, \beta, \lambda, \sigma_0^2, \mu_0, \sigma^2)
\end{aligned} \tag{55}$$

$$\begin{aligned}
&= p(s_d = 1 | \lambda) p(z_d = K + 1 | \vec{z}_{-d}, \vec{d}, \alpha, \beta) \\
&\quad + p(s_d = 0 | \lambda) p(z_d = K + 1 | \vec{z}_{-d}, \vec{d}, \sigma_0^2, \mu_0, \sigma^2)
\end{aligned} \tag{56}$$

$$\begin{aligned}
&= \lambda p(z_d = K + 1 | \vec{z}_{-d}, \vec{d}, \alpha, \beta) \\
&\quad + (1 - \lambda) p(z_d = K + 1 | \vec{z}_{-d}, \vec{d}, \sigma_0^2, \mu_0, \sigma^2)
\end{aligned} \tag{57}$$

We use the Sum Rule of Probability in Equation (55) and apply the property of *Bernoulli* distribution in Equation (57). The two parts of Equation (57) are the probability choosing a new cluster for the BoW part and text embedding part, respectively. We will first derive the first term of (57) as follows:

$$\begin{aligned}
& p(z_d = K + 1 | \vec{z}_{-d}, \vec{d}, \alpha, \beta) \\
&\propto p(z_d = K + 1 | \vec{z}_{-d}, \vec{d}_{-d}, \alpha, \beta) p(d | z_d = K + 1, \vec{z}_{-d}, \vec{d}_{-d}, \alpha, \beta)
\end{aligned} \tag{58}$$

$$= p(z_d = K + 1 | \vec{z}_{-d}, \alpha) p(d | z_d = K + 1, \beta) \tag{59}$$

Here, we use the Bayes Rule in Equation (58) and apply the property of D-Separation in Equation (59). We can derive the second term of Equation (59) as follows:

$$\begin{aligned}
& p(z_d = K + 1 | \vec{z}_{-d}, \alpha) \\
&= 1 - \sum_{k=1}^K p(z_d = k | \vec{z}_{-d}, \alpha)
\end{aligned} \tag{60}$$

$$= 1 - \frac{\sum_{k=1}^K m_{k,-d}}{D - 1 + \alpha} \tag{61}$$

$$= 1 - \frac{D - 1}{D - 1 + \alpha} \tag{62}$$

$$= \frac{\alpha}{D - 1 + \alpha} \tag{63}$$

Then, we will derive the first term of Equation (59) as follows:

$$p(d|z_d = K + 1, \beta)$$

$$= \int p(d, \Phi_{K+1}|z_d = K + 1, \beta) d\Phi_{K+1} \quad (64)$$

$$= \int p(\Phi_{K+1}|z_d = K + 1, \beta) p(d|\Phi_{K+1}, z_d = K + 1, \beta) d\Phi_{K+1} \quad (65)$$

$$= \int p(\Phi_{K+1}|\beta) p(d|\Phi_{K+1}, z_d = K + 1) d\Phi_{K+1} \quad (66)$$

$$= \int \text{Dir}(\Phi_{K+1}|\beta) \prod_{w \in d} \text{Mult}(w|\Phi_{K+1}) d\Phi_{K+1} \quad (67)$$

$$= \int \frac{1}{\Delta(\beta)} \prod_{t=1}^V \Phi_{K+1,t}^{\beta-1} \prod_{w \in d} \Phi_{K+1,w}^{N_d^w} d\Phi_{K+1} \quad (68)$$

$$= \frac{\Delta(\vec{n}_{K+1} + \beta)}{\Delta(\beta)} \quad (69)$$

$$= \frac{\sum_{t=1}^V \Gamma(n_{K+1}^t + \beta)}{\Gamma(\sum_{t=1}^V (n_{K+1}^t + \beta))} \frac{\Gamma(\sum_{t=1}^V \beta)}{\sum_{t=1}^V \Gamma(\beta)} \quad (70)$$

$$= \frac{\prod_{w \in d} \prod_{j=1}^{N_d^w} (\beta + j - 1)}{\prod_{i=1}^{N_d} (V\beta + i - 1)} \quad (71)$$

Here, we use the Sum Rule in Equation (64) and apply the Product Rule of probability in Equation (65). Equation (67) exploits the property of D-Separation. Then, we can obtain Equation (71) from Equation (70) using the following property of  $\Delta$  function:  $\frac{\Delta(x+m)}{\Delta(x)} = \prod_{i=1}^m (x + i - 1)$ . Then, we obtain the final form in Equation (57) as follows:

$$p(z_d = K + 1|\vec{z}_{-d}, \vec{d}, \alpha, \beta)$$

$$\propto \frac{\alpha}{D - 1 + \alpha} \frac{\prod_{w \in d} \prod_{j=1}^{N_d^w} (\beta + j - 1)}{\prod_{i=1}^{N_d} (V\beta + i - 1)} \quad (72)$$

now, we, will derive the second term of Equation (57) as follows:

$$p(z_d = K + 1|\vec{z}_{-d}, \sigma^2, \sigma_0^2, \mu_0, \vec{d}, \alpha)$$

$$= p(z_d = K + 1|\vec{z}_{-d}, \sigma^2, \sigma_0^2, \mu_0, \alpha, d, \vec{d}_{-d}) \quad (73)$$

$$\propto p(z_d = K + 1|\vec{z}_{-d}, \sigma^2, \sigma_0^2, \mu_0, \alpha_0, \vec{d}_{-d})$$

$$\times p(d|z_d = K + 1, \vec{z}_{-d}, \sigma^2, \sigma_0^2, \mu_0, \alpha, \vec{d}_{-d}) \quad (74)$$

$$= p(z_d = K + 1|\vec{z}_{-d}, \alpha) p(d|\vec{d}_{z=K+1, -d}, \sigma^2, \sigma_0^2, \mu_0) \quad (75)$$

Here, we use the Bayes Rule in Equation (74) and apply the Property of D-Separation in Equation (75). We can notice that the first term of Equation (75) is as same as the first term in Equation (59) and therefore it equals  $\frac{\alpha}{D-1+\alpha}$ . Besides, there is no document in the new cluster, so  $m_{z,-d} = 0$ ,  $\sum_{i \in z, i \neq d} \vec{e}_i = \vec{0}$ . Therefor, given the property of Gaussian distribution in Equation (46) and Equation (51), we can obtain the

second term in Equation (57) as follows:

$$\mu_n = \frac{1}{m_{z,\neg d}\sigma_0^2 + \sigma^2} (\sigma^2\mu_0 + \sigma_0^2 \sum_{i \in z, i \neq d} \vec{e}_i) \quad (76)$$

$$= \frac{1}{0\sigma_0^2 + \sigma^2} (\sigma^2\mu_0 + 0\sigma_0^2) \quad (77)$$

$$= \mu_0 \quad (78)$$

$$\sigma_n^2 = \left( \frac{1}{\sigma_0^2} + \frac{m_{z,\neg d}}{\sigma^2} \right)^{-1} \quad (79)$$

$$= \sigma_0^2 \quad (80)$$

the following is the final form of the second term of Equation (57)

$$\begin{aligned} & p(z_d = K + 1 | \vec{z}_{-d}, \sigma^2, \sigma_0^2, \mu_0, \vec{d}, \alpha) \\ & \propto \frac{\alpha}{D - 1 + \alpha} \mathcal{N}(d | \mu_0, \sigma_0^2 + \sigma^2) \end{aligned} \quad (81)$$

**we finally obtain the probability choosing a new cluster for a document as follows**

$$\begin{aligned} & \lambda \frac{\alpha}{D - 1 + \alpha} \frac{\prod_{w \in d} \prod_{j=1}^{N_d^w} (\beta + j - 1)}{\prod_{i=1}^{N_d} (V\beta + i - 1)} \\ & + (1 - \lambda) \frac{\alpha}{D - 1 + \alpha} \mathcal{N}(d | \mu_0, \sigma_0^2 + \sigma^2) \end{aligned} \quad (82)$$