

# Improving Topic Models with Latent Feature Word Representations

Dat Quoc Nguyen

Joint work with  
Richard Billingsley, Lan Du and Mark Johnson

Department of Computing  
Macquarie University  
Sydney, Australia

September 2015

# Introduction

- *Topic models* take a corpus of documents as input, and
  - ▶ learn a set of latent *topics* for the corpus
  - ▶ infer *document-to-topic* and *topic-to-word* distributions from co-occurrence of words within documents
- If the corpus is small and/or the documents are short, the topics will be noisy due to the limited information of word co-occurrence
- *Latent word representations* learnt from large external corpora capture various aspects of word meanings
  - ▶ We used the pre-trained word2vec (Mikolov et al., 2013) and Glove (Pennington et al., 2014) word representations

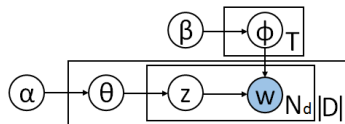
## High-level idea

- Use the word representations learnt on a large external corpus to improve the topic-word distributions in a topic model
  - ▶ we combine Latent Dirichlet Allocation (Blei et al., 2003) and Dirichlet Multinomial Mixture (Nigam et al., 2000) with the distributed representations
  - ▶ the improvement is greatest on small corpora with short documents

# LDA and DMM

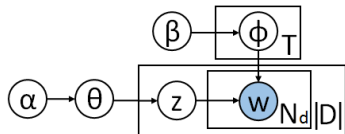
- Latent Dirichlet Allocation (LDA)

$$\begin{aligned}\theta_d &\sim \text{Dir}(\alpha) & z_{d_i} &\sim \text{Cat}(\theta_d) \\ \phi_z &\sim \text{Dir}(\beta) & w_{d_i} &\sim \text{Cat}(\phi_{z_{d_i}})\end{aligned}$$



- Dirichlet Multinomial Mixture (DMM) model: one-topic-per-document

$$\begin{aligned}\theta &\sim \text{Dir}(\alpha) & z_d &\sim \text{Cat}(\theta) \\ \phi_z &\sim \text{Dir}(\beta) & w_{d_i} &\sim \text{Cat}(\phi_{z_d})\end{aligned}$$



- Inference is typically performed with a *Gibbs sampler*, integrating out  $\theta$  and  $\phi$  (Griffiths et al., 2004; Yin and Wang, 2014)

# Latent-feature topic-to-word distributions

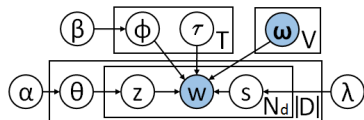
- We assume that each word  $w$  is associated with a *word vector*  $\omega_w$
- We learn a *topic vector*  $\tau_t$  for each topic  $t$
- We use these to define a latent feature topic-to-word distribution  $\text{CatE}(w)$  over words:

$$\text{CatE}(w \mid \tau_t \omega^T) \propto \exp(\tau_t \cdot \omega_w)$$

- ▶  $\tau_t \omega^T$  is a vector of unnormalized scores, one per word
- In our topic models, we *mix the CatE distribution* with a multinomial distribution over words
  - ▶ combine information from a large, general corpus (via the CatE distribution) and a smaller but more specific corpus (via the multinomial distribution)
  - ▶ use a Boolean *indicator variable* that records whether a word is generated from CatE or the multinomial distribution

# The Latent Feature LDA (LF-LDA) model

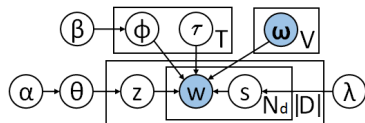
$$\begin{aligned}\theta_d &\sim \text{Dir}(\alpha) & z_{d_i} &\sim \text{Cat}(\theta_d) \\ \phi_z &\sim \text{Dir}(\beta) & s_{d_i} &\sim \text{Ber}(\lambda) \\ w_{d_i} &\sim (1 - s_{d_i})\text{Cat}(\phi_{z_{d_i}}) + s_{d_i}\text{Cat}(\tau_{z_{d_i}} \omega^\top)\end{aligned}$$



- Replace the topic-to-word Dirichlet multinomial component in LDA with a two-component mixture of a topic-to-word Dirichlet multinomial component and a latent feature topic-to-word component
- $s_{d_i}$  is the Boolean indicator variable indicating whether word  $w_{d_i}$  is generated from the latent feature component
- $\lambda$  is a user-specified hyper-parameter determining how often words are generated from the latent feature component
  - ▶ if we estimated  $\lambda$  from data, we expect it would never generate through the latent feature component

# The Latent Feature DMM (LF-DMM) model

$$\begin{aligned}\theta &\sim \text{Dir}(\alpha) & z_d &\sim \text{Cat}(\theta) \\ \phi_z &\sim \text{Dir}(\beta) & s_{d_i} &\sim \text{Ber}(\lambda) \\ w_{d_i} &\sim (1 - s_{d_i})\text{Cat}(\phi_{z_{d_i}}) + s_{d_i}\text{CatE}(\tau_{z_{d_i}} \omega^\top)\end{aligned}$$



- Replace the topic-to-word Dirichlet multinomial component in DMM with a two-component mixture of a topic-to-word Dirichlet multinomial component and a latent feature topic-to-word component
- $s_{d_i}$  is the Boolean indicator variable indicating whether word  $w_{d_i}$  is generated from the latent feature component
- $\lambda$  is a user-specified hyper-parameter determining how often words are generated from the latent feature component

## Inference for the LF-LDA model

- We integrate out  $\theta$  and  $\phi$  as in the Griffiths et al. (2004) sampler, and *interleave MAP estimation for  $\tau$  with Gibbs sweeps for the other variables*
- Algorithm outline:  
initialize the word-topic variables  $z_{d_i}$  using the LDA sampler  
repeat:
  - for each topic  $t$ :
    - use LBFGS to optimize the L2-regularized log-loss
    - $\tau_t = \arg \max_{\tau_t} P(\tau_t \mid \mathbf{z}, \mathbf{s})$
  - for each document  $d$  and each word location  $i$ :
    - sample  $z_{d_i}$  from  $P(z_{d_i} \mid \mathbf{z}_{-d_i}, \mathbf{s}_{-d_i}, \boldsymbol{\tau})$
    - sample  $s_{d_i}$  from  $P(s_{d_i} \mid \mathbf{z}, \mathbf{s}_{-d_i}, \boldsymbol{\tau})$



## Inference for the LF-DMM model

- We integrate out  $\theta$  and  $\phi$  as in the Yin and Wang (2014) sampler, and *interleave MAP estimation for  $\tau$  with Gibbs sweeps*
- Algorithm outline:  
initialize the word-topic variables  $z_{d_i}$  using the DMM sampler  
repeat:
  - for each topic  $t$ :
    - use LBFGS to optimize the L2-regularized log-loss
    - $\tau_t = \arg \max_{\tau_t} P(\tau_t \mid z, s)$
  - for each document  $d$ :
    - sample  $z_d$  and  $s_d$  from  $P(z_d, s_d \mid z_{-d}, s_{-d}, \tau)$
- Note:  $P(z_d, s_d \mid z_{-d}, s_{-d}, \tau)$  is *computationally expensive* to compute exactly, as it requires *summing over all possible values for  $s_d$*
- We approximate these probabilities by *assuming that the topic-word counts are “frozen”*, i.e., they don't increase within a document  
 $\Rightarrow$  We are able to integrate out  $s_d$

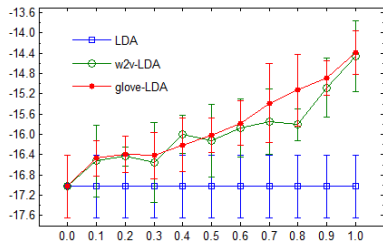
## Goals of evaluation

- A topic model learns document-topic and topic-word distributions:
  - ▶ *topic coherence* evaluates the topic-word distributions
  - ▶ *document clustering* and *document classification* evaluate the document-topic distribution
- Do the word2vec and Glove word vectors behave differently in topic modelling? (w2v-LDA, glove-LDA, w2v-LDA, glove-DMM)
- We expect that the latent feature component will have *the greatest impact on small corpora*, so our evaluation focuses on them:

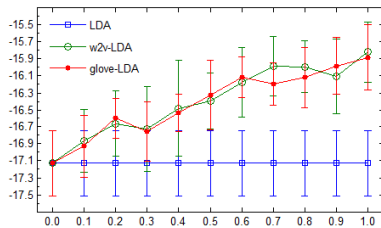
Dataset		# labels	# docs	words/doc	# types
N20	20 newsgroups	20	18,820	103.3	19,572
N20short	≤ 20 words	20	1,794	13.6	6,377
N20small	400 docs	20	400	88.0	8,157
TMN	TagMyNews	7	32,597	18.3	13,428
TMNtitle	TagMyNews titles	7	32,503	4.9	6,347
Twitter		4	2,520	5.0	1,390

## Topic coherence evaluation

- Lau et al. (2014) showed that *human scores on a word intrusion task* are highly correlated with the *normalized pointwise mutual information (NPMI)*
- We found latent feature vectors produced a *significant improvement of NPMI scores on all models and corpora*
  - ▶ greatest improvement when  $\lambda = 1$  (unsurprisingly)



20 topics



40 topics

NPMI scores on the N20short dataset, varying the mixture weight  $\lambda$  from 0.0 to 1.0.

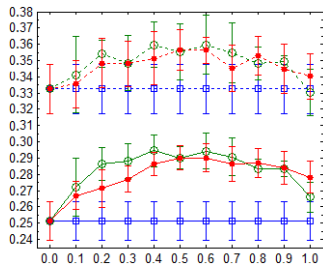
## w2v-DMM on TagMyNews titles corpus

Topic 1		Topic 3		Topic 4	
DMM	w2v-DMM	DMM	w2v-DMM	DMM	w2v-DMM
japan	japan	u.s.	prices	egypt	libya
nuclear	nuclear	oil	sales	<u>china</u>	egypt
u.s.	u.s.	japan	oil	u.s.	<b>iran</b>
crisis	plant	prices	u.s.	mubarak	<b>mideast</b>
plant	<b>quake</b>	stocks	profit	<u>bin</u>	<b>opposition</b>
<u>china</u>	radiation	sales	stocks	libya	<b>protests</b>
<u>libya</u>	<b>earthquake</b>	profit	japan	<u>laden</u>	<b>leader</b>
radiation	<b>tsunami</b>	<u>fed</u>	rise	<u>france</u>	<b>syria</b>
<u>u.n.</u>	<b>nuke</b>	rise	<b>gas</b>	bahrain	u.n.
<u>vote</u>	crisis	growth	growth	<u>air</u>	<b>tunisia</b>
<u>korea</u>	<b>disaster</b>	<u>wall</u>	shares	report	<b>chief</b>
<u>europe</u>	<b>power</b>	<u>street</u>	<b>price</b>	<u>rights</u>	<b>protesters</b>
<u>government</u>	<b>oil</b>	<u>china</u>	<b>profits</b>	<u>court</u>	mubarak
<u>election</u>	<b>japanese</b>	<u>fall</u>	<b>rises</b>	u.n.	<b>crackdown</b>
<u>deal</u>	<b>plants</b>	shares	<b>earnings</b>	<u>war</u>	bahrain

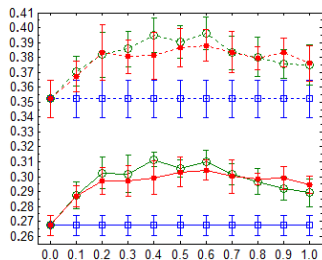
- Table shows the 15 most probable topical words found by 20-topic w2v-DMM on the TMNtitle corpus
- Words found by DMM but not by w2v-DMM are underlined
- Words found by w2v-DMM but not DMM are in bold

# Document clustering evaluation (1)

- Cluster documents by assigning them to the *highest probability topic*
- Evaluate clusterings by *purity* and *normalized mutual information* (NMI) (Manning et al., 2008)



20 topics



40 topics

Purity and NMI results on the N20short dataset, varying the mixture weight  $\lambda$  from 0.0 to 1.0.

- In general, best results with  $\lambda = 0.6$
- ⇒ Set  $\lambda = 0.6$  in all further experiments

## Document clustering evaluation (2)

Data	Method	Purity		NMI	
		T=4	T=20	T=4	T=20
Twitter	LDA	0.559 $\pm$ 0.020	0.614 $\pm$ 0.016	0.196 $\pm$ 0.018	0.174 $\pm$ 0.008
	w2v-LDA	<b>0.598</b> $\pm$ 0.023	<b>0.635</b> $\pm$ 0.016	<b>0.249</b> $\pm$ 0.021	<b>0.191</b> $\pm$ 0.011
	glove-LDA	0.597 $\pm$ 0.016	<b>0.635</b> $\pm$ 0.014	0.242 $\pm$ 0.013	<b>0.191</b> $\pm$ 0.007
	Improve.	<b>0.039</b>	<b>0.021</b>	<b>0.053</b>	<b>0.017</b>
Twitter	DMM	0.523 $\pm$ 0.011	0.619 $\pm$ 0.015	0.222 $\pm$ 0.013	0.213 $\pm$ 0.011
	w2v-DMM	<b>0.589</b> $\pm$ 0.017	0.655 $\pm$ 0.015	0.243 $\pm$ 0.014	0.215 $\pm$ 0.009
	glove-DMM	0.583 $\pm$ 0.023	<b>0.661</b> $\pm$ 0.019	<b>0.250</b> $\pm$ 0.020	<b>0.223</b> $\pm$ 0.014
	Improve.	<b>0.066</b>	<b>0.042</b>	<b>0.028</b>	<b>0.01</b>

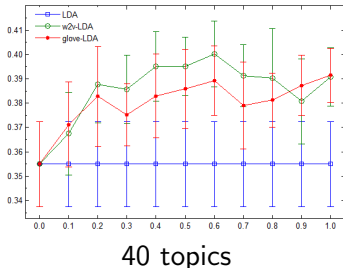
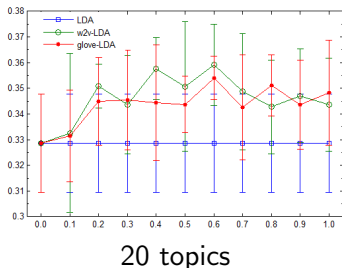
- On the short, our models obtain better clustering results than the baseline models:
  - ▶ on N20small, we get 6.0% improvement on NMI at  $T = 6$
  - ▶ on TMN and TMNtitle, we obtain 6.1% and 2.5% higher Purity at  $T = 80$

## Document clustering evaluation (3)

- For small  $T \leq 7$ , on the large datasets of N20, TMN and TMNtitle, our models and baseline models obtain similar clustering results
- With larger  $T$ , our models perform better than baselines on the short TMN and TMNtitle datasets. On the N20 dataset, the baseline LDA model obtains slightly higher clustering results than ours
- No reliable difference between word2vec and Glove vectors

# Document classification (1)

- Use SVM to predict the ground truth label from the topic-proportion vector of each document



$F_1$  scores on the N20short dataset, varying the mixture weight  $\lambda$  from 0.0 to 1.0.

Data	Method	$\lambda = 0.6$			
		T=6	T=20	T=40	T=80
N20small	LDA	0.204 $\pm$ 0.020	0.392 $\pm$ 0.029	0.459 $\pm$ 0.030	0.477 $\pm$ 0.025
	w2v-LDA	<b>0.213</b> $\pm$ 0.018	<b>0.442</b> $\pm$ 0.025	<b>0.502</b> $\pm$ 0.031	<b>0.509</b> $\pm$ 0.022
	glove-LDA	0.181 $\pm$ 0.011	0.420 $\pm$ 0.025	0.474 $\pm$ 0.029	0.498 $\pm$ 0.012
	Improve.	<b>0.009</b>	<b>0.05</b>	<b>0.043</b>	<b>0.032</b>



## Document classification (2)

Data	Method	$\lambda = 0.6$			
		T=7	T=20	T=40	T=80
TMN	DMM	0.607 $\pm$ 0.040	0.694 $\pm$ 0.026	0.712 $\pm$ 0.014	0.721 $\pm$ 0.008
	w2v-DMM	0.607 $\pm$ 0.019	0.736 $\pm$ 0.025	<b>0.760</b> $\pm$ 0.011	0.771 $\pm$ 0.005
	glove-DMM	<b>0.621</b> $\pm$ 0.042	<b>0.750</b> $\pm$ 0.011	0.759 $\pm$ 0.006	<b>0.775</b> $\pm$ 0.006
	Improve.	<b>0.014</b>	<b>0.056</b>	<b>0.048</b>	<b>0.054</b>
TMNtitle	DMM	0.500 $\pm$ 0.021	0.600 $\pm$ 0.015	0.630 $\pm$ 0.016	0.652 $\pm$ 0.005
	w2v-DMM	0.528 $\pm$ 0.028	0.663 $\pm$ 0.008	0.682 $\pm$ 0.008	<b>0.681</b> $\pm$ 0.006
	glove-DMM	<b>0.565</b> $\pm$ 0.022	<b>0.680</b> $\pm$ 0.011	<b>0.684</b> $\pm$ 0.009	<b>0.681</b> $\pm$ 0.004
	Improve.	<b>0.065</b>	<b>0.08</b>	<b>0.054</b>	<b>0.029</b>
Data	Method	$\lambda = 0.6$			
		T=4	T=20	T=40	T=80
Twitter	LDA	0.526 $\pm$ 0.021	0.636 $\pm$ 0.011	0.650 $\pm$ 0.014	0.653 $\pm$ 0.008
	w2v-LDA	<b>0.578</b> $\pm$ 0.047	0.651 $\pm$ 0.015	0.661 $\pm$ 0.011	<b>0.664</b> $\pm$ 0.010
	glove-LDA	0.569 $\pm$ 0.037	<b>0.656</b> $\pm$ 0.011	<b>0.662</b> $\pm$ 0.008	0.662 $\pm$ 0.006
	Improve.	<b>0.052</b>	<b>0.02</b>	<b>0.012</b>	<b>0.011</b>
Twitter	DMM	0.469 $\pm$ 0.014	0.600 $\pm$ 0.021	0.645 $\pm$ 0.009	0.665 $\pm$ 0.014
	w2v-DMM	<b>0.539</b> $\pm$ 0.016	0.649 $\pm$ 0.016	0.656 $\pm$ 0.007	0.676 $\pm$ 0.012
	glove-DMM	0.536 $\pm$ 0.027	<b>0.654</b> $\pm$ 0.019	<b>0.657</b> $\pm$ 0.008	<b>0.680</b> $\pm$ 0.009
	Improve.	<b>0.07</b>	<b>0.054</b>	<b>0.012</b>	<b>0.015</b>

## Conclusions and future directions

- Latent feature vectors induced from large external corpora can be used to improve topic modelling
  - ▶ latent features significantly improve topic coherence across a range of corpora with both the LDA and DMM models
  - ▶ document clustering and document classification also significantly improve, even though these depend directly only on the document-topic distribution
- The improvements were greatest for small document collections and/or for short documents
- We did not detect any reliable difference between word2vec and Glove vectors
- Retrain the word vectors to fit the topic-modeling corpus
- More sophisticated latent-feature models of topic-word distributions
- More efficient training procedures

Thank you for your attention!

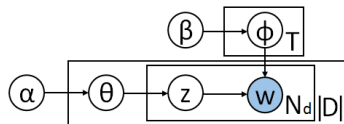
- Software:
  - ▶ <http://jldadmm.sourceforge.net>
  - ▶ <https://github.com/datquocnguyen/LFTM>

## Related work

- Phan et al. (2011) assumed that the small corpus is a sample of topics from a larger corpus like Wikipedia, and use the topics discovered in the larger corpus to help shape the topic representations in the small corpus
  - ▶ if the larger corpus has many irrelevant topics, this will “use up” the topic space of the model
- Petterson et al. (2010) proposed an extension of LDA that uses external information about word similarity, such as thesauri and dictionaries, to smooth the topic-to-word distribution
- Sahami and Heilman (2006) employed web search results to improve the information in short texts
- *Neural network topic models* of a single corpus have also been proposed (Salakhutdinov and Hinton, 2009; Srivastava et al., 2013; Cao et al., 2015).

# Latent Dirichlet Allocation (LDA)

$$\begin{aligned}\theta_d &\sim \text{Dir}(\alpha) & z_{d_i} &\sim \text{Cat}(\theta_d) \\ \phi_z &\sim \text{Dir}(\beta) & w_{d_i} &\sim \text{Cat}(\phi_{z_{d_i}})\end{aligned}$$

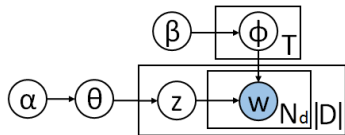


- Latent Dirichlet Allocation (LDA) is an *admixture model*, i.e., each document  $d$  is associated with a *distribution over topics*  $\theta_d$
- Inference is typically performed with a *Gibbs sampler* over the  $z_{d_i}$ , integrating out  $\theta$  and  $\phi$  (Griffiths et al., 2004)

$$P(z_{d_i}=t \mid \mathbf{Z}_{-d_i}) \propto (N_{d_{-i}}^t + \alpha) \frac{N_{-d_i}^{t, w_{d_i}} + \beta}{N_{-d_i}^t + V\beta}$$

# The Dirichlet Multinomial Mixture (DMM) model

$$\begin{aligned}\theta &\sim \text{Dir}(\alpha) & z_d &\sim \text{Cat}(\theta) \\ \phi_z &\sim \text{Dir}(\beta) & w_{d_i} &\sim \text{Cat}(\phi_{z_d})\end{aligned}$$



- The Dirichlet Multinomial Mixture (DMM) model is a *mixture model*, i.e., each document  $d$  is associated with a single topic  $z_d$  (Nigam et al., 2000)
- Inference can also be performed using a collapsed Gibbs sampler in which  $\theta$  and  $\phi_z$  are integrated out (Yin and Wang, 2014)

$$\begin{aligned}P(z_d = t \mid \mathbf{z}_{-d}) &\propto (M_{-d}^t + \alpha) \frac{\Gamma(N_{-d}^t + V\beta)}{\Gamma(N_{-d}^t + N_d + V\beta)} \\ &\prod_{w \in W} \frac{\Gamma(N_{-d}^{t,w} + N_d^w + \beta)}{\Gamma(N_{-d}^{t,w} + \beta)}\end{aligned}$$

## Estimating the topic vectors $\tau_t$

- Both the LF-LDA and LF-DMM associate each topic  $t$  with a *topic vector*  $\tau_t$ , which must be learnt from the training corpus
- After each Gibbs sweep:
  - ▶ the topic variables  $z$  identify which topic each word is generated from
  - ▶ the indicator variables  $s$  identify which words are generated from the latent feature distributions CatE
- We use LBFGS to optimize the L2-regularized log-loss (MAP estimation)

$$L_t = - \sum_{w \in W} K^{t,w} \left( \tau_t \cdot \omega_w - \log \left( \sum_{w' \in W} \exp(\tau_t \cdot \omega_{w'}) \right) \right) + \mu \| \tau_t \|_2^2$$

$$\text{NPMI-Score}(t) = \sum_{1 \leq i < j \leq N} \frac{\log \frac{P(w_i, w_j)}{P(w_i)P(w_j)}}{-\log P(w_i, w_j)}$$

- Sampling equations for inference in LF-LDA

$$\begin{aligned} &P(z_{d_i} = t \mid \mathbf{Z}_{\neg d_i}, \boldsymbol{\tau}, \boldsymbol{\omega}) \\ &\propto (N_{d_{\neg i}}^t + K_{d_{\neg i}}^t + \alpha) \\ &\quad \left( (1 - \lambda) \frac{N_{\neg d_i}^{t, w_{d_i}} + \beta}{N_{\neg d_i}^t + V\beta} + \lambda \text{CatE}(w_{d_i} \mid \boldsymbol{\tau}_t \boldsymbol{\omega}^\top) \right) \end{aligned} \quad (1)$$

$$P(s_{d_i} = s \mid z_{d_i} = t) \propto \begin{cases} (1 - \lambda) \frac{N_{\neg d_i}^{t, w_{d_i}} + \beta}{N_{\neg d_i}^t + V\beta} & \text{for } s = 0 \\ \lambda \text{CatE}(w_{d_i} \mid \boldsymbol{\tau}_t \boldsymbol{\omega}^\top) & \text{for } s = 1 \end{cases} \quad (2)$$



- Sampling equations for inference in LF-DMM

$$\begin{aligned}
 & P(z_d = t, s_d \mid \mathbf{Z}_{-d}, \mathbf{S}_{-d}, \boldsymbol{\tau}, \boldsymbol{\omega}) \\
 & \propto \lambda^{K_d} (1 - \lambda)^{N_d} (M_{-d}^t + \alpha) \frac{\Gamma(N_{-d}^t + V\beta)}{\Gamma(N_{-d}^t + N_d + V\beta)} \\
 & \quad \prod_{w \in W} \frac{\Gamma(N_{-d}^{t,w} + N_d^w + \beta)}{\Gamma(N_{-d}^{t,w} + \beta)} \prod_{w \in W} \text{CatE}(w \mid \boldsymbol{\tau}_t \boldsymbol{\omega}^\top)^{K_d^w}
 \end{aligned} \tag{3}$$

$$\begin{aligned}
 & Q(z_d = t, s_d \mid \mathbf{Z}_{-d}, \mathbf{S}_{-d}, \boldsymbol{\tau}, \boldsymbol{\omega}) \\
 & \propto \lambda^{K_d} (1 - \lambda)^{N_d} (M_{-d}^t + \alpha) \\
 & \quad \prod_{w \in W} \left( \frac{N_{-d}^{t,w} + \beta}{N_{-d}^t + V\beta} \right)^{N_d^w} \prod_{w \in W} \text{CatE}(w \mid \boldsymbol{\tau}_t \boldsymbol{\omega}^\top)^{K_d^w}
 \end{aligned} \tag{4}$$

$$\begin{aligned}
& Q(z_d = t \mid \mathbf{Z}_{-d}, \boldsymbol{\tau}, \boldsymbol{\omega}) \\
& \propto (M_{-d}^t + \alpha) \prod_{w \in W} \left( (1 - \lambda) \frac{N_{-d}^{t,w} + \beta}{N_{-d}^t + V\beta} + \lambda \text{CatE}(w \mid \boldsymbol{\tau}_t \boldsymbol{\omega}^\top) \right)^{(N_d^w + K_d^w)} \quad (5)
\end{aligned}$$

$$Q(s_{d_i} = s \mid z_d = t) \propto \begin{cases} (1 - \lambda) \frac{N_{-d}^{t, w_{d_i}} + \beta}{N_{-d}^t + V\beta} & \text{for } s = 0 \\ \lambda \text{CatE}(w_{d_i} \mid \boldsymbol{\tau}_t \boldsymbol{\omega}^\top) & \text{for } s = 1 \end{cases} \quad (6)$$