

# Supplementary Materials: Modeling Heterogeneous Statistical Patterns In High-dimensional Data By Adversarial Distributions: An Unsupervised Generative Framework

Han Zhang<sup>1</sup>, Wenhao Zheng<sup>3</sup>, Charley Chen<sup>1</sup>, Kevin Gao<sup>1</sup>, Yao Hu<sup>3</sup>, Ling Huang<sup>2</sup> and Wei Xu<sup>1</sup>

<sup>1</sup>Tsinghua University, Beijing, China

<sup>2</sup>AHI Fintech, China, <sup>3</sup>Youku Cognitive and Intelligent Lab, Alibaba Group

## 1 Derivation of EM Updates

In E step we approximate the likelihood by finding the best variational distribution, which is given by

$$p(d_n, \mathbf{f}_n | \mathbf{x}_n, \hat{\boldsymbol{\theta}}) = \frac{\prod_{g=1}^G \left\{ \hat{\pi}_g \prod_{m=1}^M \gamma_{ngm}^{f_{nm}} \cdot \bar{\gamma}_{ngm}^{1-f_{nm}} \right\}^{d_{ng}}}{\sum_{g'=1}^G \phi_{ng'}}.$$

The information of this distribution is summarized in Eq. (6). In M step we optimize the log-likelihood with respect to the model parameters:

$$\begin{aligned} \max_{\boldsymbol{\theta}} Q(\boldsymbol{\theta}, \hat{\boldsymbol{\theta}}) &= \sum_{n=1}^N \sum_{d_n, \mathbf{f}_n} p(d_n, \mathbf{f}_n | \mathbf{x}_n, \hat{\boldsymbol{\theta}}) \log p(\mathbf{x}_n, \mathbf{f}_n, d_n | \boldsymbol{\theta}) \\ &\quad - \sum_{g=1}^G \lambda_g^{(1)} \log \pi_g - \sum_{gmi} \lambda_{gmi}^{(2)} (\log \alpha_{gmi} - \log \beta_{gmi}), \end{aligned}$$

where the parameters lie in probability simplexes. For  $\boldsymbol{\mu}_g$  we set

$$\begin{aligned} \frac{\partial Q}{\partial \mu_{gm}} &= \sum_{n=1}^N \sum_{d_n, \mathbf{f}_n} p(d_n, \mathbf{f}_n | \mathbf{x}_n, \hat{\boldsymbol{\theta}}) \frac{\partial \log p(\mathbf{f}_n | d_n, \boldsymbol{\mu})}{\partial \mu_{gm}} \\ &= \sum_{n=1}^N \frac{\hat{\pi}_g \prod_{m' \neq m} (\gamma_{ngm'} + \bar{\gamma}_{ngm'})}{\sum_{g'=1}^G \phi_{ng'}} \left( \frac{\gamma_{ngm}}{\mu_{gm}} - \frac{\bar{\gamma}_{ngm}}{1 - \mu_{gm}} \right) \end{aligned}$$

to zero and we obtain the update equation in Eq. (7). For  $\beta_{gmi}$  we need to introduce Lagrange multiplier  $\omega$  to eliminate the probability simplex constraint:

$$\begin{aligned} \frac{\partial Q}{\partial \beta_{gmi}} &= \sum_{n=1}^N \sum_{d_n, \mathbf{f}_n} p(d_n, \mathbf{f}_n | \mathbf{x}_n, \hat{\boldsymbol{\theta}}) \frac{\partial \log p(x_{nm} | f_{nm}, d_n, \boldsymbol{\alpha}, \boldsymbol{\beta})}{\partial \beta_{gmi}} + \lambda_{gmi}^{(2)} - \omega \\ &= \sum_{n=1}^N \sum_{n=1}^N \frac{\hat{\pi}_g \prod_{m' \neq m} (\gamma_{ngm'} + \bar{\gamma}_{ngm'})}{\sum_{g'=1}^G \phi_{ng'}} \cdot \frac{x_{nmi}}{\beta_{gmi}} + \lambda_{gmi}^{(2)} - \omega. \end{aligned}$$

Setting  $\partial Q / \partial \beta_{gmi}$  to zero and using  $\sum_i \beta_{gmi} = 1$ ,  $\sum_i x_{nmi} = 1$  we obtain the result in Eq. (7).

Optimization w.r.t.  $\alpha_{gmi}$  and  $\pi_g$  is difficult since the objective function is no longer convex. [1] proposed a technique to optimize such concave problems with convex regularizers in a probability simplex. Given a variable  $\boldsymbol{\theta}$  in a  $(K+1)$ -simplex we can introduce a pseudo-Dirichlet prior to promote sparsity

$$Q(\boldsymbol{\theta}; \lambda, \epsilon) = \sum_{k=1}^K c_k \log \theta_k - \sum_{k=1}^K \lambda \log(\theta_k + \epsilon) + \text{constant},$$

where  $c_k$  is the observation counts for  $k$ -th possible value and  $\lambda > 0$ . [1] proved that  $Q(\boldsymbol{\theta}; \lambda, \epsilon)$  has one global maximum, and we can use the fixed point iteration to calculate this maximum. Suppose  $\sum_k c_k = c$ , the update rule is

$$\theta_k = \frac{c_k + \theta_k \cdot \lambda \cdot \sum_{k=1}^K \theta_k / (\theta_k + \epsilon)}{c + \lambda / (\theta_k + \epsilon)}.$$

Using the fact that the sum of observation counts for  $\boldsymbol{\pi}$  is  $N$  and  $\sum_{n=1}^N \tilde{\gamma}_{ngm} \tilde{\phi}_{ng}$  for  $\boldsymbol{\alpha}_{gm}$  and letting  $\epsilon \rightarrow 0$  we obtain the update rules in Eq. (8).

## References

- [1] Martin O. Larsson and Johan Ugander. A concave regularization technique for sparse mixture models. In *Advances in Neural Information Processing Systems 24: 25th Annual Conference on Neural Information Processing Systems 2011. Proceedings of a meeting held 12-14 December 2011, Granada, Spain.*, pages 1890–1898, 2011.