# Supplementary Information for "Paragraph-citation Topic Models for Corpora with Citations: An Application to the United States Supreme Court"[*]

ByungKoo Kim[†][‡]    Saki Kuzushima[†][§]    Yuki Shiraito[¶]

This draft: July 13, 2022

# Contents

# A  Data preprocessing

Results of topic models can be highly sensitive to how data is preprocessed (Denny and Spirling, 2018). In addition to the simple preprocessing steps we introduced in Section 2, we removed words that appear very commonly across documents. The list of these words are "Statue","Supp","Ann","Rev","Stat","Judgment","Reverse","Follow","Certiorari" and "Opinion". While words such as "Follow" or "Reverse" could convey certain contexts, in legal opinions they are typically used to define how the drafted opinion stands in relation to precedents, and we believe they do not contain useful information with respect to topic discovery. In addition, words such as "Supp" or "Ann" are short words for Supplementary and Annex, which are specific collection of legal documents and thus removed for a better detection of topics.

Since common terms can vary by different subsets, we made additional preprocessing for each subset we used for application of our model. For each subset, we removed terms that appear too frequently as well as terms that appear too infrequently. Terms too common across documents for Privacy subset include "agent", "month","level" and "unfair" and for Voting Rights subset the removed words include"Vote", "Voter","Elect" and "Candid". For both subsets, terms that were too uncommon turned out to be simple typos or names of people or institutions such as "Rawlinson". The above process removed about 40% of the terms.

# B  Model inference: collapsed Gibbs sampler

$$\boldsymbol{\eta}_i \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

$$z_{ip} \sim \mathrm{Multinom}(1, \mathrm{softmax}(\boldsymbol{\eta}_i))$$

$$\boldsymbol{\Psi}_k \sim \mathrm{Dirichlet}(\boldsymbol{\beta})$$

$$\mathbf{w}_{ip} \sim \mathrm{Multinom}(N_{ip}, \boldsymbol{\Psi}_{z_{ip}})$$

$$D^*_{ipj} \sim \mathcal{N}(\boldsymbol{\tau}^T \mathbf{x}_{ipj}, 1)$$

$$D_{ipj} = \begin{cases} 1 \text{ if } D^*_{ipj} \geq 0 \\ 0 \text{ if } D^*_{ipj} < 0 \end{cases} \tag{1}$$

where $\mathbf{x}_{ipj} = \{1, \kappa_j^{(i)}, \eta_{j, z_{ip}}\}$

$\boldsymbol{\mu}$ and $\boldsymbol{\tau}$ each have hyperpriors assigned as

$$\boldsymbol{\mu} \sim \mathrm{MVN}(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0)$$

$$\boldsymbol{\tau} \sim \mathrm{MVN}(\boldsymbol{\mu}_\tau, \boldsymbol{\Sigma}_\tau) \tag{2}$$

The full posterior is denoted as follows.

$$p(\boldsymbol{\eta}, \boldsymbol{\Psi}, \mathbf{Z}, \boldsymbol{\tau} | \mathbf{W}, \mathbf{D}) \propto p(\boldsymbol{\mu}|\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0) p(\boldsymbol{\tau}|\boldsymbol{\mu}_\tau, \boldsymbol{\Sigma}_\tau) p(\boldsymbol{\eta}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) p(\boldsymbol{\Psi}|\boldsymbol{\beta}) p(\mathbf{Z}|\boldsymbol{\eta}) p(\mathbf{W}|\boldsymbol{\Psi}, \mathbf{Z}) p(\mathbf{D}|\mathbf{D}^*) p(\mathbf{D}^*|\boldsymbol{\tau}, \boldsymbol{\eta}, \mathbf{Z}, \mathbf{D}) \tag{3}$$

Unfortunately, the inference of the given posterior distribution is hard due to the non-conjugacy between normal prior for $\boldsymbol{\eta}$ and the logistic transformation function (Blei and Lafferty, 2007). Variational inference is the most frequently employed tool to address this problem, with an additional advantage of computational speed. However, obtained parameters are for the variational distribution which is an approximation to the target posterior. Moreover, the quality of the approximation is often not sufficiently explored (Add citations here).

To remedy this problem, we follow the recent advances in the inference of CTM models (Held and Holmes, 2006; Chen et al., 2013; Linderman et al., 2015). We first partially collapse the posterior distribution by integrating out $\boldsymbol{\Psi}$. Then we introduce an auxiliary Polya-Gamma variable $\boldsymbol{\lambda}$ and augment the collapsed posterior. Partial collapsing and data augmentation enables us to use Gibbs sampling which is known to produce samples that converge to the exact posterior.

With $\boldsymbol{\Psi}$ integrated out, our new posterior is proportional to

$$\int_{\boldsymbol{\Psi}} p(\boldsymbol{\eta},\boldsymbol{\Psi},\mathbf{Z},\boldsymbol{\tau}|\mathbf{W},\mathbf{D}) \propto p(\boldsymbol{\mu}|\boldsymbol{\mu}_0,\boldsymbol{\Sigma}_0)p(\boldsymbol{\tau}|\boldsymbol{\mu}_\tau,\boldsymbol{\Sigma}_\tau)p(\boldsymbol{\eta}|\boldsymbol{\mu},\boldsymbol{\Sigma})p(\mathbf{Z}|\boldsymbol{\eta})p(\mathbf{W}|\mathbf{Z})p(\mathbf{D}|\mathbf{D}^*)p(\mathbf{D}^*|\boldsymbol{\tau},\boldsymbol{\eta},\mathbf{Z},\mathbf{D})$$
(4)

## B.1    Derivation of the conditional distribution for Z

For $ip$th paragraph, the conditional distribution of $z_{ip}$ is

$$p(z_{ip}^k=1|\mathbf{Z}_{-ip},\boldsymbol{\eta},\mathbf{W},\mathbf{D}^*) \propto p(z_{ip}^k=1|\boldsymbol{\eta}_i)p(\mathbf{W}_{ip}|z_{ip}^k=1,\mathbf{Z}_{-ip},\mathbf{W}_{-ip})\prod_{j=1}^{i-1}p(D_{ipj}^*|z_{ip}^k=1,\mathbf{Z}_{-ip},\boldsymbol{\tau},\boldsymbol{\eta},\kappa)$$
(5)

The first term is $\frac{e^{\eta_{ik}}}{\sum_l e^{\eta_{il}}}$ which is proportional to $e^{\eta_{ik}}$.

The form of second term warrants further elaboration. Integrating out $\boldsymbol{\Psi}$ as

$$\begin{aligned}
p(\mathbf{W}|\mathbf{Z}) &= \int_{\boldsymbol{\Psi}} p(\mathbf{W},\boldsymbol{\Psi}|\mathbf{Z})d\boldsymbol{\Psi} \\
&= \int_{\boldsymbol{\Psi}} p(\mathbf{W}|\boldsymbol{\Psi},\mathbf{Z})p(\boldsymbol{\Psi}|\mathbf{Z})d\boldsymbol{\Psi} \\
&= \int_{\boldsymbol{\Psi}} p(\mathbf{W}|\boldsymbol{\Psi},\mathbf{Z})p(\boldsymbol{\Psi})d\boldsymbol{\Psi}
\end{aligned}$$
(6)

for $ip$th paragraph with $k$th topic yields the following.

$$\begin{aligned}
p(\mathbf{W}_{ip}|z_{ip}^k=1,\mathbf{Z}_{-ip},\mathbf{W}_{-ip}) \propto \int_{\boldsymbol{\Psi}_k} &\Psi_{k1}^{\beta_1-1}\Psi_{k2}^{\beta_2-1}...\Psi_{kV}^{\beta_V-1}\prod_v \Psi_{kv}^{\sum_{l=1}^{n_{ip}}\mathbb{I}(W_{ipl}=v)} \\
&\times \prod_v \prod_{(i',p')\neq(i,p)} \Psi_{kv}^{\sum_{l=1}^{n_{i'p'}}\mathbb{I}(W_{i'p'l}=v)\mathbb{I}(z_{i'p'}^k=1)}d\boldsymbol{\Psi}_k
\end{aligned}$$
(7)

Here, $N_{ip}$ denotes the total number of words in $ip$th paragraph, and $n_{ip}$ denotes the total number of unique words in $ip$th paragraph. Let $C_k^v = \sum_{i=1}^N \sum_{p=1}^{N_{ip}} \sum_{l=1}^{n_{ip}} \mathbb{I}(W_{ipl}=v)\mathbb{I}(z_{ip}^k=1)$, and $c_{k,ip}^v = \sum_{l=1}^{n_{ip}} \mathbb{I}(W_{ipl}=v)\mathbb{I}(z_{ip}^k=1)$ then the above can be simplified as

$$\begin{aligned}
p(\mathbf{W}_{ip}|z_{ip}^k=1,\mathbf{Z}_{-ip},\mathbf{W}_{-ip}) &\propto \int_{\boldsymbol{\Psi}_k} \Psi_{k1}^{\beta_1+c_{k,ip}^1+c_{k,-ip}^1-1}\Psi_{k2}^{\beta_2+c_{k,ip}^2+c_{k,-ip}^2-1}...\Psi_{kV}^{\beta_V+c_{k,ip}^V+c_{k,-ip}^V-1}d\boldsymbol{\Psi}_k \\
&= \frac{\prod_v \Gamma(\beta_v+c_{k,ip}^v+c_{k,-ip}^v)}{\Gamma(\sum_v \beta_v+c_{k,ip}^v+c_{k,-ip}^v)}
\end{aligned}$$
(8)

Imagine a paragraph of 3 words $\mathbf{W}_{ip}=\{1,1,3\}$, two of the first word and one of the third

word. Then

$$p(\mathbf{W}_{ip}|z_{ip}^k = 1, \mathbf{Z}_{-ip}, \mathbf{W}_{-ip}) \propto \frac{\prod_v \Gamma(\beta_v + c_{k,ip}^v + c_{k,-ip}^v)}{\Gamma(\sum_v \beta_v + c_{k,ip}^v + c_{k,-ip}^v)} \tag{9}$$

The numerator is

$$\Gamma(\beta_1 + 2 + c_{k,-ip}^1)\Gamma(\beta_3 + 1 + c_{k,-ip}^3) \times \prod_{v \neq (1,3)} \Gamma(\beta_v + c_{k,-ip}^v)$$

$$= (\beta_1 + 1 + c_{k,-ip}^1)(\beta_1 + c_{k,-ip}^1)(\beta_3 + c_{k,-ip}^3) \times \prod_v \Gamma(\beta_v + c_{k,-ip}^v) \tag{10}$$

In the same sense, the denominator is

$$\Gamma(3 + \sum_v \beta_v + c_{k,-ip}^v) = (2 + \sum_v \beta_v + c_{k,-ip}^v)(1 + \sum_v \beta_v + c_{k,-ip}^v)(\sum_v \beta_v + c_{k,-ip}^v)\Gamma(\sum_v \beta_v + c_{k,-ip}^v) \tag{11}$$

Rearrange the above and we have

$$\frac{(\beta_1 + 1 + c_{k,-ip}^1)(\beta_1 + c_{k,-ip}^1)(\beta_3 + c_{k,-ip}^3)}{(2 + \sum_v \beta_v + c_{k,-ip}^v)(1 + \sum_v \beta_v + c_{k,-ip}^v)(\sum_v \beta_v + c_{k,-ip}^v)} \times \frac{\prod_v \Gamma(\beta_v + c_{k,-ip}^v)}{\Gamma(\sum_v \beta_v + c_{k,-ip}^v)} \tag{12}$$

The second term does not depend on $z_{ip}^k$. Then for $\mathbf{W}_{ip} = \{1,1,3\}$, we have

$$p(\mathbf{W}_{ip}|z_{ip}^k = 1, \mathbf{Z}_{-ip}, \mathbf{W}_{-ip}) \propto \frac{(\beta_1 + 1 + c_{k,-ip}^1)(\beta_1 + c_{k,-ip}^1)(\beta_3 + c_{k,-ip}^3)}{(2 + \sum_v \beta_v + c_{k,-ip}^v)(1 + \sum_v \beta_v + c_{k,-ip}^v)(\sum_v \beta_v + c_{k,-ip}^v)} \tag{13}$$

If a paragraph consists of only one word such that $W_{ip} = l$, the above changes to

$$p(\mathbf{W}_{ip}|z_{ip}^k = 1, \mathbf{Z}_{-ip}, \mathbf{W}_{-ip}) \propto \frac{\beta_l + c_{k,-ip}^l}{\sum_v \beta_v + c_{k,-ip}^v} \tag{14}$$

which matches with the form for the equivalent part in collapsed Gibbs for LDA (Porteous et al., 2008; Xiao and Stibor, 2010; Asuncion et al., 2012).

The third term $p(D_{ipj}^*|z_{ip}^k = 1, \mathbf{Z}_{-ip}, \boldsymbol{\tau}, \boldsymbol{\eta}, \boldsymbol{\kappa}) = \exp\{-\frac{1}{2}(D_{ipj}^* - (\tau_0 + \tau_1\kappa_j^{(i)} + \tau_2\eta_{j,z_{ip}}))^2\}$ is proportional to

$$\exp\left\{-\frac{1}{2}\left(\tau_2^2\eta_{jk}^2 + 2(\tau_0\tau_2 + \tau_1\tau_2\kappa_j^{(i)} - \tau_2 D_{ipj}^*)\eta_{jk}\right)\right\} \tag{15}$$

4

## B.2 Derivation of the conditional distribution for $\boldsymbol{\eta}$

$$p(\boldsymbol{\eta}|\mathbf{Z}, \mathbf{W}, \mathbf{D}) = \prod_{i=1}^{N} \left( \prod_{p=1}^{N_i} p(z_{ip}|\boldsymbol{\eta}_i) \right) \mathcal{N}(\boldsymbol{\eta}_i|\boldsymbol{\mu}, \boldsymbol{\Sigma}) \prod_{p=1}^{N_i} \prod_{j=1}^{i-1} p(D_{ipj}^*|\kappa, \boldsymbol{\eta}_i, \mathbf{Z})$$

$$= \prod_{i=1}^{N} \left( \prod_{p=1}^{N_i} \frac{e^{\eta_{i,z_{ip}}}}{\sum_{j=1}^{K} e^{\eta_{ij}}} \right) \mathcal{N}(\boldsymbol{\eta}_i|\boldsymbol{\mu}, \boldsymbol{\Sigma}) \prod_{p=1}^{N_i} \prod_{j=1}^{i-1} p(D_{ipj}^*|\kappa, \boldsymbol{\eta}_i, \mathbf{Z}) \tag{16}$$

Following Held and Holmes (2006), the likelihood for $\eta_{ik}$ conditioned on $\eta_{i,-k}$ is

$$\ell(\eta_{ik}|\eta_{i,-k}) = \prod_{p=1}^{N_i} \left( \frac{e^{\rho_{ik}}}{1 + e^{\rho_{ik}}} \right)^{z_{ip,k}} \left( \frac{1}{1 + e^{\rho_{ik}}} \right)^{1-z_{ip,k}}$$

$$= \frac{(e^{\rho_{ik}})^{t_{ik}}}{(1 + e^{\rho_{ik}})^{N_i}} \tag{17}$$

where $\rho_{ik} = \eta_{ik} - \log(\sum_{l \neq k} e^{\eta_{il}})$ and $t_{ik} = \sum_{p=1}^{N_i} \mathbb{I}(z_{ip} = k)$.

Then

$$p(\eta_{ik}|\eta_{i,-k}, \mathbf{Z}, \mathbf{W}, \mathbf{D}, \boldsymbol{\tau}) \propto \ell(\eta_{ik}|\eta_{i,-k}) \mathcal{N}(\eta_{ik}|\nu_{ik}, \sigma_k^2) p(D^*|\boldsymbol{\eta}, \boldsymbol{\tau}, \mathbf{Z}) \tag{18}$$

where

$$\nu_{ik} = \mu_k - \Lambda_{kk}^{-1} \boldsymbol{\Lambda}_{k,-k}(\boldsymbol{\eta}_{i,-k} - \boldsymbol{\mu}_{i,-k})$$
$$\sigma_k^2 = \boldsymbol{\Lambda}_{kk}^{-1}$$
$$\boldsymbol{\Lambda} = \boldsymbol{\Sigma}^{-1} \tag{19}$$

The third term can be rewritten with respect to $\boldsymbol{\eta}$ as

$$p(\mathbf{D}^*|\boldsymbol{\eta}, \boldsymbol{\tau}, \mathbf{Z}) = \prod_i \prod_p \prod_{j=1}^{i-1} \exp\left\{ -\frac{1}{2}\big(D_{ipj}^* - (\tau_0 + \tau_1 \kappa_j^{(i)} + \tau_2 \eta_{j,z_{ip}})\big)^2 \right\}$$

$$\propto \prod_i \prod_p \prod_{j=1}^{i-1} \exp\left\{ -\frac{1}{2(1/\tau_2^2)}\Big(\eta_{j,z_{ip}}^2 - 2\frac{D_{ipj}^* - \tau_0 - \tau_1 \kappa_j^{(i)}}{\tau_2} \eta_{j,z_{ip}}\Big) \right\}$$

$$\propto \prod_i \prod_p \prod_{j=1}^{i-1} \mathcal{N}(\eta_{j,z_{ip}}|\mu_{ipj}^*, \frac{1}{\tau_2^2})$$

$$= \prod_i \prod_p \prod_{j=1}^{i-1} \prod_k \mathcal{N}(\eta_{jk}|\mu_{ipj}^*, \frac{1}{\tau_2^2})^{\mathbb{I}(z_{ip}=k)} \tag{20}$$

where $\mu_{ipj}^* = \frac{D_{ipj}^* - \tau_0 - \tau_1 \kappa_j^{(i)}}{\tau_2}$. We notice that the above can be rewritten as a product of univariate

normal distributions such that

$$\prod_k \prod_{s=i+1}^{N} \prod_{p=1}^{N_s} \mathcal{N}(\eta_{ik}|\mu^*_{spi}, \sigma^{2*})^{\mathbb{I}(z_{sp}=k)}$$

$$\equiv \prod_{k=1}^{K} \mathcal{N}(\eta_{ik}|m_{ik}, V_{i,kk}) \tag{21}$$

$\mathbf{V}_i$ is a diagnoal matrix with the $k$th diagonal entry of the inverse of $\mathbf{V}_i$ (or $\mathbf{V}_i^{-1}$) as

$$V_{i,kk}^{-1} = \frac{1}{\sigma^{2*}} \sum_{s=i+1}^{N} \sum_{p=1}^{N_s} \mathbb{I}(z_{sp} = k)$$

$$= \tau_2^2 \sum_{s=i+1}^{N} \sum_{p=1}^{N_s} \mathbb{I}(z_{sp} = k) \tag{22}$$

The $k$th entry of $\mathbf{m}_i$ then is

$$m_{ik} = \frac{\tau_2^2 \sum_{s=i+1}^{N} \sum_{p=1}^{N_s} \mu^*_{spi} \mathbb{I}(z_{sp} = k)}{V_{i,kk}^{-1}}$$

$$= \frac{\sum_s \sum_p \mu^*_{spi} \mathbb{I}(z_{sp} = k)}{\sum_s \sum_p \mathbb{I}(z_{sp} = k)} \tag{23}$$

Then the $\eta$ conditional is

$$p(\eta_{ik}|\eta_{i,-k}, \mathbf{Z}, \mathbf{W}, \mathbf{D}, \boldsymbol{\tau}) \propto \ell(\eta_{ik}|\eta_{i,-k}) \mathcal{N}(\eta_{ik}|\nu_{ik}, \sigma_k^2) \mathcal{N}(\eta_{ik}|m_{ik}, V_{i,kk}) \tag{24}$$

We now introduce Polya-Gamma augmentation such that

$$p(\eta_{ik}|\eta_{i,-k}, \mathbf{Z}, \mathbf{W}, \mathbf{D}, \boldsymbol{\tau}, \lambda_{ik}) \propto \exp\{(t_{ik} - \frac{N_i}{2})\rho_{ik} - \frac{\lambda_{ik}}{2}\rho_{ik}^2\} \mathcal{N}(\eta_{ik}|\nu_{ik}, \sigma_k^2) \mathcal{N}(\eta_{ik}|m_{ik}, V_{i,kk})$$

$$\propto \mathcal{N}(\eta_{ik}|\frac{t_{ik} - N_i/2}{\lambda_{ik}} + \log(\sum_{l \neq k} e^{\eta_{il}}), 1/\lambda_{ik}) \mathcal{N}(\eta_{ik}|\nu_{ik}, \sigma_k^2) \mathcal{N}(\eta_{ik}|m_{ik}, V_{i,kk}) \tag{25}$$

Summing all of the above, the conditional distribution of $\eta_{ik}$ is

$$p(\eta_{ik}|\eta_{i,-k}, \mathbf{Z}, \mathbf{W}, \mathbf{D}, \boldsymbol{\tau}, \lambda_{ik}) \propto \mathcal{N}(\eta_{ik}|\tilde{\mu}_{ik}, \tilde{\sigma}_k^2) \tag{26}$$

where

$$\tilde{\sigma}_k^2 = (\sigma_k^{-2} + \lambda_{ik} + v_{i,kk}^{-1})^{-1}$$

$$\tilde{\mu}_{ik} = \tilde{\sigma}_k^2 \left( v_{i,kk}^{-1} m_{ik} + \sigma_k^{-2} \nu_{ik} + t_{ik} - \frac{N_i}{2} + \lambda_{ik} \log(\sum_{l \neq k} e^{\eta_{il}}) \right) \tag{27}$$

## Derivation of conditional distribution for $\boldsymbol{\lambda}$

The Gibbs sampling for the augmentation variable $\boldsymbol{\lambda}$ is obtained by collecting terms that include $\boldsymbol{\lambda}_i$ in the joint of $\boldsymbol{z}_i$ and $\boldsymbol{\eta}_i$.

$$p(\lambda_{ik}|\mathbf{Z}, \mathbf{W}, \boldsymbol{\eta}) \propto PG(N_i, \rho_{ik}) \tag{28}$$

## B.3  Derivation of conditional distribution for $\mathbf{D}^*$

$$p(D_{ipj}^*|\boldsymbol{\eta}, \mathbf{Z}, \boldsymbol{\tau}, \mathbf{D}) \propto \begin{cases} TN_{(0,\infty)}(\tau_0 + \tau_1\kappa_j^{(i)} + \tau_2\eta_{j,z_{ip}}, 1) & \text{if } D_{ipj} = 1 \\ TN_{(-\infty,0]}(\tau_0 + \tau_1\kappa_j^{(i)} + \tau_2\eta_{j,z_{ip}}, 1) & \text{if } D_{ipj} = 0 \end{cases} \tag{29}$$

## B.4  Derivation of conditional distribution for $\boldsymbol{\tau}$

Let $\mathbf{x}_{ipj} = [1, \kappa_j^{(i)}, \eta_{j,z_{ip}}]^T$ and $\boldsymbol{\tau} = [\tau_0, \tau_1, \tau_2]^T$

$$p(\boldsymbol{\tau}|\boldsymbol{\eta}, \mathbf{Z}, \mathbf{D}^*) \propto exp\left\{ -\frac{1}{2}\sum_{ipj}\left(D_{ipj}^* - \mathbf{x}_{ipj}^T\boldsymbol{\tau}\right)^2 \right\} N(\boldsymbol{\mu_\tau}, \Sigma_\tau)$$

$$\propto N(\tilde{\boldsymbol{\tau}}, \tilde{\Sigma}_\tau) \tag{30}$$

where $\tilde{\Sigma}_\tau = \left(\left(\sum_{ipj}\mathbf{x}_{ipj}\mathbf{x}_{ipj}^T\right) + \Sigma_\tau^{-1}\right)^{-1}$ and $\tilde{\tau} = \tilde{\Sigma}_\tau\left(\left(\sum_{ipj}\mathbf{x}_{ipj}^T D_{ipj}^*\right) + \Sigma_\tau^{-1}\boldsymbol{\mu_\tau}\right)$

## B.5  Recovering $\boldsymbol{\Psi}$

We estimate the integrated out parameter $\boldsymbol{\Psi}$ from our posterior samples as follows.

$$\hat{\Psi}_{kv} = \frac{\sum_i \sum_p \left(\beta_v + \mathbb{I}(z_{ip}^k = 1)\mathrm{W}_{ip,v}\right)}{\sum_i \sum_p \sum_l \left(\beta_l + \mathbb{I}(z_{ip}^k = 1)\mathrm{W}_{ip,l}\right)} \tag{31}$$

# C Initialization strategy for collapsed Gibbs sampler

Similar to any other Gibbs samplers, the quality of our collapsed Gibbs sampler depends highly on the initial values of the parameters. We propose to use LDA to generate stable initial values for $\boldsymbol{\eta}$, then use it to generate reasonable initial values for other parameters.

We first fit LDA with variational EM on document-level document-feature matrix to obtain $\hat{\boldsymbol{\theta}}$. For $i$th document,

$$
\begin{aligned}
z_{ip}^{(0)} &\sim \text{Categorical}(\hat{\boldsymbol{\theta}}_i) \quad \forall p = 1, 2, ..., N_i \\
\boldsymbol{\eta}_i^{(0)} &= \log(\hat{\boldsymbol{\theta}}_i / \hat{\theta}_{iK})
\end{aligned}
\tag{32}
$$

Set $\tilde{\tau}_0$, or the sparsity parameter, using the observed density of the citation matrix and randomly draw the other two parameters as

$$
\tilde{\tau}_0 = \frac{1}{2}\log(\text{density}(\mathbf{D}))
$$

$$
\tilde{\tau}_1, \tilde{\tau}_2 \sim \text{unif}(0, 1)
\tag{33}
$$

Sample $\mathbf{D}^*$ using the above parameters

$$
\begin{aligned}
{D_{ipj}^*}^{(0)} &\sim TN_{(-\infty,0)}(\tilde{\tau}_0 + \tilde{\tau}_1 \kappa_j^{(i)} + \tilde{\tau}_2 \eta_{j,z_{ip}^{(0)}}^{(0)}, 1) \quad \text{if } D_{ipj} = 0 \\
{D_{ipj}^*}^{(0)} &\sim TN_{[0,\infty)}(\tilde{\tau}_0 + \tilde{\tau}_1 \kappa_j^{(i)} + \tilde{\tau}_2 \eta_{j,z_{ip}^{(0)}}^{(0)}, 1) \quad \text{if } D_{ipj} = 1
\end{aligned}
\tag{34}
$$

Then set $\boldsymbol{\tau}^{(0)}$ again using MLE

$$
\boldsymbol{\tau}^{(0)} = (\sum_{ipj} \mathbf{x}_{ipj}^{(0)} {\mathbf{x}_{ipj}^{(0)}}^T)^{-1}(\sum_{ipj} {\mathbf{x}_{ipj}^{(0)}}^T {D_{ipj}^*}^{(0)})
\tag{35}
$$

where $\mathbf{x}_{ipj}^{(0)} = \{1, \kappa_j^{(i)}, \eta_{j,z_{ip}^{(0)}}^{(0)}\}$

Finally, set the values of $\boldsymbol{\lambda}^{(0)}$ by

$$
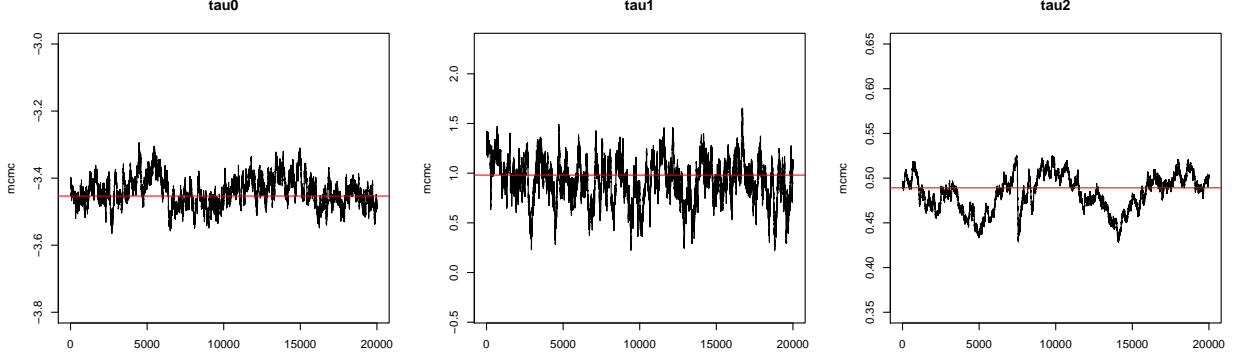\lambda_i^{(0)} \sim \text{PG}(N_i, \boldsymbol{\eta}_i^{(0)})
\tag{36}
$$

Figure D.1: MCMC convergence of $\boldsymbol{\tau}$ posterior samples in simulation. Horizontal red line indicates the true values of $\boldsymbol{\tau}$.



Figure D.2: MCMC convergence of $\boldsymbol{\theta}$ parameters for the first document. $\boldsymbol{\theta}$ values are obtained by transforming the posterior samples of $\boldsymbol{\eta}$ of the corresponding document. Horizontal red line indicates the true values of $\boldsymbol{\theta}$ for the first document for each topic. We do not display the MCMC convergence for other documents, but all documents show similar level of convergence to the true value of $\boldsymbol{\theta}$.

# D   More results on simulation
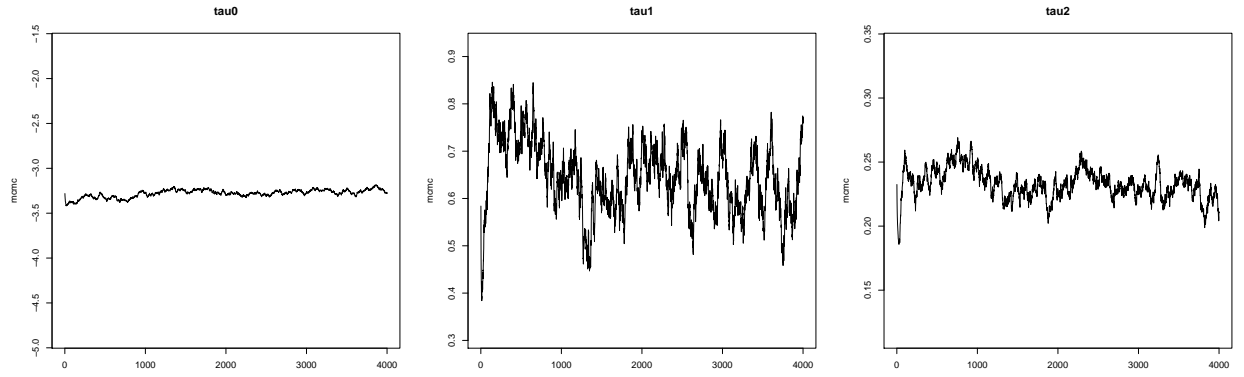
# E   More results on USSC application

Figure E.1: MCMC convergence of $\boldsymbol{\tau}$ posterior samples for the USSC application on Privacy issue area. Horizontal red line indicates the true values of $\boldsymbol{\tau}$.
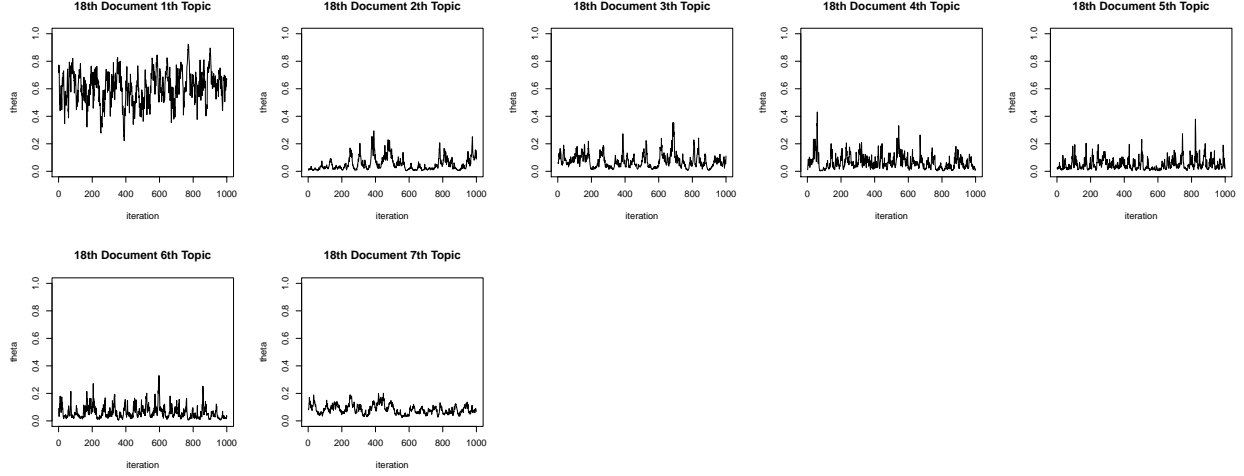
Figure E.2: MCMC convergence of $\boldsymbol{\theta}$ parameters for the 18th document in the subset of Privacy issue area. $\boldsymbol{\theta}$ values are obtained by transforming the posterior samples of $\boldsymbol{\eta}$ of the corresponding document. Horizontal red line indicates the true values of $\boldsymbol{\theta}$ for the 18th document for each topic. We do not display the MCMC convergence for other documents, but all documents show similar level of convergence to the true value of $\boldsymbol{\theta}$.
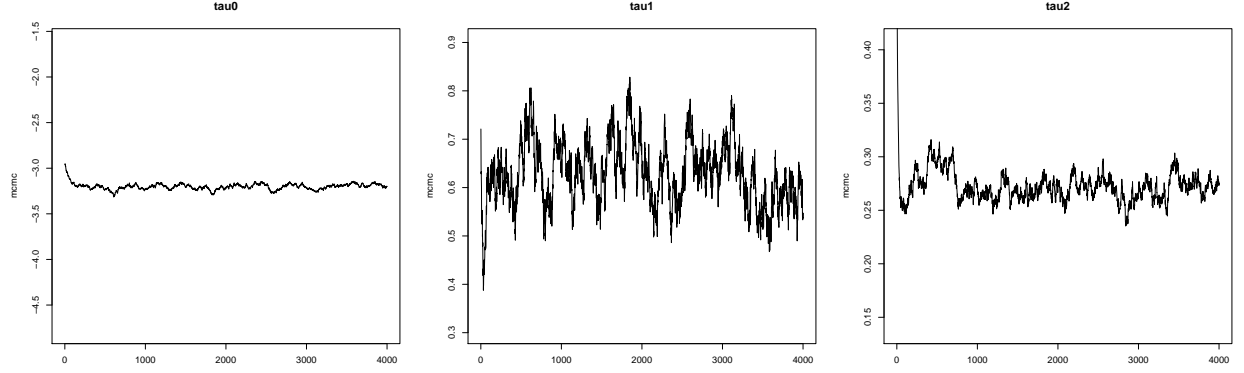
Figure E.3: MCMC convergence of $\boldsymbol{\tau}$ posterior samples for the USSC application on Voting Rights issue area. Horizontal red line indicates the true values of $\boldsymbol{\tau}$.



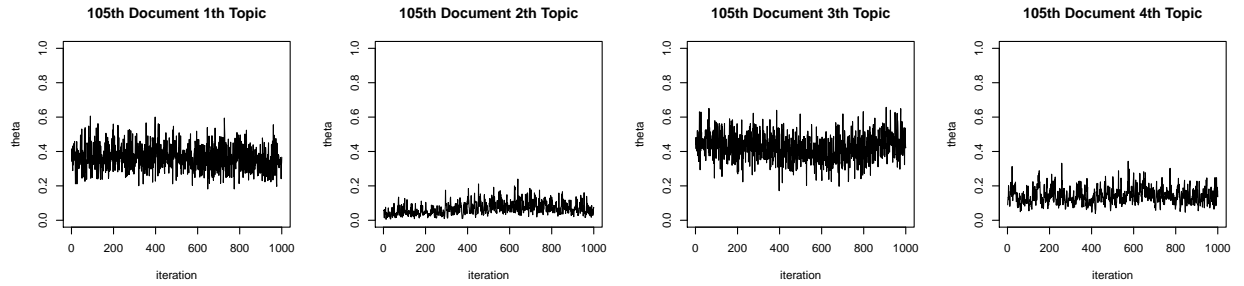Figure E.4: MCMC convergence of $\boldsymbol{\theta}$ parameters for the 105th document in the subset of Voting Rights issue area. $\boldsymbol{\theta}$ values are obtained by transforming the posterior samples of $\boldsymbol{\eta}$ of the corresponding document. Horizontal red line indicates the true values of $\boldsymbol{\theta}$ for the 105th document for each topic. We do not display the MCMC convergence for other documents, but all documents show similar level of convergence to the true value of $\boldsymbol{\theta}$.

# References

Asuncion, A., Welling, M., Smyth, P., and Teh, Y. W. (2012). On smoothing and inference for topic models. *arXiv preprint arXiv:1205.2662*.

Blei, D. M. and Lafferty, J. D. (2007). A correlated topic model of science. *The annals of applied statistics*, 1(1):17–35.

Chen, J., Zhu, J., Wang, Z., Zheng, X., and Zhang, B. (2013). Scalable inference for logistic-normal topic models. *Advances in neural information processing systems*, 26.

Denny, M. J. and Spirling, A. (2018). Text preprocessing for unsupervised learning: Why it matters, when it misleads, and what to do about it. *Political Analysis*, 26(2):168–189.

Held, L. and Holmes, C. C. (2006). Bayesian auxiliary variable models for binary and multinomial regression. *Bayesian analysis*, 1(1):145–168.

Linderman, S., Johnson, M. J., and Adams, R. P. (2015). Dependent multinomial models made easy: Stick-breaking with the pólya-gamma augmentation. *Advances in Neural Information Processing Systems*, 28.

Porteous, I., Newman, D., Ihler, A., Asuncion, A., Smyth, P., and Welling, M. (2008). Fast collapsed gibbs sampling for latent dirichlet allocation. In *Proceedings of the 14th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 569–577.

Xiao, H. and Stibor, T. (2010). Efficient collapsed gibbs sampling for latent dirichlet allocation. In *Proceedings of 2nd asian conference on machine learning*, pages 63–78. JMLR Workshop and Conference Proceedings.