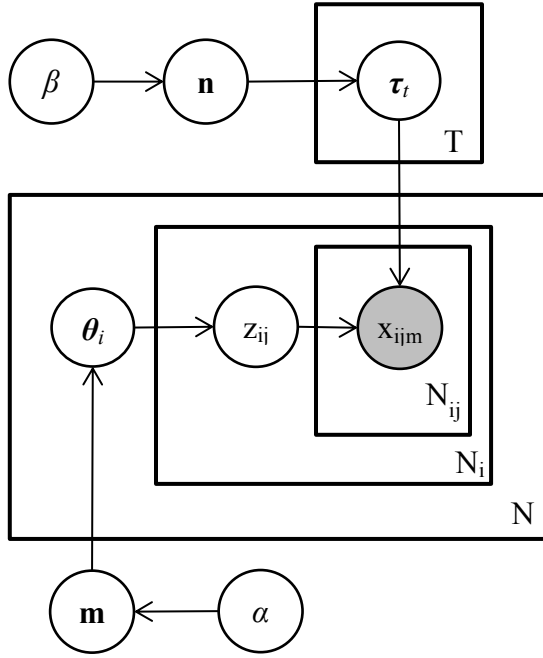


## Derivation of Posterior Probabilities for TpT-LDA



### Definitions:

1. Topic – A multinomial probability distribution over all words in the vocabulary. The total number of topics in an LDA model needs to be either pre-determined or an appropriate number needs to be selected by cross-validation.
2. Document – Documents are represented as bag-of-words. In context of Twitter, each *tweet* treated as a document.
3. User – This is a unique person sending out messages. Typically, each person will have a particular distribution over topics based on personal preferences. In Twitter, a person is identified by the unique Twitter *handle*.

### Notation:

$N$  = Number of users.

$N_i$  = Number of documents for  $i$ -th user.

$N_{ij}$  = Number of words within the  $j$ -th document of the  $i$ -th user.

$T$  = Number of topics.

$V$  = Number of unique words in vocabulary.

$W$  = Size of the entire corpus, i.e. the total number of words across all document and all users.

$\mathbf{X}$  = The entire corpus – the complete set of words across all document and all users.

$\Theta = \{\theta_i\}$  = Distribution over topics for  $i$ -th user;  $i = 1, \dots, N$

$\mathbf{T} = \{\tau_t\}$  = Distribution over words for  $t$ -th topic;  $t = 1, \dots, T$

$\mathbf{Z} = \{z_{ij}\}$  = Topic for  $j$ -th document of  $i$ -th user.

$\mathbf{x}_{ij} = j$ -th document of  $i$ -th user;  $j = 1, \dots, N_i$

$x_{ijm} = m$ -th word in the  $j$ -th document of  $i$ -th user.

$\beta, \alpha$  – Dirichlet concentration parameters

$\mathbf{n}, \mathbf{m}$  – uniform Dirichlet hyper-priors

Unlike the regular LDA model, each document in our setup belongs to only one topic. This is reasonable in our context because we expect the documents to be short and therefore be generated from only one topic.

The generative process is:

1. For each user  $u_i$ :
  - a. Draw the topic distribution  $\theta_i \sim \text{Dir}(\alpha \mathbf{m})$
  - b. For each document:
    - i. Draw topic  $z_{ij} \sim \text{Multi}(\theta_i)$
    - ii. Draw the words within the document  $\mathbf{x}_{ij} \sim \text{Multi}(\tau_{z_{ij}})$

The complete data likelihood is:

$$p(\mathbf{X}, \mathbf{Z}, \boldsymbol{\theta}, \mathbf{T} | \beta \mathbf{n}, \alpha \mathbf{m}) = p(\mathbf{T} | \beta \mathbf{n}) p(\boldsymbol{\theta} | \alpha \mathbf{m}) p(\mathbf{Z} | \boldsymbol{\theta}) p(\mathbf{X} | \mathbf{Z}, \mathbf{T})$$

Marginalizing out  $\mathbf{Z}$  to get the observed data likelihood leads to the following for each document:

$$p(\mathbf{x}_{ij} | \theta_i, \mathbf{T}) = \sum_{t=1}^T \left[ \prod_m^{N_{ij}} p(x_{ijm} | \tau_t) \right] p(z_{ij} = t | \theta_i)$$

Now, the likelihood for the entire corpus is:

$$p(\mathbf{X} | \boldsymbol{\theta}, \mathbf{T}) = \prod_{i=1}^N \prod_{j=1}^{N_i} p(\mathbf{x}_{ij} | \theta_i, \mathbf{T}) \quad (1)$$

where:

$$\begin{aligned} P(\mathbf{T} | \beta \mathbf{n}) &= \text{Dir}(\mathbf{T} | \beta \mathbf{n}) \\ P(\theta_i | \alpha \mathbf{m}) &= \text{Dir}(\theta_i | \alpha \mathbf{m}) \\ P(z_{ij} | \theta_i) &= \text{Multi}(z_{ij} | \theta_i) \\ P(x_{ijm} | z_{ij}, \tau_{z_{ij}}) &= \text{Multi}(x_{ijm} | z_{ij}, \tau_{z_{ij}}) \end{aligned}$$

Exact inference of the parameters is hard for the LDA model. Therefore, we use Gibbs sampling. Most of the derivations below are similar to that of the original LDA model in (Heinrich, 2008).

The Collapsed Gibbs sampler is often more efficient than the regular Gibbs sampler where we try to avoid inference on some parameters. In our case, these parameters are  $\boldsymbol{\theta}$  and  $\mathbf{T}$ . These parameters can be interpreted as statistics of the associations between the observed  $x_{ijm}$  and the corresponding  $z_{ij}$ . The target inference is the distribution  $p(\mathbf{Z} | \mathbf{X})$ :

$$p(\mathbf{Z} | \mathbf{X}) = \frac{p(\mathbf{Z}, \mathbf{X})}{p(\mathbf{X})} = \frac{\prod_{i=1}^W p(z_i, x_i)}{\prod_{i=1}^W \sum_{t=1}^T p(z_i = t, x_i)}$$

The Gibbs sampler runs a Markov chain that uses the full conditional  $p(z_i | \mathbf{Z}_{-i}, \mathbf{X})$  in order to simulate  $p(\mathbf{Z} | \mathbf{X})$ .  $i$  refers to the tuple  $(n, m)$  (i.e.,  $n$ -th user's  $m$ -th document). ' $\neg i$ ' indicates that

document  $i$  is excluded. The joint distribution  $p(\mathbf{Z}, \mathbf{X})$  can be factored as:

$$p(\mathbf{Z}, \mathbf{X} | \alpha, \beta) = p(\mathbf{X} | \mathbf{Z}, \beta) p(\mathbf{Z} | \alpha) \quad (2)$$

$p(\mathbf{X} | \mathbf{Z}, \beta)$  is a multinomial over the word counts given the associated topics.

$$p(\mathbf{X} | \mathbf{Z}, \beta) = \prod_{i=1}^W p(x_i | z_i) = \prod_{i=1}^W \tau_{z_i, x_i} \quad (3)$$

The above equation (3) can also be written as:

$$p(\mathbf{X} | \mathbf{Z}, \mathbf{T}) = \prod_{t=1}^T \prod_{m=1}^V \tau_{t,m}^{n_t^{(m)}} \quad (4)$$

In equation (4),  $n_t^{(m)}$  refers to the number of times that the word  $m$  occurs in the corpus across all documents within topic  $t$ . The distribution  $p(\mathbf{X} | \mathbf{Z}, \beta)$  can be computed by integrating over  $\mathbf{T}$ .

$$\begin{aligned} p(\mathbf{X} | \mathbf{Z}, \beta) &= \int p(\mathbf{X} | \mathbf{Z}, \mathbf{T}) p(\mathbf{T} | \beta) d\mathbf{T} \\ &= \int \prod_{t=1}^T \frac{1}{\Delta(\boldsymbol{\beta})} \prod_{m=1}^V \tau_{t,m}^{n_t^{(m)} + \beta_m - 1} d\boldsymbol{\tau}_t \\ &= \prod_{t=1}^T \frac{\Delta(\mathbf{n}_t + \boldsymbol{\beta})}{\Delta(\boldsymbol{\beta})}, \mathbf{n}_t = \{n_t^{(m)}\}_{m=1}^V, \Delta(\boldsymbol{\beta}) = \frac{\prod_{m=1}^V \Gamma(\beta_m)}{\Gamma(\sum_{m=1}^V \beta_m)} \end{aligned}$$

Similar to the above, we will now derive  $p(\mathbf{Z} | \alpha)$ . We first start with  $p(\mathbf{Z} | \boldsymbol{\theta})$ :

$$p(\mathbf{Z} | \boldsymbol{\theta}) = \prod_{i=1}^W p(z_i | \theta_i) = \prod_{n=1}^N \prod_{t=1}^T p(z_i = t | u_i, \theta_i) = \prod_{n=1}^N \prod_{t=1}^T \theta_{n,t}^{n_n^{(t)}} \quad (5)$$

In equation (5),  $u_i$  is the user to which  $z_i$  belongs;  $\theta_i$  refers to  $u_i$ 's topic distribution.  $n_n^{(t)}$  refers to the number of times that topic  $t$  has been observed with a document of user  $n$ . We will now integrate out  $\boldsymbol{\theta}$ .

$$\begin{aligned} p(\mathbf{Z} | \alpha) &= \int p(\mathbf{Z} | \boldsymbol{\theta}) p(\boldsymbol{\theta} | \alpha) d\boldsymbol{\theta} \\ &= \int \prod_{n=1}^N \frac{1}{\Delta(\boldsymbol{\alpha})} \prod_{t=1}^T \theta_{n,t}^{n_n^{(t)} + \alpha_t - 1} d\boldsymbol{\theta}_n \\ &= \prod_{n=1}^N \frac{\Delta(\mathbf{n}_n + \boldsymbol{\alpha})}{\Delta(\boldsymbol{\alpha})}, \mathbf{n}_n = \{n_n^{(t)}\}_{t=1}^T, \Delta(\boldsymbol{\alpha}) = \frac{\prod_{t=1}^T \Gamma(\alpha_t)}{\Gamma(\sum_{t=1}^T \alpha_t)} \end{aligned}$$

The joint distribution (2) now becomes:

$$p(\mathbf{Z}, \mathbf{X} | \alpha, \beta) = \prod_{t=1}^T \frac{\Delta(\mathbf{n}_t + \boldsymbol{\beta})}{\Delta(\boldsymbol{\beta})} \cdot \prod_{n=1}^N \frac{\Delta(\mathbf{n}_n + \boldsymbol{\alpha})}{\Delta(\boldsymbol{\alpha})}$$

For Gibbs sampling, we need to sample  $p(z_i = t | \mathbf{Z}_{\neg i}, \mathbf{X})$ .

$$\begin{aligned}
p(z_i = t | \mathbf{Z}_{\neg i}, \mathbf{X}) &= \frac{p(\mathbf{X}, \mathbf{Z})}{p(\mathbf{X}, \mathbf{Z}_{\neg i})} = \frac{p(\mathbf{X} | \mathbf{Z})}{p(\mathbf{X}_{\neg i} | \mathbf{Z}_{\neg i}) p(w_i)} \cdot \frac{p(\mathbf{Z})}{p(\mathbf{Z}_{\neg i})} \\
&\propto \frac{\Delta(\mathbf{n}_t + \boldsymbol{\beta})}{\Delta(\mathbf{n}_{t, \neg i} + \boldsymbol{\beta})} \cdot \frac{\Delta(\mathbf{n}_n + \boldsymbol{\alpha})}{\Delta(\mathbf{n}_{n, \neg i} + \boldsymbol{\alpha})} \\
&= \frac{\prod_{m \in \mathbf{X}_i} \Gamma(n_t^{(m)} + \beta_m) \Gamma(\sum_{m=1}^V n_{t, \neg i}^{(m)} + \beta_m)}{\prod_{m \in \mathbf{X}_i} \Gamma(n_{t, \neg i}^{(m)} + \beta_m) \Gamma(\sum_{m=1}^V n_t^{(m)} + \beta_m)} \\
&\quad \cdot \frac{\Gamma(n_n^{(t)} + \alpha_t) \Gamma(\sum_{t'=1}^T n_{n, \neg i}^{(t')} + \alpha_{t'})}{\Gamma(n_{n, \neg i}^{(t)} + \alpha_t) \Gamma(\sum_{t'=1}^T n_n^{(t')} + \alpha_{t'})} \\
&= \frac{\prod_{m \in \mathbf{X}_i} \prod_{j=0}^{f_m} (n_{t, \neg i}^{(m)} + \beta_m + j)}{\prod_{j=0}^{length(\mathbf{X}_i)} (\sum_{m=1}^V n_{t, \neg i}^{(m)} + \beta_m + j)} \cdot \frac{n_{n, \neg i}^{(t)} + \alpha_t}{[\sum_{t'=1}^T n_n^{(t')} + \alpha_{t'}] - 1} \\
&\propto \frac{\prod_{m \in \mathbf{X}_i} \prod_{j=0}^{f_m} (n_{t, \neg i}^{(m)} + \beta_m + j)}{\prod_{j=0}^{length(\mathbf{X}_i)} (\sum_{m=1}^V n_{t, \neg i}^{(m)} + \beta_m + j)} (n_{n, \neg i}^{(t)} + \alpha_t) \quad (6)
\end{aligned}$$

In the above, counts ' $\mathbf{n}_{\cdot, \neg i}^{(\cdot)}$ ' indicate that the document  $i$  was excluded.  $f_m$  is the count of occurrences of word  $m$  in document  $i$ .  $length(\mathbf{X}_i)$  is the total count of words in document  $i$ .  $\{m \in \mathbf{X}_i\}$  is the set of unique words in document  $i$ .

The multinomial parameters  $(\boldsymbol{\theta}, \mathbf{T})$  can be computed as:

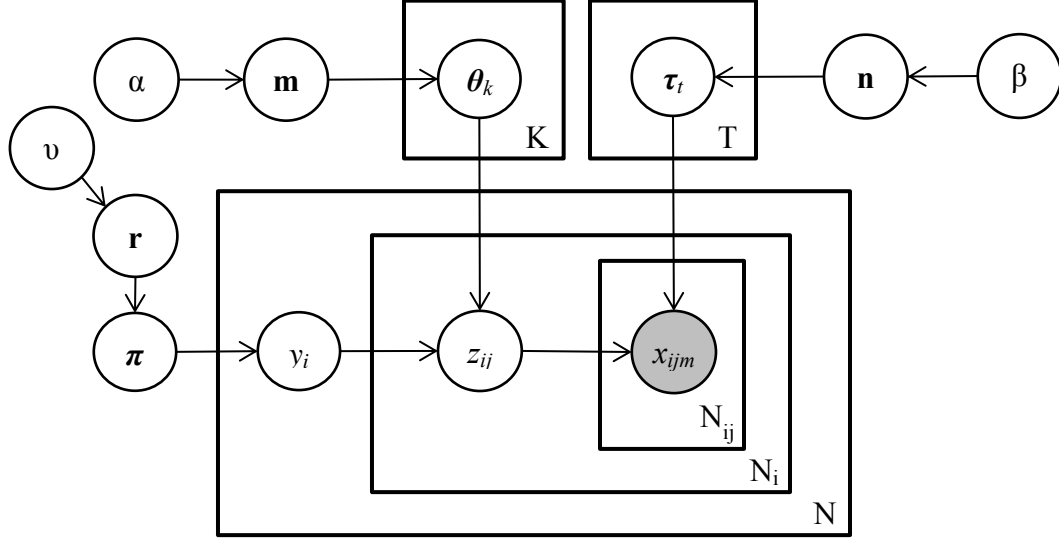
$$\begin{aligned}
p(\boldsymbol{\theta}_n | \mathbf{z}_n, \boldsymbol{\alpha}) &= \frac{1}{Z_{\boldsymbol{\theta}_n}} \prod_{m=1}^{N_n} p(z_{n,m} | \boldsymbol{\theta}_n) \cdot p(\boldsymbol{\theta}_n | \boldsymbol{\alpha}) = Dir(\mathbf{n}_n + \boldsymbol{\alpha}) \\
p(\boldsymbol{\tau}_t | \mathbf{Z}, \mathbf{X}, \boldsymbol{\beta}) &= \frac{1}{Z_{\boldsymbol{\tau}_t}} \prod_{\{i: z_i = t\}} p(w_i | \boldsymbol{\tau}_t) \cdot p(\boldsymbol{\tau}_t | \boldsymbol{\beta}) = Dir(\mathbf{n}_t + \boldsymbol{\beta})
\end{aligned}$$

Since both the above are Dirichlet distributions, we can compute the expectations easily:

$$\tau_{t,m} = \frac{n_t^{(m)} + \beta_m}{\sum_{m=1}^V n_t^{(m)} + \beta_m} \quad (7)$$

$$\theta_{n,t} = \frac{n_n^{(t)} + \alpha_t}{\sum_{t'=1}^T n_n^{(t')} + \alpha_{t'}} \quad (8)$$

### Derivation of Posterior Probabilities for GTpT-LDA



#### Definitions:

In this version of the model, we continue with the same three definitions as in the previous model and add one more:

4. Group – We define a group as an entity whose users share the same topic distribution. In our model, we allow one user to belong to only one group.

The major change here is that we now define the topic distributions for groups of users so that we can impose a clustering over those users.

#### Notation:

$N$  = Number of users.

$N_i$  = Number of documents for  $i$ -th user.

$N_{ij}$  = Number of words within the  $j$ -th document of the  $i$ -th user.

$T$  = Number of topics.

$K$  = Number of user groups.

$V$  = Number of unique words in vocabulary.

$W$  = Size of the entire corpus, i.e. the total number of words across all documents and all users.

$X$  = The entire corpus – the complete set of words across all documents and all users.

$\Theta = \{\theta_k\}$  = Topic proportions for  $k$ -th user-group;  $k = 1, \dots, K$ . This is a  $T$ -by- $K$  matrix.

$T = \{\tau_t\}$  = Distribution over words for  $t$ -th topic;  $t = 1, \dots, T$

$Y = \{y_i\}$  = Group to which the  $i$ -th user belongs;  $i = 1, \dots, N$

$Z = \{z_{ij}\}$  = Topic for  $j$ -th document of  $i$ -th user.

$x_{ij}$  =  $j$ -th document of  $i$ -th user;  $j = 1, \dots, N_i$

$x_{ijm}$  –  $m$ -th word in the  $j$ -th document of  $i$ -th user.

$\alpha, \beta, \nu$  – Dirichlet concentration parameters

$\mathbf{n}, \mathbf{m}, \mathbf{r}$  – uniform Dirichlet hyper-priors

The generative process is:

1. For each user  $u_i$ :
  - a. Draw the group membership  $y_i \sim \text{Multi}(\boldsymbol{\pi})$
  - b. Draw a topic distribution for the group  $\boldsymbol{\theta}_{y_i} \sim \text{Dir}(\alpha \mathbf{m})$
  - c. For each document of user  $u_i$ :
    - i. Draw a topic  $z_{ij} \sim \text{Multi}(\boldsymbol{\theta}_{y_i})$
    - ii. For each words within the document
$$\text{Draw } x_{ijm} \sim \text{Multi}(\boldsymbol{\tau}_{z_{ij}})$$

The complete data likelihood is:

$$\begin{aligned} p(\mathbf{X}, \mathbf{Y}, \mathbf{Z}, \boldsymbol{\Theta}, \mathbf{T}, \boldsymbol{\pi} | \beta \mathbf{n}, \alpha \mathbf{m}) \\ = p(\mathbf{T} | \beta \mathbf{n}) p(\boldsymbol{\pi} | \nu \mathbf{r}) p(\mathbf{Y} | \boldsymbol{\pi}) p(\boldsymbol{\Theta} | \alpha \mathbf{m}) p(\mathbf{Z} | \mathbf{Y}, \boldsymbol{\Theta}) p(\mathbf{X} | \mathbf{Z}, \mathbf{T}) \end{aligned}$$

Marginalizing out  $\mathbf{Z}, \mathbf{Y}$  to get the observed data likelihood leads to the following for the  $i$ -th user where  $\mathbf{x}_i$  is the complete set of words for the  $i$ -th user across all documents of that user:

$$p(\mathbf{x}_i | \boldsymbol{\Theta}, \mathbf{T}, \boldsymbol{\pi}) = \sum_{k=1}^K p(y_i = \pi_k) \left\{ \prod_{j=1}^{N_i} \left[ \sum_{t=1}^T p(z_{ij} = t | \boldsymbol{\theta}_k) \left( \prod_{m=1}^{N_{ij}} p(x_{ijm} | \boldsymbol{\tau}_t) \right) \right] \right\}$$

Now, the likelihood for the entire corpus is (where  $N$  is the total number of users):

$$p(\mathbf{X} | \boldsymbol{\Theta}, \mathbf{T}, \boldsymbol{\pi}) = \prod_{i=1}^N p(\mathbf{x}_i | \boldsymbol{\Theta}, \mathbf{T}, \boldsymbol{\pi}) \quad (9)$$

where:

$$\begin{aligned} P(\boldsymbol{\pi} | \nu \mathbf{r}) &= \text{Dir}(\boldsymbol{\pi} | \nu \mathbf{r}) \\ P(\mathbf{T} | \beta \mathbf{n}) &= \text{Dir}(\mathbf{T} | \beta \mathbf{n}) \\ P(\boldsymbol{\theta}_i | \alpha \mathbf{m}) &= \text{Dir}(\boldsymbol{\theta}_i | \alpha \mathbf{m}) \\ P(y_i | \boldsymbol{\pi}) &= \text{Multi}(y_i | \boldsymbol{\pi}) \\ P(z_{ij} | y_i, \boldsymbol{\theta}_{y_i}) &= \text{Multi}(z_{ij} | y_i, \boldsymbol{\theta}_{y_i}) \\ P(x_{ijm} | z_{ij}, \boldsymbol{\tau}_{z_{ij}}) &= \text{Multi}(x_{ijm} | z_{ij}, \boldsymbol{\tau}_{z_{ij}}) \end{aligned}$$

Here again, we use collapsed Gibbs sampling instead of exact inference. In this case, these parameters are  $\boldsymbol{\Theta}, \mathbf{T}$ , and  $\boldsymbol{\pi}$ . These parameters can be interpreted as statistics of the associations between the observed  $x_{ijm}$  and the corresponding  $y_i$  and  $z_{ij}$ . The target inference is the distribution  $p(\mathbf{Y}, \mathbf{Z} | \mathbf{X})$ :

$$p(\mathbf{Y}, \mathbf{Z} | \mathbf{X}) = \frac{p(\mathbf{Y}, \mathbf{Z}, \mathbf{X})}{p(\mathbf{X})}$$

The Gibbs sampler runs a Markov chain that uses the full conditional  $p(y_i, \mathbf{z}_i | \mathbf{Y}_{-i}, \mathbf{Z}_{-i}, \mathbf{X})$  in

order to simulate  $p(\mathbf{Y}, \mathbf{Z}|\mathbf{X})$ . Here, ‘ $i$ ’ refers to all documents of the  $i$ -th user. ‘ $\neg i$ ’ indicates that all documents of the  $i$ -th are excluded. The joint distribution  $p(\mathbf{Y}, \mathbf{Z}, \mathbf{X})$  can be factored as:

$$p(\mathbf{Y}, \mathbf{Z}, \mathbf{X}|\alpha, \beta, v) = p(\mathbf{X}|\mathbf{Z}, \beta)p(\mathbf{Z}|\mathbf{Y}, \alpha)p(\mathbf{Y}|v) \quad (10)$$

$p(\mathbf{X}|\mathbf{Z}, \beta)$  is a multinomial over the word counts given the associated topics.

$$p(\mathbf{X}|\mathbf{Z}, \beta) = \prod_{i=1}^W p(x_i|z_i) = \prod_{i=1}^W \tau_{z_i, x_i} \quad (11)$$

The above equation (11) can also be written as (once  $\mathbf{Z}$  is known,  $\mathbf{X}$  is independent of  $\mathbf{Y}$ ):

$$p(\mathbf{X}|\mathbf{Z}, \boldsymbol{\theta}, \mathbf{T}, \boldsymbol{\pi}) = \prod_{t=1}^T \prod_{m=1}^V \tau_{t,m}^{n_t^{(m)}} \quad (12)$$

In equation (12),  $n_t^{(m)}$  refers to the number of times that the word  $m$  occurs in the corpus across all documents belonging to topic  $t$ . The distribution  $p(\mathbf{X}|\mathbf{Z}, \beta)$  can be computed by integrating over  $\mathbf{T}$ .

$$\begin{aligned} p(\mathbf{X}|\mathbf{Z}, \beta) &= \int p(\mathbf{X}|\mathbf{Z}, \mathbf{T})p(\mathbf{T}|\beta)d\mathbf{T} \\ &= \int \prod_{t=1}^T \frac{1}{\Delta(\boldsymbol{\beta})} \prod_{m=1}^V \tau_{t,m}^{n_t^{(m)} + \beta_m - 1} d\boldsymbol{\tau}_t \\ &= \prod_{t=1}^T \frac{\Delta(\mathbf{n}_t + \boldsymbol{\beta})}{\Delta(\boldsymbol{\beta})}, \mathbf{n}_t = \{n_t^{(m)}\}_{m=1}^V, \quad \Delta(\boldsymbol{\beta}) = \frac{\prod_{m=1}^V \Gamma(\beta_m)}{\Gamma(\sum_{m=1}^V \beta_m)} \end{aligned}$$

Similar to the above, we can derive  $p(\mathbf{Z}|\mathbf{Y}, \alpha)$  as:

$$p(\mathbf{Z}|\mathbf{Y}, \boldsymbol{\theta}, \alpha) = \prod_{i=1}^W p(z_i|y_i) = \prod_{n=1}^N \prod_{k=1}^K \prod_{t=1}^T p(z_i = t|y_i = k) = \prod_{n=1}^N \prod_{k=1}^K \prod_{t=1}^T \theta_{k,t}^{m_{k,n}^{(t)}} \quad (13)$$

In equation (13),  $i \equiv (n, k, t)$ ; and  $y_i$  refers to the group for the  $i$ -th instance.  $m_{k,n}^{(t)}$  refers to the number of times the topic  $t$  has been observed within the documents of user  $n$  who belongs to group  $k$ . Note that each user can belong to only one group. We will now integrate out  $\boldsymbol{\theta}$ .

$$\begin{aligned} p(\mathbf{Z}|\mathbf{Y}, \alpha) &= \int p(\mathbf{Z}|\mathbf{Y}, \boldsymbol{\theta})p(\boldsymbol{\theta}|\alpha)d\boldsymbol{\theta} \\ &= \int \prod_{n=1}^N \prod_{k=1}^K \frac{1}{\Delta(\boldsymbol{\alpha})} \prod_{t=1}^T \theta_{k,t}^{m_{k,n}^{(t)} + \alpha_t - 1} d\boldsymbol{\theta}_k \\ &= \int \prod_{k=1}^K \frac{1}{\Delta(\boldsymbol{\alpha})} \prod_{t=1}^T \theta_{k,t}^{m_k^{(t)} + \alpha_t - 1} d\boldsymbol{\theta}_k \\ &= \prod_{k=1}^K \frac{\Delta(\mathbf{m}_k + \boldsymbol{\alpha})}{\Delta(\boldsymbol{\alpha})}, \mathbf{m}_k = \{m_k^{(t)}\}_{t=1}^T, \quad \Delta(\boldsymbol{\alpha}) = \frac{\prod_{t=1}^T \Gamma(\alpha_t)}{\Gamma(\sum_{t=1}^T \alpha_t)} \end{aligned}$$

Here,  $m_k^{(t)}$  is the number of times that topic  $t$  has been observed in group  $k$  across all users belonging to the group  $k$ .

Now for evaluating  $p(\mathbf{Y}|\mathbf{v})$  we first start with:

$$p(\mathbf{Y}|\boldsymbol{\pi}, \mathbf{v}) = \prod_{k=1}^K \pi_k^{n^{(k)}}$$

Next, we get  $p(\mathbf{Y}|\mathbf{v})$  by integrating out  $\boldsymbol{\pi}$ :

$$p(\mathbf{Y}|\mathbf{v}) = \int p(\mathbf{Y}|\boldsymbol{\pi})p(\boldsymbol{\pi}|\mathbf{v})d\boldsymbol{\pi} = \frac{\Delta(\mathbf{n} + \mathbf{v})}{\Delta(\mathbf{v})}, \mathbf{n} = \{n^{(k)}\}_{k=1}^K, \Delta(\mathbf{v}) = \frac{\prod_{k=1}^K \Gamma(v_k)}{\Gamma(\sum_{k=1}^K v_k)}$$

Here,  $n^{(k)}$  is the number of users in  $k$ -th group.

The joint distribution (10) now becomes:

$$p(\mathbf{Y}, \mathbf{Z}, \mathbf{X}|\alpha, \beta) = \left( \prod_{t=1}^T \frac{\Delta(\mathbf{n}_t + \boldsymbol{\beta})}{\Delta(\boldsymbol{\beta})} \right) \cdot \left( \prod_{k=1}^K \frac{\Delta(\mathbf{m}_k + \boldsymbol{\alpha})}{\Delta(\boldsymbol{\alpha})} \right) \cdot \left( \frac{\Delta(\mathbf{n}_k + \mathbf{v})}{\Delta(\mathbf{v})} \right) \quad (14)$$

For Gibbs sampling, we need to sample  $p(y_i|\mathbf{Y}_{\neg i}, \mathbf{z}_i, \mathbf{Z}_{\neg i}, \mathbf{X})$ , followed by  $p(z_i|y_i, \mathbf{Y}_{\neg i}, \mathbf{Z}_{\neg i}, \mathbf{X})$ .

While computing  $p(y_i|\mathbf{Y}_{\neg i}, \mathbf{z}_i, \mathbf{Z}_{\neg i}, \mathbf{X})$ , ‘ $i$ ’ refers to all documents of the  $i$ -th user. ‘ $\neg i$ ’ means that all documents of user ‘ $i$ ’ are excluded.

$$\begin{aligned} p(y_i = k|\mathbf{Y}_{\neg i}, \mathbf{z}_i, \mathbf{Z}_{\neg i}, \mathbf{X}) &= \frac{p(\mathbf{Y}, \mathbf{Z}, \mathbf{X})}{p(\mathbf{Y}_{\neg i}, \mathbf{z}_i, \mathbf{Z}_{\neg i}, \mathbf{X})} \\ &= \frac{p(\mathbf{X}|\mathbf{Y}, \mathbf{Z})}{p(\mathbf{X}_{\neg i}|\mathbf{Y}_{\neg i}, \mathbf{Z}_{\neg i})p(\mathbf{z}_i, \mathbf{x}_i)} \cdot \frac{p(\mathbf{Z}|\mathbf{Y})}{p(\mathbf{Z}_{\neg i}|\mathbf{Y}_{\neg i})} \cdot \frac{p(\mathbf{Y})}{p(\mathbf{Y}_{\neg i})} \end{aligned}$$

Now, since  $p(\mathbf{z}_i, \mathbf{x}_i)$  does not depend on  $y_i$ ,

$$p(y_i = k|\mathbf{Y}_{\neg i}, \mathbf{z}_i, \mathbf{Z}_{\neg i}, \mathbf{X}) \propto \frac{p(\mathbf{X}|\mathbf{Y}, \mathbf{Z})}{p(\mathbf{X}_{\neg i}|\mathbf{Y}_{\neg i}, \mathbf{Z}_{\neg i})} \cdot \frac{p(\mathbf{Z}|\mathbf{Y})}{p(\mathbf{Z}_{\neg i}|\mathbf{Y}_{\neg i})} \cdot \frac{p(\mathbf{Y})}{p(\mathbf{Y}_{\neg i})} \quad (15)$$

$$= \left( \prod_{t=1}^T \frac{\Delta(\mathbf{n}_t + \boldsymbol{\beta})}{\Delta(\mathbf{n}_{t, \neg i} + \boldsymbol{\beta})} \right) \cdot \left( \frac{\Delta(\mathbf{m}_k + \boldsymbol{\alpha})}{\Delta(\mathbf{m}_{k, \neg i} + \boldsymbol{\alpha})} \right) \cdot \left( \frac{\Delta(\mathbf{n} + \mathbf{v})}{\Delta(\mathbf{n}_{\neg i} + \mathbf{v})} \right) \quad (16)$$

Equation (16) follows from substituting (14) into (15). Also, note that  $\Delta(\mathbf{m}_k + \boldsymbol{\alpha})$  is obtained by removing all documents of user  $i$  from group  $y_i$  and adding them to group  $k$ . All product terms other than those for group  $k$  cancel out from the middle term of equation (14).

$$\begin{aligned} &= \left( \prod_{t=1}^T \frac{\prod_{m=1}^V \Gamma(n_t^{(m)} + \beta_m) \Gamma(\sum_{m=1}^V n_{t, \neg i}^{(m)} + \beta_m)}{\prod_{m=1}^V \Gamma(n_{t, \neg i}^{(m)} + \beta_m) \Gamma(\sum_{m=1}^V n_t^{(m)} + \beta_m)} \right) \\ &\quad \times \left( \frac{\prod_{t=1}^T \Gamma(m_k^{(t)} + \alpha_t) \Gamma(\sum_{t=1}^T m_{k, \neg i}^{(t)} + \alpha_t)}{\prod_{t=1}^T \Gamma(m_{k, \neg i}^{(t)} + \alpha_t) \Gamma(\sum_{t=1}^T m_k^{(t)} + \alpha_t)} \right) \\ &\quad \times \left( \frac{\prod_{k=1}^K \Gamma(n^{(k)} + v_k) \Gamma(\sum_{l=1}^K n_{\neg i}^{(l)} + v_l)}{\prod_{k=1}^K \Gamma(n_{\neg i}^{(k)} + v_k) \Gamma(\sum_{l=1}^K n^{(l)} + v_l)} \right) \end{aligned}$$



$$\begin{aligned}
&= \left( \prod_{t=1}^T \frac{\prod_{m=1}^V \prod_{j=0}^{n_m-1} (n_{t,\neg i}^{(m)} + \beta_m + j)}{\prod_{j=0}^{\#actions_i^{t-1}} (\sum_{m=1}^V n_{t,\neg i}^{(m)} + \beta_m + j)} \right) \\
&\quad \times \left( \frac{\prod_{t=1}^T \prod_{j=0}^{m_{y_i,i}^{(t)}-1} (m_{k,\neg i}^{(t)} + \alpha_t + j)}{\prod_{j=0}^{(\sum_{t=1}^T m_{y_i,i}^{(t)})-1} (\sum_{t=1}^T m_{k,\neg i}^{(t)} + \alpha_t + j)} \right) \\
&\quad \times \left( \frac{n_{\neg i}^{(k)} + v_k}{[\sum_{l=1}^K n_{\neg i}^{(l)} + v_l] - 1} \right)
\end{aligned}$$

In the above expression,  $\sum_{t=1}^T m_{y_i,i}^{(t)}$  is the total number of documents for user  $i$  (who belongs to group  $y_i$ ) and  $m_{y_i,i}^{(t)}$  is the total number of documents the user  $i$  has in topic  $t$ . After removing all the terms that do not depend on ' $k$ ', we have:

$$p(y_i = k | \mathbf{Y}_{\neg i}, \mathbf{z}_i, \mathbf{Z}_{\neg i}, \mathbf{X}) \propto \left( \frac{\prod_{t=1}^T \prod_{j=0}^{m_{y_i,i}^{(t)}-1} (m_{k,\neg i}^{(t)} + \alpha_t + j)}{\prod_{j=0}^{(\sum_{t=1}^T m_{y_i,i}^{(t)})-1} (\sum_{t=1}^T m_{k,\neg i}^{(t)} + \alpha_t + j)} \right) \times (n_{\neg i}^{(k)} + v_k) \quad (17)$$

We now compute  $p(z_i = t | y_i, \mathbf{Z}_{\neg i}, \mathbf{X})$ . In this case,  $i \equiv (p, q)$  refers to the  $p$ -th document of the  $q$ -th user. All documents of a user belong to the same group. For ease of notation, now  $z_i$  and  $y_i$  now refer to the topic and user-group associated the  $i$ -th document respectively. We observe that once  $y_i$  is given, we know the group to which  $z_i$  belongs and it is independent of the topic distributions of other groups.  $z_i$  is also independent of all documents in other groups. Therefore,  $p(z_i = t | y_i, \mathbf{Z}_{\neg i}, \mathbf{X}) = p(z_i = t | y_i, \mathbf{Z}_{\neg i}^{y_i}, \mathbf{X}^{y_i})$  where  $(\mathbf{Z}_{\neg i}^{y_i}, \mathbf{X}^{y_i})$  refers to the data in the group  $y_i$  only. The following derivation is analogous to (6) with only difference being the superscripts  $y_i$ .

$$\begin{aligned}
p(z_i = t | y_i, \mathbf{Z}_{\neg i}^{y_i}, \mathbf{X}^{y_i}) &= \frac{p(y_i, \mathbf{Z}_{\neg i}^{y_i}, \mathbf{X}^{y_i})}{p(y_i, \mathbf{Z}_{\neg i}^{y_i}, \mathbf{X}^{y_i})} = \frac{p(\mathbf{X}^{y_i} | \mathbf{Z}_{\neg i}^{y_i})}{p(\mathbf{X}_{\neg i}^{y_i} | \mathbf{Z}_{\neg i}^{y_i}) p(\mathbf{w}_i^{y_i})} \cdot \frac{p(\mathbf{Z}_{\neg i}^{y_i})}{p(\mathbf{Z}_{\neg i}^{y_i})} \\
&\propto \frac{\Delta(\mathbf{n}_t^{y_i} + \boldsymbol{\beta})}{\Delta(\mathbf{n}_{t,\neg i}^{y_i} + \boldsymbol{\beta})} \cdot \frac{\Delta(\mathbf{m}_{y_i} + \boldsymbol{\alpha})}{\Delta(\mathbf{m}_{y_i,\neg i} + \boldsymbol{\alpha})} \\
&= \frac{\prod_{m \in X_i} \Gamma(n_t^{y_i,m} + \beta_m) \Gamma(\sum_{m=1}^V n_{t,\neg i}^{y_i,m} + \beta_m)}{\prod_{m \in X_i} \Gamma(n_{t,\neg i}^{y_i,m} + \beta_m) \Gamma(\sum_{m=1}^V n_t^{y_i,m} + \beta_m)} \\
&\quad \times \frac{\Gamma(m_{y_i}^{(t)} + \alpha_t) \Gamma(\sum_{t=1}^T m_{y_i,\neg i}^{(t)} + \alpha_t)}{\Gamma(m_{y_i,\neg i}^{(t)} + \alpha_t) \Gamma(\sum_{t=1}^T m_{y_i}^{(t)} + \alpha_t)} \\
&= \frac{\prod_{m \in X_i} \prod_{j=0}^{f_m} (n_{t,\neg i}^{y_i,m} + \beta_m + j)}{\prod_{j=0}^{length(X_i)} (\sum_{m=1}^V n_{t,\neg i}^{y_i,m} + \beta_m + j)} \times \frac{m_{y_i,\neg i}^{(t)} + \alpha_t}{[\sum_{t=1}^T m_{y_i}^{(t)} + \alpha_t] - 1}
\end{aligned}$$

$$\propto \frac{\prod_{m \in X_i} \prod_{j=0}^{f_m} (n_{t,-i}^{y_i, m} + \beta_m + j)}{\prod_{j=0}^{length(X_i)} (\sum_{m=1}^V n_{t,-i}^{y_i, m} + \beta_m + j)} \times (m_{y_i, -i}^{(t)} + \alpha_t) \quad (18)$$

In the above,  $n_t^{y_i, m}$  is the count of occurrences of word  $m$  in topic  $t$  in group  $y_i$ .  $m_{y_i}^{(t)}$  is the total number of documents of topic  $t$  in group  $y_i$ . ' $m_{y_i, -i}^{(t)}$ ' indicates that the document  $i$  was excluded.  $f_m$  is the count of occurrences of word  $m$  in document  $i$ .  $length(X_i)$  is the total count of words in document  $i$ .  $\{m \in X_i\}$  is the set of unique words in document  $i$ .

The multinomial parameters  $\theta_k$  for the group  $k$  can be computed as:

$$\begin{aligned} p(\theta_k | \mathbf{Z}^{(k)}, \alpha) \\ &= \frac{1}{Z_{\theta_k}} \prod_{i=1}^{N^{(k)}} p(z_i^{(k)} | \theta_k) \cdot p(\theta_k | \alpha) = Dir(\mathbf{m}_k + \alpha), \\ \mathbf{m}_k &= \{m_k^{(t)}\}_{t=1}^T \end{aligned}$$

Here,  $N^{(k)}$  is the number of users in group  $k$ , and  $\mathbf{Z}^{(k)}$  are the membership assignments for all users in group  $k$ . As already mentioned earlier,  $m_k^{(t)}$  is the number of times topic  $t$  has been observed in group  $k$  across all users belonging to the group  $k$ .

Now, we compute the parameters  $\mathbf{T}$ . The word distribution within each topic is dependent on the entire dataset since they are shared by all groups.

$$\begin{aligned} p(\tau_t | \mathbf{Y}, \mathbf{Z}, \mathbf{X}, \alpha, \beta, v) \\ &= \frac{1}{Z_{\tau_t}} \prod_{\{i: z_i = t\}} p(x_i | \tau_t) \cdot p(\tau_t | \beta) = Dir(\mathbf{n}_t + \beta), \\ \mathbf{n}_t &= \{n_t^{(m)}\}_{m=1}^V \end{aligned}$$

The group distribution can now be computed as:

$$\begin{aligned} p(\pi | \mathbf{Y}, \mathbf{Z}, \mathbf{X}, \alpha, \beta, v) &= p(\pi | \mathbf{Y}, v) = \frac{1}{Z_\pi} \prod_{\{i=1..N\}} p(y_i | \pi) \cdot p(\pi | v) = Dir(\mathbf{n} + v), \\ \mathbf{n} &= \{n^{(k)}\}_{k=1}^K \end{aligned}$$

Since both the above are Dirichlet distributions, we can compute the expectations easily:

$$\theta_k^t = \frac{m_k^t + \alpha_t}{\sum_{t'=1}^T m_k^{t'} + \alpha_{t'}} \quad (19)$$

$$\tau_{t,m} = \frac{n_t^{(m)} + \beta_m}{\sum_{m'=1}^V n_t^{(m')} + \beta_{m'}} \quad (20)$$

$$\pi_k = \frac{n^{(k)} + v_k}{\sum_{k'=1}^K n^{(k')} + v_{k'}} \quad (21)$$

## **References**

- [1] Heinrich, Gregor. Parameter Estimation for Text Analysis, 2008