

Name : Tang Sophearin

Affiliation : Department of Intelligence Science and Technology, Nishida Laboratory

Student ID : 6930-30-3963

# Latent Dirichlet Allocation: Parameter Estimation Using Gibbs Sampling

## Latent Dirichlet Allocation

Latent Dirichlet Allocation is a generative probabilistic model used for topic modelling for documents with sets of topics. LDA generate topics for each document through the topic distribution and words from each topic through that distribution.

Each document consists of mixture topics and each topic is a probability distribution over words. Therefore, the probability of  $i$ th word in a given document having  $T$  topics is

$$P(w_i) = \sum_{j=1}^T P(w_i|z_i = j)P(z_i = j), \quad [1]$$

where:

$P(w_i|z_i = j)$  : probability of  $i$ th word generated from topics  $j$ ;

$P(z_i = j)$  : probability of choosing a word from topics  $j$  in the current document.

LDA is different from a standard Bayesian classifier model. Instead of considering all words in a given document is generated from a single topic, LDA assumes words can be generated from multiple topics in a single document (soft classification).

There are two parameters to be estimated in LDA:  $\phi$  and  $\theta$ . Given  $D$  documents with  $T$  topics consisting of  $W$  words,  $\phi$  is  $T$  multinomial distributions over  $W$  words and  $\theta$  is  $D$  multinomial distributions over  $T$  topics. The main goal here is to find optimal  $\phi$  distributions which gives high probability to the words appearing in some corpus. One might consider obtaining this estimate by maximizing  $P(w|\phi, \theta)$  applying Expectation-Maximization to equation 1. However, this technique is slow to converge and has local maxima problem. To address this problem, LDA makes assumption about the prior probability distribution of  $\theta$ . Therefore, in LDA, to generate a new document, first we decide  $\theta$  a Dirichlet distribution and pick some topics; then, words are generated from  $\phi^{(j)}$  for each topic. Therefore, we need to maximize  $P(w|\phi, \alpha) = \int P(w|\phi, \theta)P(\theta|\alpha) d\theta$ , where  $P(\theta)$  is a Dirichlet( $\alpha$ ) distribution.

## Gibbs Sampling

Instead of estimating  $\phi$  and  $\theta$  explicitly, Gibbs sampling method uses Monte Carlos procedure to find posterior distribution  $P(z|w)$ ; then, by examining this posterior, we can obtain  $\phi$  and  $\theta$ .

Probability model is:

$$\begin{aligned} w_i | z_i, \phi^{(z_i)} &\sim \text{Discrete}(\phi^{(z_i)}) \\ \phi &\sim \text{Dirichlet}(\beta) \\ z_i | \theta^{(d_i)} &\sim \text{Discrete}(\theta^{(d_i)}) \\ \theta &\sim \text{Dirichlet}(\alpha) \end{aligned}$$

Here,  $\alpha$  and  $\beta$  are hyperparameters. The  $\alpha$  controls the mixture of topics for any given document. Turn it down for less of a mixture of topics; turn it up for more mixture. The  $\beta$  hyperparameter controls the distribution of words per topic. Turn it down and the topics will have less words; turn it up and the topics will have more words.

Unfortunately, the following posterior distribution cannot be calculated in closed form solution because the denominators does not factorize and involves  $T^n$  terms ( $n$  is number of words).

$$P(z|w) = \frac{P(w, z)}{\sum_z P(w, z)}.$$

To deal with this problem, we use Markov chain Monte Carlo to sample (Gibbs sampling). Target distribution is obtained by iterating through,

$$P(z_i = j | z_{-i}, w) \propto \frac{n_{-i,j}^{(w)} + \beta}{n_{-i,j}^{(\cdot)} + W\beta} \frac{n_{-i,j}^{(d)} + \alpha}{n_{-i,\cdot}^{(d)} + T\alpha},$$

where  $n_{-i,j}^{(\cdot)}$  is a word count not including  $z_i$ . The first part represents the probability of words under each topic, and the second part represents the probability of topics under each document.  $n$  are number of word counts under topics/documents;  $W$  is number of unique instances for words;  $T$  is number topics. For each iteration, we use state from previous iteration and sample  $z_i$  to obtain another new state. Finally, the posterior distribution converges to desired target distribution.

$\phi$  and  $\theta$  can then be computed by

$$\begin{aligned} \hat{\phi}_j^{(w)} &= \frac{n_j^{(w)} + \beta}{n_j^{(\cdot)} + W\beta} \\ \hat{\theta}_j^{(d)} &= \frac{n_j^{(d)} + \alpha}{n_{\cdot}^{(d)} + T\alpha}. \end{aligned}$$

Gibbs sampling in LDA algorithm performs quite well against other parameter estimation algorithms such as variational Bayes and expectation propagation while it is also computationally efficient because it only caches the number of times a word assigned to a topic and the number of times a topic occurs in a document for each iteration of Gibbs sampling.