

DEEP UNFOLDING INFERENCE FOR SUPERVISED TOPIC MODEL

Chao-Hsi Lee and Jen-Tzung Chien

Department of Electrical and Computer Engineering
National Chiao Tung University, Taiwan, ROC

ABSTRACT

Conventional supervised topic model for multi-class classification is inferred via the variational inference algorithm where the model parameters are estimated by maximizing the lower bound of the logarithm of marginal likelihood function over input documents and labels. The classification accuracy is constrained by the variational lower bound. In this study, we aim to improve the classification accuracy by relaxing this constraint through directly maximizing the negative cross entropy error function via a deep unfolding inference (DUI). The inference procedure for class posterior is treated as the layer-wise learning in a deep neural network. The classification accuracy in DUI is accordingly increased by using the estimated topic parameters according to the exponentiated updates. Deep learning of supervised topic model is achieved through an error back-propagation algorithm. Experimental results show the superiority of DUI to variational Bayes inference in supervised topic model.

Index Terms— Deep unfolding, variational inference, supervised topic model

1. INTRODUCTION

Probabilistic topic model has been successfully developed for document categorization [1], image annotation [2], text segmentation [3], speech recognition [4, 5], document summarization [6], information retrieval [7], speech separation [8] and many other applications. Using topic model, the latent semantic topics are learned from a bag of words to capture the salient aspects embedded in data collection. The topic proportions is further used as the feature for classification [1]. To improve the classification accuracy, several methods [9, 10, 11, 12] have been studied to develop the supervised topic models to jointly capture the relationship of the documents and class labels.

Traditionally, the those topic models were inferred by the variational Bayes or the Gibbs sampling which are seen as the approximate inference procedures with different degrees of smoothing. It is because that the exact inference does not exist due to the coupling of multiple latent variables including topic assignments and topic proportions. Instead of direct maximization of marginal likelihood, the variational infer-

ence is developed by indirectly maximizing the lower bound of marginal likelihood through the variational expectation-maximization (EM) algorithm [13, 14].

However, the end performance is constrained by this lower bound. In general, the approximate inference may lead to errors at both estimation and prediction phases due to two reasons [15]. First, the expression of an underlying model may be deteriorated. Second, the approximation with parameter change may mislead the standard learning algorithm. In [16], the minimization of empirical risk was used instead of maximization of the approximate likelihood of training data. In [17], the errors due to inconsistent parameter estimator were shown to be partially compensated by using an approximate prediction method. In [18, 19, 20, 12], the discriminative training was demonstrated to outperform standard generative training in prediction tasks.

Recently, deep neural networks (DNNs) have been recognized as one of the most successful machine learning approaches to speech-related applications. This study aims to incorporate the power of DNN architecture into the inference of topic model and maximize the negative cross entropy error function directly. A deep unfolding inference is then developed accordingly. We derive the updating rules for parameter estimation which are implemented in a style of error back-propagation [16] similar to that in DNNs. The DNN-inferred topic model is then established. In the experiments, the classification performance is increased by using this deep unfolding inference compared to that using standard variational inference.

2. SUPERVISED TOPIC MODEL

This study investigates the DNN-inferred supervised topic model based on the supervised latent Dirichlet allocation (sLDA) [21] for multi-class classification (sLDAC) [10]. We first survey the model construction and inference based on sLDAC with graphical representation in Figure 1.

2.1. Model construction

sLDAC [10] is the extension of sLDA [21] which characterizes both documents and their labels. Given a set of Documents $\mathbf{w} = \{\mathbf{w}_d\} = \{w_{dn}\}$, the topic $z_{dn} = k$ of a

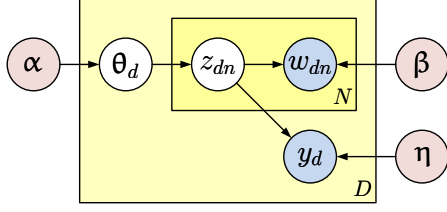


Fig. 1: Graphical representation for supervised topic model

word $w_{dn} = v$ is drawn by a multinomial distribution over V vocabulary words with parameters $\beta = \{\beta_k\} = \{\beta_{kv}\}$. sLDAC introduces a Dirichlet prior with hyper-parameters α for document-dependent topic $\theta_d = \{\theta_{dk}\} \sim \text{Dir}(\alpha)$ over K topics with hyper-parameters $\alpha = \{\alpha_k\}$ where $\alpha_k > 0$. Each document is treated as a “random mixture” over different topics. Each word w_{dn} and its topic assignment z_{dn} in a document d are multinomial variables expressed by $w_{dn} = v \sim \text{Mult}(\beta_{z_{dn}})$ and $z_{dn} = k \sim \text{Mult}(\theta_d)$, respectively. Accordingly, the latent variables in sLDAC consist of topic proportions $\theta = \{\theta_{dk}\}$ and topic assignments $\mathbf{z} = \{z_{dn}\}$. The class label outputs $\mathbf{y} = \{y_d\}$ with $y_d = m \sim \text{softmax}(\boldsymbol{\eta}^\top \bar{\mathbf{z}}_d)$ have the distributions

$$p(y_d = m | \mathbf{z}_d, \boldsymbol{\eta}) = \frac{\exp(\boldsymbol{\eta}_m^\top \bar{\mathbf{z}}_d)}{\sum_{i=1}^M \exp(\boldsymbol{\eta}_i^\top \bar{\mathbf{z}}_d)} \quad (1)$$

where $\bar{\mathbf{z}}_d = \{\bar{z}_{dk}\}$ denotes the empirical topic frequency, which satisfies $\bar{z}_{dk} = (1/N_d) \sum_{n=1}^{N_d} \delta(z_{dn}, k)$, $\boldsymbol{\eta} = \{\boldsymbol{\eta}_{mk}\}$ denotes the class label coefficients, and M is the number of class labels.

The model parameters $\{\alpha, \beta, \boldsymbol{\eta}\}$ are estimated by maximizing the marginal likelihood of $\{\mathbf{w}, \mathbf{y}\}$

$$p(\mathbf{w}, \mathbf{y} | \alpha, \beta, \boldsymbol{\eta}) = \prod_{d=1}^D p(y_d | \mathbf{z}_d, \boldsymbol{\eta}) \int_{\theta_d} p(\theta_d | \alpha) \times \prod_{n=1}^{N_d} \sum_{k=1}^K p(z_{dn} = k | \theta_d) p(w_{dn} | z_{dn} = k, \beta) d\theta_d \quad (2)$$

where N_d denotes the number of words in document d . However, the exact solution to model inference based on Eq. (2) does not exist due to the coupling of multiple latent variables θ and \mathbf{z} in posterior distribution $p(\theta, \mathbf{z} | \mathbf{w}, \mathbf{y}, \alpha, \beta, \boldsymbol{\eta})$.

2.2. Variational Bayes inference

Variational Bayes (VB) is implemented to resolve the intractable posterior distribution $p(\theta, \mathbf{z} | \mathbf{w}, \mathbf{y}, \alpha, \beta, \boldsymbol{\eta})$ [10] by using a factorizable variational distribution

$$q(\theta, \mathbf{z} | \boldsymbol{\gamma}, \boldsymbol{\phi}) = \prod_{d=1}^D q(\theta_d | \boldsymbol{\gamma}_d) \prod_{n=1}^{N_d} q(z_{dn} = k | \phi_{dnk}) \quad (3)$$

through maximizing a lower bound of the logarithm of marginal likelihood $\mathcal{L}(\boldsymbol{\gamma}, \boldsymbol{\phi}; \alpha, \beta, \boldsymbol{\eta})$ where $\boldsymbol{\gamma}$ and $\boldsymbol{\phi}$ denote the variational Dirichlet and multinomial parameters, respectively. We have the relation

$$\ln p(\mathbf{w}, \mathbf{y} | \alpha, \beta, \boldsymbol{\eta}) = \mathcal{L}(\boldsymbol{\gamma}, \boldsymbol{\phi}; \alpha, \beta, \boldsymbol{\eta}) + \text{KL}(q(\theta, \mathbf{z} | \boldsymbol{\gamma}, \boldsymbol{\phi}) \| p(\theta, \mathbf{z} | \mathbf{w}, \mathbf{y}, \alpha, \beta, \boldsymbol{\eta})). \quad (4)$$

Therefore, maximizing lower bound \mathcal{L} with respect to $\{\boldsymbol{\gamma}, \boldsymbol{\phi}\}$ is equivalent to estimating the new variational distribution q which is closest to the true posterior p with the smallest Kullback-Leibler divergence $\text{KL}(\cdot \| \cdot)$. The variational EM is solved by [10]. When the label is unknown, lower bound \mathcal{L} with variational parameters $\{\boldsymbol{\gamma}, \boldsymbol{\phi}\}$ degenerates into the expectation step of LDA. The class label output is estimated by

$$\hat{\mathbf{y}}_d = \mathbb{E}_q[y_d = m | \boldsymbol{\phi}_d, \mathbf{w}_d, \alpha, \beta, \boldsymbol{\eta}] = \frac{\boldsymbol{\eta}_m^\top \bar{\boldsymbol{\phi}}_d}{\sum_{i=1}^M \boldsymbol{\eta}_i^\top \bar{\boldsymbol{\phi}}_d} \quad (5)$$

where $\bar{\boldsymbol{\phi}}_d = (1/N_d) \sum_{n=1}^{N_d} \boldsymbol{\phi}_{dn}$ denotes the variational empirical topic frequency. We summarize the label prediction rule for sLDAC in Algorithm 1.

Algorithm 1 Prediction rule for supervised topic model

repeat

$\gamma_{dk} = \alpha_k + \sum_v N_{dv} \phi_{dvk}$
 $\phi_{dvk} \propto \exp\{\mathbb{E}[\ln \beta_{kv}] + \mathbb{E}[\ln \theta_{dk}]\}$

until γ_{dk} and ϕ_{dvk} converged

return $\hat{\mathbf{y}}_d = \text{softmax}(\boldsymbol{\eta}^\top \bar{\boldsymbol{\phi}}_d)$

3. DEEP UNFOLDING INFERENCE

We first address a general solution to deep unfolding inference (DUI) and then present how this new algorithm is realized to develop a deep inference for sLDAC.

3.1. General framework

Consider a topic model with parameters $\boldsymbol{\Theta}$. For each instance n , this model specifies an interesting output y_n for an observation x_n . In the inference procedure, we first estimate the data dependent parameters $\boldsymbol{\Psi}_n$ that maximize the objective function \mathcal{F} with respect to $\boldsymbol{\Theta}$. Once the optimal parameters $\hat{\boldsymbol{\Psi}}_n$ are estimated, we compute the interesting output \hat{y}_n through an estimator $\mathcal{G}_{\boldsymbol{\Theta}}$ as expressed by [20]

$$\hat{\boldsymbol{\Psi}}_n(x_n | \boldsymbol{\Theta}) = \arg \max_{\boldsymbol{\Psi}_n} \mathcal{F}_{\boldsymbol{\Theta}}(x_n, \boldsymbol{\Psi}_n) \quad (6)$$

$$\hat{y}_n(x_n | \boldsymbol{\Theta}) = \mathcal{G}_{\boldsymbol{\Theta}}(x_n, \hat{\boldsymbol{\Psi}}_n(x_n | \boldsymbol{\Theta})).$$

Typically, the optimal parameters $\hat{\boldsymbol{\Psi}}_n$ are obtained by running an iterative updating algorithm

$$\boldsymbol{\Psi}^{(l)} = f_{\boldsymbol{\Theta}}(x_n, \boldsymbol{\Psi}^{(l-1)}). \quad (7)$$

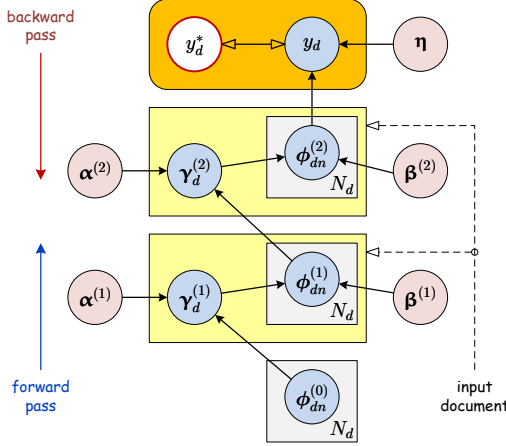


Fig. 2: Deep unfolding inference for supervised topic model

Here, the iteration index l can be treated as the layer in deep unfolding. The performance of a model can be measured by the evaluator \mathcal{J}_Θ with the optimal parameters Θ for each \hat{y}_n

$$\max_{\Theta} \mathcal{J}_\Theta(\{\hat{y}_n, \forall n\}). \quad (8)$$

Accordingly, maximizing the performance of a model is seen as a bi-level optimization problem. Maximizing \mathcal{J} in Eq. (8) with respect to Θ involves in maximizing \mathcal{F} with respect to parameter Ψ_n , and vice versa. Deep unfolding inference provides a way to this general framework where an error back-propagation algorithm is implemented to calculate the differentiations

$$\begin{aligned} \frac{\partial \mathcal{J}}{\partial \Theta^{(L)}} &= \sum_n \frac{\partial \mathcal{J}}{\partial \hat{y}_n} \frac{\partial \hat{y}_n}{\partial \Theta^{(L)}} & \frac{\partial \mathcal{J}}{\partial \Psi_n^{(L)}} &= \frac{\partial \mathcal{J}}{\partial \hat{y}_n} \frac{\partial \hat{y}_n}{\partial \Psi_n^{(L)}} \\ \frac{\partial \mathcal{J}}{\partial \Theta^{(l)}} &= \sum_n \frac{\partial \mathcal{J}}{\partial \Psi_n^{(l+1)}} \frac{\partial \Psi_n^{(l+1)}}{\partial \Theta^{(l)}} & \frac{\partial \mathcal{J}}{\partial \Psi_n^{(l)}} &= \frac{\partial \mathcal{J}}{\partial \Psi_n^{(l+1)}} \frac{\partial \Psi_n^{(l+1)}}{\partial \Psi_n^{(l)}}. \end{aligned} \quad (9)$$

3.2. DUI for supervised topic model

In this study, we would like to carry out deep unfolding inference for supervised latent Dirichlet allocation on classification task. Optimization problem is formulated as

$$\{\hat{\alpha}, \hat{\beta}, \hat{\eta}\} = \arg \max_{\{\alpha, \beta, \eta\}} \mathcal{E}(\mathbf{y}^*, \text{softmax}(\eta^\top \bar{\phi})) \quad (10)$$

$$\text{where } \{\hat{\gamma}, \hat{\phi}\} = \arg \max_{\{\phi, \gamma\}} \mathcal{L}(\gamma(\alpha, \phi), \phi(\beta, \gamma)) \quad (11)$$

where $\mathbf{y}^* = \{y_d^*\}$ denotes the truth labels and $\mathcal{E}(\cdot)$ is the negative cross entropy error function

$$\mathcal{E}(\mathbf{y}^*, \hat{\mathbf{y}}) = \sum_{d=1}^D \sum_{m=1}^M \delta(y_d^*, m) \ln \hat{y}_{dm}. \quad (12)$$

To maximize \mathcal{E} directly, we resort to DUI method. The model parameters Θ consist of $\mathbf{A} = \{\alpha^{(l)}\}$, $\mathbf{B} = \{\beta^{(l)}\}$ and η while the variational parameters Ψ are composed of $\{\gamma^{(l)}, \phi^{(l)}\}$ with L layers. $\mathcal{F} \triangleq \mathcal{L}$ and $\mathcal{J} \triangleq \mathcal{E}$. These parameters are estimated through unfolding the maximum step in Eq. (10) via

$$\{\hat{\mathbf{A}}, \hat{\mathbf{B}}, \hat{\eta}\} = \arg \max_{\{\mathbf{A}, \mathbf{B}, \eta\}} \mathcal{E}(\mathbf{y}^*, \text{softmax}(\eta^\top \bar{\phi}(\phi^{(L-1)}))) \quad (13)$$

with L layers

$$\phi_d^{(l)}(\beta^{(l)}, \gamma_d^{(l)}(\alpha^{(l)}, \phi_d^{(l-1)}(\dots \gamma_d^{(1)}(\alpha^{(1)}, \phi_d^{(0)}))) \dots \quad (14)$$

where $\phi_{dnk}^{(0)} = 1/K$. The nested structure is illustrated as shown in Figure 2, which can be realized by a procedure similar to Algorithm 1.

3.3. Error back-propagation procedure

The error back-propagation can be expressed in forward and backward pass. In forward pass, we calculate the variational parameters $\Psi^{(l)}$ for different layers l based on the model parameters $\Theta^{(l)}$ estimated in previous learning epochs. In backward pass, the model parameters are updated layer by layer by directly maximizing the negative cross entropy error function in Eq. (13) through a back-propagation algorithm. Similar to the error back-propagation in standard DNN training.

To deal with back-propagation of variational parameters, we define a set of auxiliary variables $\{a^{(l)}\}$ and re-express the forward pass as

$$\phi_{dvk}^{(l)} = \exp(a_{dvk}^{(l)}) / \sum_j \exp(a_{dvj}^{(l)}) \quad (15)$$

where $l = L - 1, \dots, 1$ and

$$a_{dvk}^{(L)} = \ln \beta_{vk}^{(L)} + \ln \gamma_{dk}^{(L)} - \ln \sum_j \gamma_{dj}^{(L)} \quad (16)$$

$$a_{dvk}^{(l)} = \ln \beta_{vk}^{(l)} + \psi(\gamma_{dk}^{(l)}). \quad (17)$$

The class label output is then estimated by

$$\bar{\phi}_{dk} = (1/N_d) \sum_v N_{dv} \phi_{dvk}^{(L-1)} \quad (18)$$

$$a_{dm}^{(L)} = \sum_k \eta_{mk} \bar{\phi}_{dk} \quad (19)$$

$$\hat{y}_{dm} = \exp(a_{dm}) / \sum_{i=1}^M \exp(a_{di}^{(L)}). \quad (20)$$

Next, the backward pass can be derived by applying chain rule and find the layer-wise differentiation of negative cross entropy. We summarize the procedure of back-propagation and updating rule in Algorithm 2, where $\psi'(\cdot)$ denotes the trigamma function. Since model parameters $\{\mathbf{A}, \mathbf{B}\}$ are non-negative, we adopt the *exponentiated gradient* method [22] to solve the constrained optimization problem subject to the non-negativity constraints. As a result, we implement the DNN-styled inference for topic model which is seen as an incorporation of the deterministic DNN into the inference of model-based method [20]. The powers of deep learning using DNN and stochastic learning using model-based method are assured in the proposed DUI.

Algorithm 2 Error back-propagation for DUI

/Back-propagation/

for all document d **do**

$$\frac{\partial \mathcal{E}}{\partial a_{dm}} = \hat{y}_{dm} - y_{dm}^*$$

$$\frac{\partial \mathcal{E}}{\partial \phi_{dk}} = \sum_m \frac{\partial \mathcal{E}}{\partial a_{dm}^{(L)}} \eta_{mk}$$

$$\frac{\partial \mathcal{E}}{\partial \phi_{dvk}^{(L-1)}} = \frac{N_{dv}}{N_d} \left(\frac{\partial \mathcal{E}}{\partial \phi_{dk}} - \frac{N_{dv}}{N_d} \sum_j \phi_{dvj}^{(L-1)} \frac{\partial \mathcal{E}}{\partial \phi_{dj}} \right)$$

for $l = L - 1$ **to** 1 **do**

$$\frac{\partial \mathcal{E}}{\partial a_{dvk}^{(l)}} = \phi_{dvk}^{(l)} \frac{\partial \mathcal{E}}{\partial \phi_{dvk}^{(l)}} - \phi_{dvk}^{(l)} \sum_j \phi_{dvj}^{(l)} \frac{\partial \mathcal{E}}{\partial \phi_{dj}^{(l)}}$$

$$\frac{\partial \mathcal{E}}{\partial \phi_{dvk}^{(l-1)}} = \frac{\partial \mathcal{E}}{\partial a_{dvk}^{(l)}} N_{dv} \psi'(\gamma_{dk}^{(l)})$$

end for

end for

/Updating rule/

Update η with $\frac{\partial \mathcal{E}}{\partial \eta_{km}} = \frac{\partial \mathcal{E}}{\partial a_{dm}^{(L)}} \bar{\phi}_{dk}$

for $l = 1$ **to** L **do**

Update $\beta^{(l)}$ with $\frac{\partial \mathcal{E}}{\partial \beta_{vk}^{(l)}} = \sum_d \frac{\partial \mathcal{E}}{\partial a_{dvk}^{(l)}} \frac{1}{\beta_{vk}^{(l)}}$

Update $\alpha^{(l)}$ with $\frac{\partial \mathcal{E}}{\partial \alpha_k^{(l)}} = \sum_{dv} \frac{\partial \mathcal{E}}{\partial \phi_{dvk}^{(l-1)}}$

end for

4. EXPERIMENTS

This section investigates the proposed DUI on document classification and compares the result with VB.

4.1. Experimental setup

For document classification, we evaluated the proposed method by using 20 newsgroups data set. It contains approximately 20,000 newsgroup documents, partitioned across 20 different newsgroups. In this task, we removed a list of English stop words provided by SMART information retrieval system, and the words that document frequencies were less than 5. Stemming was processed to reduce the inflected words into their word stem. Finally, we kept 5000 most frequent words in dictionary for document modeling. The original training set was randomly divided into a smaller training set and test set with 9000 documents and 6000 documents, respectively.

4.2. Experimental result

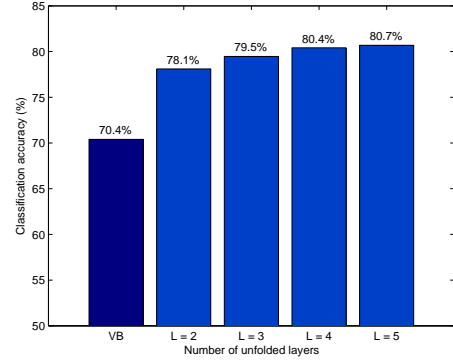
We investigated how the number of unfolded layers affects the accuracy of classification. We train a DUI model with $K = 60$, $\alpha = 0.1$, and vary $L = \{2, 3, 4, 5\}$. Mini-batch size is set to 1500. The model parameters \mathbf{A} and \mathbf{B} are random initialized and tied for all layers. η is initialized with a block diagonal matrix and trained with stochastic gradient descent

$$\eta = \{\eta_{mk}\} = \begin{bmatrix} 1 & & & \\ & \ddots & & \\ & & 1 & \\ & & & 1 \end{bmatrix} \quad (21)$$

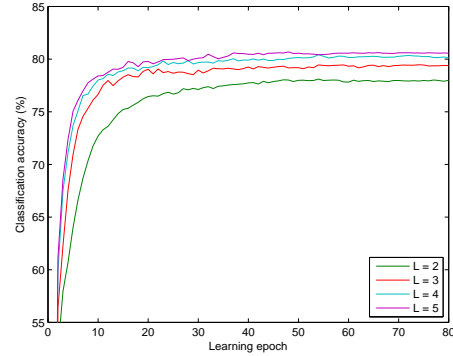
where $\mathbf{1} = (1, \dots, 1)$ in M/K dimensions, non-diagonal blocks are set to zero. Learning rate is scheduled as

$$\frac{\exp(\text{no. of epochs} - 1)}{200}. \quad (22)$$

Figure 3(a) shows the classification accuracy versus the number of layers using DUI. The result shows that DUI can achieve better classification accuracy around 80.7%, outperforming the result of VB, 70.4%. This implies that DUI can improve the accuracy significantly by this DNN training style. Figure 3(b) shows the convergence of classification accuracy. The deeper the DUI network, the better the accuracy and the convergence rate we can achieve.



(a) Classification accuracy vs. number of unfolded layers



(b) Convergence of classification accuracy

Fig. 3: Evaluation of classification accuracy using DUI

5. CONCLUSIONS

We proposed the DUI for sLDAC by unfolding the prediction rule to obtain a deep model. We derived the error back-propagation algorithm for learning parameters via DUI. Such learning algorithm is able to seek a deep model to improve classification performance. Experiments showed that maximizing the negative cross entropy error function instead of lower bound of marginal likelihood can outperform the conventional VB in terms of classification accuracy for test documents. Such DUI could be extended to the other models.

6. REFERENCES

- [1] D. M. Blei, A. Y. Ng, and M. I. Jordan, "Latent Dirichlet allocation," *Journal of Machine Learning Research*, vol. 3, pp. 993–1022, Mar. 2003.
- [2] D. Blei, L. Carin, and D. Dunson, "Probabilistic topic models," *IEEE Signal Processing Magazine*, vol. 27, no. 6, pp. 55–65, Nov. 2010.
- [3] J.-T. Chien and C.-H. Chueh, "Topic-based hierarchical segmentation," *IEEE Transactions on Audio, Speech and Language Processing*, vol. 20, no. 1, pp. 55–66, 2012.
- [4] J.-T. Chien and C.-H. Chueh, "Dirichlet class language models for speech recognition," *IEEE Transactions on Audio, Speech and Language Processing*, vol. 19, no. 3, pp. 482–495, 2011.
- [5] J.-T. Chien, "Hierarchical Pitman-Yor-Dirichlet language model," *IEEE/ACM Transactions on Audio, Speech and Language Processing*, vol. 23, no. 8, pp. 1259–1272, 2015.
- [6] Y.-L. Chang and J.-T. Chien, "Latent Dirichlet learning for document summarization," in *Proc. of International Conference on Acoustics, Speech, and Signal Processing*, 2009, pp. 1689–1692.
- [7] J.-T. Chien and M.-S. Wu, "Adaptive Bayesian latent semantic analysis," *IEEE Transactions on Audio, Speech and Language Processing*, vol. 16, no. 1, pp. 198–207, Jan 2008.
- [8] P. Smaragdis, B. Raj, and M. Shashanka, "A probabilistic latent variable model for acoustic modeling," *Advances in Models for Acoustic Processing Workshop (NIPS)*, 2006.
- [9] S. Lacoste-Julien, F. Sha, and M. I. Jordan, "DiscLDA: discriminative learning for dimensionality reduction and classification," in *Advances in Neural Information Processing Systems (NIPS)*, D. Koller, D. Schuurmans, Y. Bengio, and L. Bottou, Eds., 2008, pp. 897–904.
- [10] C. Wang, D. Blei, and F.-F. Li, "Simultaneous image classification and annotation," in *Proc. of IEEE Conference on Computer Vision and Pattern Recognition*, 2009, pp. 1903–1910.
- [11] J. Zhu, A. Ahmed, and E. P. Xing, "MedLDA: maximum margin supervised topic models," *Journal of Machine Learning Research*, vol. 13, no. 1, pp. 2237–2278, 2012.
- [12] J. Chen, J. He, Y. Shen, L. Xiao, X. He, J. Gao, X. Song, and L. Deng, "End-to-end learning of latent Dirichlet allocation by mirror-descent back propagation," *arXiv:1508.03398*, 2015.
- [13] A. P. Dempster, N. M. Laird, and D. B. Rubin, "Maximum likelihood from incomplete data via the EM algorithm," *Journal of the Royal Statistical Society. Series B (Methodological)*, vol. 39, no. 1, pp. 1–38, Jan. 1977.
- [14] M. J. Beal, *Variational Algorithms for Approximate Bayesian Inference*, Ph.D. thesis, University of London, 2003.
- [15] A. Kulesza and F. Pereira, "Structured learning with approximate inference," in *Advances in Neural Information Processing Systems (NIPS)*, 2007, pp. 785–792.
- [16] V. Stoyanov, A. Ropson, and J. Eisner, "Empirical risk minimization of graphical model parameters given approximate inference, decoding, and model structure," in *Proc. of International Conference on Artificial Intelligence and Statistics*, 2011, pp. 725–733.
- [17] M. J. Wainwright, "Estimating the "wrong" graphical model: benefits in the computation-limited setting," *Journal of Machine Learning Research*, vol. 7, pp. 1829–1859, 2006.
- [18] J. Lasserre, C. M. Bishop, and T. P. Minka, "Principled hybrids of generative and discriminative models," in *Proc. of IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR)*, 2006, vol. 1, pp. 87–94.
- [19] A. Holub and P. Perona, "A discriminative framework for modelling object classes," in *Proc. of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2005, vol. 1, pp. 664–671.
- [20] J. R. Hershey, J. Le Roux, and F. Weninger, "Deep unfolding: model-based inspiration of novel deep architectures," *arXiv:1409.2574*, 2014.
- [21] J. D. Mcauliffe and D. M. Blei, "Supervised topic models," in *Advances in Neural Information Processing Systems (NIPS)*, 2007, pp. 121–128.
- [22] J. Kivinen and M. K. Warmuth, "Exponentiated gradient versus gradient descent for linear predictors," *Information and Computation*, vol. 132, no. 1, pp. 1–63, 1997.