

# Grounding Topic Models with Knowledge Bases: Supplementary Material

## A Inference via Collapsed Gibbs Sampling

Here we describe the inference algorithm for LGSA based on collapsed Gibbs Sampling.

Given a document corpus  $\mathcal{D}$ , the informative priors over entities  $\{\mathbf{p}^\eta, \mathbf{p}^\zeta, \lambda^\eta, \lambda^\zeta\}$ , and the hyperparameters  $\{\alpha, \beta\}$ , LGSA specifies the following full posterior distribution:

$$p(\Lambda, \eta, \zeta, \theta, \mathbf{z}, \mathbf{r}, \mathbf{e}, \mathbf{y} | \mathcal{D}, \mathbf{p}^\eta, \mathbf{p}^\zeta, \lambda^\eta, \lambda^\zeta, \alpha, \beta) \propto \left( p(\Lambda | \beta) p(\eta | \mathbf{p}^\eta, \lambda^\eta) p(\zeta | \mathbf{p}^\zeta, \lambda^\zeta) p(\theta | \alpha) p(\mathbf{z} | \theta) \right. \quad (\text{A.1})$$

$$\left. p(\mathbf{r}, \mathbf{e} | \mathbf{z}, \Lambda) p(\mathbf{m}_{\mathcal{D}} | \mathbf{e}, \eta) p(\mathbf{y} | \mathbf{M}_{\mathcal{D}}) p(\mathbf{w}_{\mathcal{D}} | \mathbf{y}, \mathbf{e}, \zeta) \right).$$

where  $\mathbf{m}_{\mathcal{D}}$  and  $\mathbf{w}_{\mathcal{D}}$  are the mentions and words in the document collections, respectively;  $\mathbf{M}_{\mathcal{D}} = \{M_1, \dots, M_D\}$  is the counts of mentions of the documents. The constant of proportionality is the marginal likelihood of the observed data.

The task of posterior inference for LGSA is to determine the probability distribution of the hidden variables given the observed mentions and words. However, exact inference is intractable due to the difficulty of calculating the normalizing constant in the above posterior distribution.

We use collapsed Gibbs Sampling, a well-established Markov chain Monte Carlo (MCMC) technique for approximate inference. In collapsed Gibbs Sampling, the distributions  $\Phi = \{\Lambda, \eta, \zeta, \theta\}$  are first marginalized (collapsed), a Markov chain over the latent indicators  $\{\mathbf{z}, \mathbf{r}, \mathbf{e}, \mathbf{y}\}$  is then constructed, whose stationary distribution is the posterior. We obtain samples of latent variables from the Markov chain. Point estimates for the collapsed distributions  $\Phi$  can then be computed given the samples, and predictive distributions are computed by averaging over multiple samples.

### Sampling Procedure

Gibbs Sampler repeatedly samples each latent variable conditioned on the current states of other hidden variables and observations; a configuration of latent states of the system is then obtained. Next we provide the derivation of the sampling formulas (i.e. Eqs.(1-4) in the paper).

By marginalizing out  $\Phi$  in Eq.(A.1), we obtain:

$$p(\mathbf{z}, \mathbf{r}, \mathbf{e}, \mathbf{y} | \cdot) \propto p(\mathbf{z} | \alpha) P(\mathbf{r}, \mathbf{e} | \mathbf{z}, \beta) p(\mathbf{m}_{\mathcal{D}} | \mathbf{e}, \mathbf{p}^\eta, \lambda^\eta) p(\mathbf{y} | \mathbf{M}_{\mathcal{D}}) \cdot p(\mathbf{w}_{\mathcal{D}} | \mathbf{y}, \mathbf{e}, \mathbf{p}^\zeta, \lambda^\zeta) = \int p(\theta | \alpha) P(\mathbf{z} | \theta) d\theta \int p(\Lambda | \beta) p(\mathbf{r}, \mathbf{e} | \mathbf{z}, \Lambda) d\Lambda \quad (\text{A.2})$$

$$\cdot \int p(\eta | \mathbf{p}^\eta, \lambda^\eta) p(\mathbf{m}_{\mathcal{D}} | \mathbf{e}, \eta) d\eta \cdot p(\mathbf{y} | \mathbf{M}_{\mathcal{D}}) \cdot \int p(\zeta | \mathbf{p}^\zeta, \lambda^\zeta) p(\mathbf{w}_{\mathcal{D}} | \mathbf{y}, \mathbf{e}, \zeta) d\zeta.$$

The conditional of  $z_{dj}$  can be computed as:

$$p(z_{dj} = z | e_{dj} = e, \mathbf{r}_{-dj}, \mathbf{z}_{-dj}, \cdot) \propto p(z_{dj} = z | \mathbf{z}_{-dj}, \alpha) p(e_{dj} = e | z_{dj} = z, \mathbf{r}_{-dj}, \cdot) \quad (\text{A.3})$$

The first term of Eq.(A.3) is:

$$p(z_{dj} = z | \mathbf{z}_{-dj}, \alpha) = \frac{p(z_{dj} = z, \mathbf{z}_{-dj} | \alpha)}{p(\mathbf{z}_{-dj} | \alpha)}. \quad (\text{A.4})$$

As we assume each  $z$  is generated from a multinomial distribution  $\theta$ , and the hyperparameter for conjugate Dirichlet prior is  $\alpha$ , we have:

$$p(\mathbf{z} | \alpha) = \int P(\theta | \alpha) P(\mathbf{z} | \theta) d\theta = \int \prod_d \frac{\Gamma(K\alpha)}{\prod_z \Gamma(\alpha)} \prod_z \theta_{dz}^{\alpha-1} \cdot \prod_d \prod_z \theta_{dz}^{n_d^{(z)}} d\theta = \prod_d \frac{\Gamma(K\alpha)}{\prod_z \Gamma(\alpha)} \cdot \frac{\prod_z \Gamma(n_d^{(z)} + \alpha)}{\Gamma(n_d^{(\cdot)} + K\alpha)},$$

where  $n_d^{(z)}$  is the number of times that topic  $z$  has been associated with a mention of document  $d$ . Marginal counts are represented with dots (i.e.  $n_d^{(\cdot)}$  is obtained by marginalizing  $n_d^{(z)}$  over  $z$ ). Combining the above equation with Eq.(A.4) leads to:

$$\frac{p(z_{dj} = z, \mathbf{z}_{-dj} | \alpha)}{p(\mathbf{z}_{-dj} | \alpha)} = \frac{\Gamma(n_d^{(z)} + \alpha) \Gamma(n_{d,-dj}^{(\cdot)} + K\alpha)}{\Gamma(n_{d,-dj}^{(z)} + \alpha) \Gamma(n_d^{(\cdot)} + K\alpha)} \quad (\text{A.5})$$

$$= \frac{n_{d,-dj}^{(z)} + \alpha}{n_{d,-dj}^{(\cdot)} + K\alpha},$$

where the count with subscript  $-ij$  denotes a quantity with the current instance (i.e. mention  $m_{dj}$ ) excluded. Here we use the identity  $\Gamma(x+1) = x\Gamma(x)$ .

The second term of Eq.(A.3) is the probability of generating entity  $e$  conditioned on topic  $z$ , which requires summing over the probabilities of all paths in  $z$  that could have generated  $u$ :

$$p(e_{dj} = e | z_{dj} = z, \mathbf{r}_{-dj}, \cdot) = \sum_{\mathbf{r}(e \in \mathbf{r})} p(\mathbf{r} | \mathbf{r}_{-dj}, z_{dj} = z, \cdot). \quad (\text{A.6})$$

The probability of a path  $\mathbf{r}$  is the product of the topic-specific transition probabilities along the path from root  $c_0$  to leaf  $c_{|\mathbf{r}|-1}$  (i.e. entity  $e$ ):

$$p(\mathbf{r} | \mathbf{r}_{-dj}, z_{dj} = z, \cdot) = \prod_{h=0}^{|\mathbf{r}|-2} p(c_{h+1} | c_h, z_{dj} = z, \mathbf{r}_{-dj} \cdot).$$

Here  $p(c_{h+1} | c_h, z_{dj} = z, \mathbf{r}_{-dj} \cdot)$  can be derived analogously to Eqs.(A.4-A.5), where the Dirichlet-Multinomial conjugates ensure the tractability of the integrals. We then obtain:

$$p(c_{h+1} | c_h, z_{dj} = z, \mathbf{r}_{-dj} \cdot) = \frac{n_{c_h, c_{h+1}}^{(z), -dj} + \beta}{n_{c_h, \cdot}^{(z), -dj} + |C(c_h)|\beta}, \quad (\text{A.7})$$

where  $n_{c_h, c_{h+1}}^{(z), -dj}$  is the number of paths in topic  $z$  that go from  $c_h$  to  $c_{h+1}$ , with the path of mention  $m_{dj}$  excluded.

Finally, by combining Eqs.(A.3-A.7) we obtain the sampling formula for  $z_{dj}$  as Eq.(1-2) in the paper. Note that we omit the subscripts/superscripts  $-ij$  in Eq.(1-2) to avoid cluttering of notation. Eq.(3) and Eq.(4) are derived in a similar manner.

### Distribution Estimation

After a sufficient number of sampling iterations as described above, we obtain a set of samples. The unknown distributions can then be computed by integrating across the samples. Specifically, for any single sample we can estimate  $\theta$ ,  $\theta'$ ,  $\Lambda$ ,  $\eta$  and  $\zeta$  as:

paths:

$$\hat{\phi}_{ze} = \sum_{\mathbf{r}(e \in \mathbf{r})} \prod_{h=0}^{|\mathbf{r}|-2} \hat{\Lambda}_{z c_h c_{h+1}},$$

$$\hat{\tau}_{zc} = \sum_{\mathbf{r}(e \in \mathbf{r})} \prod_{h=0}^{c_{h+1}=c} \hat{\Lambda}_{z c_h c_{h+1}}.$$

$$\hat{\theta}_{dz} = \frac{n_d^{(z)} + \alpha}{n_d^{(\cdot)} + K\alpha},$$

$$\hat{\theta}'_{de} = \frac{n_d^{(e)}}{n_d^{(\cdot)}},$$

$$\hat{\Lambda}_{zcc'} = \frac{n_{c, c'}^{(z)} + \beta}{n_{c, \cdot}^{(z)} + |C(c)|\beta},$$

$$\hat{\eta}_{em} = \frac{n_e^{(m)} + \lambda^\eta p_{em}^\eta}{n_e^{(\cdot)} + \lambda^\eta},$$

$$\hat{\zeta}_{ew} = \frac{n_e^{(w)} + \lambda^\zeta p_{ew}^\zeta}{n_e^{(\cdot)} + \lambda^\zeta}.$$

Finally, based on the estimated  $\hat{\Lambda}$ , we can compute the topic representations  $(\phi, \tau)$  by summing over all possible

It is straightforward to see that  $\sum_e \hat{\phi}_{ze} = 1$ . That is,  $\hat{\phi}_z$  is a distribution over entities.

### Inference on New Documents

Given a new-arriving document  $d$ , we can infer its topic distribution  $\theta_d$  and entity distribution  $\theta'_d$  to reveal its major themes and entities. The inference can be carried out using the Gibbs Sampling described above, but this time with the topic and entity statistics (i.e.  $\Lambda$ ,  $\eta$  and  $\zeta$ ) fixed.