

Scalable Deep Poisson Factor Analysis for Topic Modeling

Zhe Gan, Changyou Chen,
Ricardo Henao, David Carlson, Lawrence Carin

Duke University

July 9th, 2015

Outline

- 1 Introduction
- 2 Model Formulation
- 3 Scalable Posterior Inference
- 4 Experiments
- 5 Summary

Problem of interest: How to develop **deep generative** models for documents that are represented in **bag-of-words** form?

- **Directed Graphical Models:**

- Latent Dirichlet Allocation (**LDA**) (Blei et al., 2003)
- Focused Topic Model (**FTM**) (Williamson et al., 2010)
- Poisson Factor Analysis (**PFA**) (Zhou et al., 2012)

- **Going “Deep”?**

- Hierarchical **tree-structured** topic models
- nested Chinese Restaurant Process (**nCRP**) (Blei et al., 2004)
- Hierarchical Dirichlet Process (**HDP**) (Teh et al., 2006)
- nested Hierarchical Dirichlet Process (**nHDP**) (Paisley et al., 2015)

- How about we want to model **general topic correlations**?

- **Undirected Graphical Models:**
 - Replicated Softmax Model (**RSM**) (Salakhutdinov and Hinton, 2009b)
 - One generalization of the Restricted Boltzmann Machine (**RBM**) (Hinton, 2002)
- **Going Deep?**
 - Deep Belief Networks (**DBN**) (Hinton et al., 2006; Hinton and Salakhutdinov, 2011)
 - Deep Boltzmann Machines (**DBM**) (Salakhutdinov and Hinton, 2009a; Srivastava et al., 2013)
- Topics are not defined “properly”.

Introduction

Main idea:

- Poisson Factor Analysis (PFA) + Deep Sigmoid Belief Network (SBN) or Restricted Boltzmann Machine (RBM).
- PFA is employed to interact with data at the bottom layer.
- Deep SBN or RBM serve as a flexible prior for revealing topic structure.

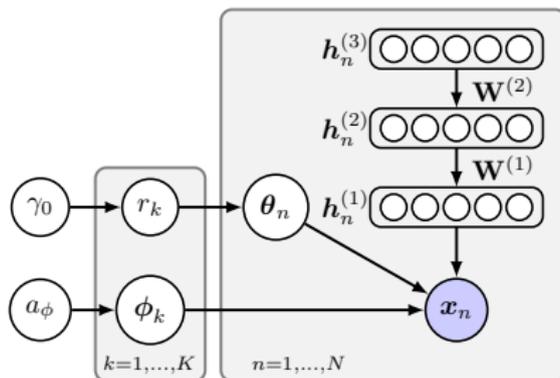


Figure: Graphical model for the Deep Poisson Factor Analysis with three layers of hidden binary hierarchies. The directed binary hierarchy may be replaced by a *deep Boltzmann machine*.

Poisson Factor Analysis: (Zhou et al., 2012)

- We represent a discrete matrix $\mathbf{X} \in \mathbb{Z}_+^{P \times N}$ containing counts from N documents and P words as

$$\mathbf{X} = \text{Pois}(\Phi(\Theta \circ \mathbf{H}^{(1)})). \quad (1)$$

- Each column of Φ , ϕ_k , encodes the relative **importance of each word** in topic k .
- Each column of Θ , θ_n , contains relative **topic intensities** specific to document n .
- Each column of $\mathbf{H}^{(1)}$, $\mathbf{h}_n^{(1)}$, defines a **sparse** set of topics associated with each document.

Poisson Factor Analysis: (Zhou et al., 2012)

- We construct PFAs by placing **Dirichlet** priors on ϕ_k and **gamma** priors on θ_n .

$$x_{pn} = \sum_{k=1}^K x_{pnk}, \quad x_{pnk} \sim \text{Pois}(\phi_{pk} \theta_{kn} h_{kn}^{(1)}), \quad (2)$$

with priors specified as $\phi_k \sim \text{Dir}(\mathbf{a}_\phi, \dots, \mathbf{a}_\phi)$,
 $\theta_{kn} \sim \text{Gamma}(r_k, p_n / (1 - p_n))$, $r_k \sim \text{Gamma}(\gamma_0, 1/c_0)$, and
 $\gamma_0 \sim \text{Gamma}(e_0, 1/f_0)$.

- Previously, a **beta-Bernoulli process** prior is defined on $\mathbf{h}_n^{(1)}$, assuming **topic independence** (Zhou and Carin, 2015).
- The **novelty** in our models comes from the prior for $\mathbf{h}_n^{(1)}$.

Structured Priors on the Latent Binary matrix:

- Assume $\mathbf{h}_n^{(1)} \in \{0, 1\}^{K_1}$, we define another hidden set of units $\mathbf{h}_n^{(2)} \in \{0, 1\}^{K_2}$ placed at a layer “above” $\mathbf{h}_n^{(1)}$.
- **Modeling with the RBM: (Undirected)**

$$- E(\mathbf{h}_n^{(1)}, \mathbf{h}_n^{(2)}) = (\mathbf{h}_n^{(1)})^\top \mathbf{c}^{(1)} + (\mathbf{h}_n^{(1)})^\top \mathbf{W}^{(1)} \mathbf{h}_n^{(2)} + (\mathbf{h}_n^{(2)})^\top \mathbf{c}^{(2)}. \quad (3)$$

- **Modeling with the SBN (Neal, 1992): (Directed)**

$$p(h_{k_2 n}^{(2)} = 1) = \sigma(c_{k_2}^{(2)}), \quad (4)$$

$$p(h_{k_1 n}^{(1)} = 1 | \mathbf{h}_n^{(2)}) = \sigma \left((\mathbf{w}_{k_1}^{(1)})^\top \mathbf{h}_n^{(2)} + c_{k_1}^{(1)} \right). \quad (5)$$

Going Deep?

- Add multiple layers of SBNs or RBMs.

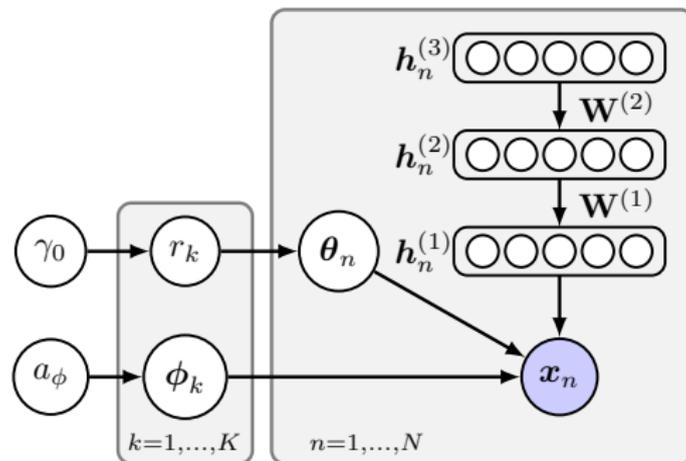


Figure: Graphical model for the Deep Poisson Factor Analysis with three layers of hidden binary hierarchies. The directed binary hierarchy may be replaced by a *deep Boltzmann machine*.

Challenge: Designing **scalable Bayesian inference** algorithms.

Solutions: Scaling up inference by **stochastic** algorithms.

- Applying **Bayesian conditional density filtering** algorithm (Guhaniyogi et al., 2014).
- Extending recently proposed work on **stochastic gradient thermostats** (Ding et al., 2014).

Bayesian conditional density filtering (BCDF):

- Repeatedly updating the surrogate conditional sufficient statistics (SCSS) using the current mini-batch.
- Drawing samples from the conditional posterior distributions of model parameters, based on SCSS.
- “stochastic Gibbs-style” updates.

Input: text documents, *i.e.*, a count matrix \mathbf{X} .
Initialize $\Psi_g^{(0)}$ randomly and set $\mathbf{s}_g^{(0)}$ all to zero.
for $t = 1$ **to** ∞ **do**
 Get one mini-batch $\mathbf{X}^{(t)}$.
 Initialize $\Psi_g^{(t)} = \Psi_g^{(t-1)}$, and $\mathbf{s}_g^{(t)} = \mathbf{s}_g^{(t-1)}$.
 Initialize $\Psi_j^{(t)}$ randomly.
 for $s = 1$ **to** S **do**
 Gibbs sampling for DPFA on $\mathbf{X}^{(t)}$.
 Collect samples $\Psi_g^{1:S}, \Psi_j^{1:S}$ and $\mathbf{s}_g^{1:S}$.
 end for
 Set $\Psi_g^{(t)} = \text{mean}(\Psi_g^{1:S})$, and $\mathbf{s}_g^{(t)} = \text{mean}(\mathbf{s}_g^{1:S})$.
end for

- Ψ_g : global parameters
- Ψ_j : local hidden variables
- \mathbf{s}_g : SCSS for Ψ_g

Stochastic Gradient Nose-Hoover Thermostats (SGNHT):

- Extending *Hamiltonian Monte Carlo* using **stochastic gradient**.
- Introducing **thermostat** to maintain system temperature.
- Adaptively **absorbing** stochastic gradient noise.
- The motion of the particles in the system are defined by the stochastic differential equations (**SDE**)

$$\begin{aligned}d\boldsymbol{\Psi}_g &= \mathbf{v}dt, & d\mathbf{v} &= \tilde{f}(\boldsymbol{\Psi}_g)dt - \xi \mathbf{v}dt + \sqrt{D}d\mathcal{W}, \\d\xi &= \left(\frac{1}{M} \mathbf{v}^T \mathbf{v} - 1 \right) dt,\end{aligned}\tag{6}$$

where $\boldsymbol{\Psi}_g \in \mathbb{R}^M$ are model parameters, $\mathbf{v} \in \mathbb{R}^M$ are the momentum variables, $\tilde{f}(\boldsymbol{\Psi}_g) \triangleq -\nabla_{\boldsymbol{\Psi}_g} \tilde{U}(\boldsymbol{\Psi}_g)$, and $\tilde{U}(\boldsymbol{\Psi}_g)$ is the negative log-posterior.

Scalable Posterior Inference

Extension:

- Extending the SGNHT by introducing **multiple thermostat variables** (ξ_1, \dots, ξ_M) into the system such that each ξ_i controls one degree of the particle momentum.
- The proposed SGNHT is defined by the following SDEs

$$\begin{aligned} d\boldsymbol{\Psi}_g &= \mathbf{v}dt, & d\mathbf{v} &= \tilde{f}(\boldsymbol{\Psi}_g)dt - \boldsymbol{\Xi}\mathbf{v}dt + \sqrt{D}d\mathcal{W}, \\ d\boldsymbol{\Xi} &= (\mathbf{q} - \mathbf{I})dt, \end{aligned} \tag{7}$$

where $\boldsymbol{\Xi} = \text{diag}(\xi_1, \xi_2, \dots, \xi_M)$, $\mathbf{q} = \text{diag}(v_1^2, \dots, v_M^2)$

Theorem

The equilibrium distribution of the SDE system in (7) is

$$p(\boldsymbol{\Psi}_g, \mathbf{v}, \boldsymbol{\Xi}) \propto \exp\left(-\frac{1}{2}\mathbf{v}^\top \mathbf{v} - U(\boldsymbol{\Psi}_g) - \frac{1}{2}\text{tr}\left\{(\boldsymbol{\Xi} - D)^\top (\boldsymbol{\Xi} - D)\right\}\right).$$

Stochastic Gradient Noise-Hoover Thermostats (SGNHT):

Input: text documents, *i.e.*, a count matrix \mathbf{X} .

Random Initialization.

for $t = 1$ **to** ∞ **do**

$$\boldsymbol{\Psi}_g^{(t+1)} = \boldsymbol{\Psi}_g^{(t)} + \mathbf{v}^{(t)} h.$$

$$\mathbf{v}^{(t+1)} = \tilde{f}(\boldsymbol{\Psi}_g^{(t+1)})h - \Xi^{(t)} \mathbf{v}^{(t)} h + \sqrt{2Dh} \mathcal{N}(0, \mathbf{I}).$$

$$\Xi^{(t+1)} = \Xi^{(t)} + (\mathbf{q}^{(t+1)} - \mathbf{I})h, \text{ where } \mathbf{q} = \text{diag}(v_1^2, \dots, v_M^2).$$

end for

Stochastic Gradient Noise-Hoover Thermostats (SGNHT):

Input: text documents, *i.e.*, a count matrix \mathbf{X} .

Random Initialization.

for $t = 1$ **to** ∞ **do**

$$\boldsymbol{\Psi}_g^{(t+1)} = \boldsymbol{\Psi}_g^{(t)} + \mathbf{v}^{(t)} h.$$

$$\mathbf{v}^{(t+1)} = \tilde{f}(\boldsymbol{\Psi}_g^{(t+1)})h - \boldsymbol{\Xi}^{(t)}\mathbf{v}^{(t)}h + \sqrt{2Dh}\mathcal{N}(0, \mathbf{I}).$$

$$\boldsymbol{\Xi}^{(t+1)} = \boldsymbol{\Xi}^{(t)} + (\mathbf{q}^{(t+1)} - \mathbf{I})h, \text{ where } \mathbf{q} = \text{diag}(v_1^2, \dots, v_M^2).$$

end for

Discussion:

- **BCDF:** ease of implementation, but prefers the conditional densities for all the parameters.
- **SGNHT:** more general and robust, fast convergence.

Datasets:

- **20 Newsgroups:** 20K documents with a vocabulary size of 2K.
- **RCV1-v2:** 800K documents with a vocabulary size of 10K.
- **Wikipedia:** 10M documents with a vocabulary size of 8K.

Quantitative Evaluation:

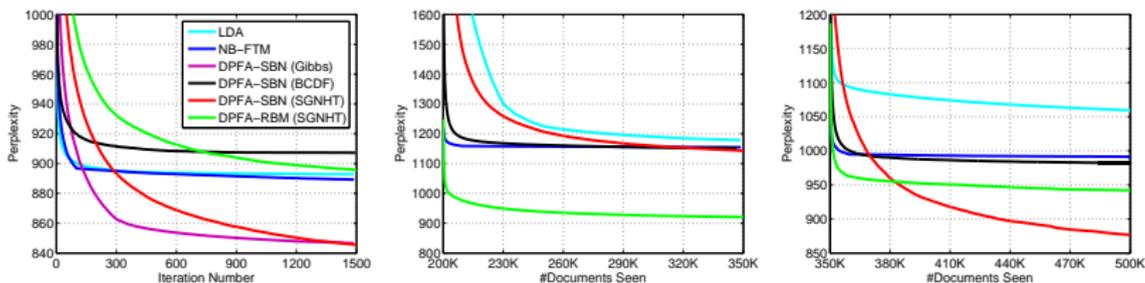
Table: 20 Newsgroups.

MODEL	METHOD	DIM	PERP.
DPFA-SBN- <i>t</i>	GIBBS	128-64-32	827
DPFA-SBN	GIBBS	128-64-32	846
DPFA-SBN	SGNHT	128-64-32	846
DPFA-RBM	SGNHT	128-64-32	896
DPFA-SBN	BCDF	128-64-32	905
DPFA-SBN	GIBBS	128-64	851
DPFA-SBN	SGNHT	128-64	850
DPFA-RBM	SGNHT	128-64	893
DPFA-SBN	BCDF	128-64	896
LDA	GIBBS	128	893
NB-FTM	GIBBS	128	887
RSM	CD5	128	877
NHDP	SVB	(10,10,5) [◊]	889

Table: RCV1-v2 & Wikipedia.

MODEL	METHOD	DIM	RCV	Wiki
DPFA-SBN	SGNHT	1024-512-256	964	770
DPFA-SBN	SGNHT	512-256-128	1073	799
DPFA-SBN	SGNHT	128-64-32	1143	876
DPFA-RBM	SGNHT	128-64-32	920	942
DPFA-SBN	BCDF	128-64-32	1149	986
LDA	BCDF	128	1179	1059
NB-FTM	BCDF	128	1155	991
RSM	CD5	128	1171	1001
NHDP	SVB	(10,5,5) [◊]	1041	932

Quantitative Evaluation:



Sensitivity Analysis:

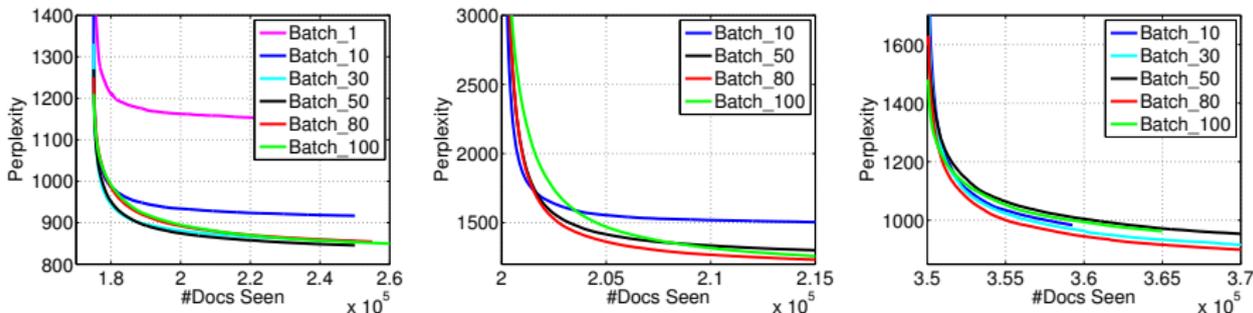


Figure: Perplexities. (Left) *20 News*. (Middle) *RCV1-v2*. (Right) *Wikipedia*.

Experiments

Topics we learned on 20 Newsgroups:

T1	T3	T8	T9	T10	T14	T15	T19	T21	T24
year hit runs good season	people real simply world things	group groups reading newsgroup pro	world country countries germany nazi	evidence claim people argument agree	game games win cup hockey	israel israeli jews arab jewish	software modem port mac serial	files file ftp program format	team players player play teams
T25	T26	T29	T40	T41	T43	T50	T54	T55	T64
god existence exist human atheism	fire fbi koresh children batf	people life death kill killing	wrong doesn jim agree quote	image program application widget color	boston toronto montreal chicago pittsburgh	problem work problems system fine	card video memory mhz bit	windows dos file win ms	turkish armenian armenians turks armenia
T65	T69	T78	T81	T91	T94	T112	T118	T120	T126
truth true point fact body	window server display manager client	drive disk scsi hard drives	makes power make doesn part	question answer means true people	code mit comp unix source	children father child mother son	people make person things feel	men women man hand world	sex sexual cramer gay homosexual

Experiments

Visualization:

Sports, Computers, and Politics/Law.

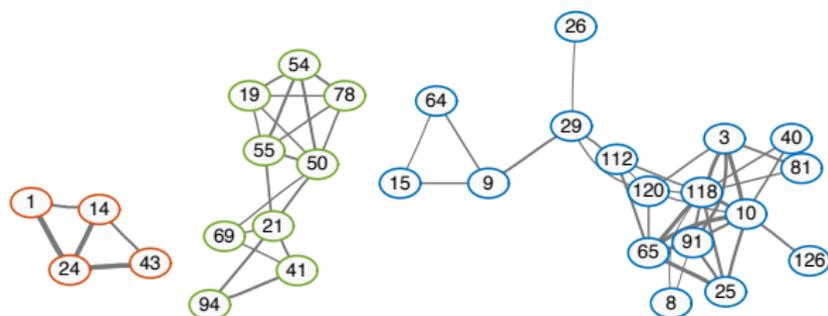


Figure: Graphs induced by the correlation structure learned by DPFA-SBN for the *20 Newsgroups*.

- **Model:** Deep Poisson Factor Analysis
 - PFA is employed to interact with data at the bottom layer.
 - Deep SBN or RBM serve as a flexible prior for revealing topic structure.
- **Scalable Inference:**
 - Bayesian conditional density filtering.
 - Stochastic gradient thermostats.



https://github.com/zhegan27/dpfa_icml2015

Questions?

References I

- Blei, D. M., Griffiths, T., Jordan, M. I., and Tenenbaum, J. B. (2004). Hierarchical topic models and the nested Chinese restaurant process. *NIPS*.
- Blei, D. M., Ng, A. Y., and Jordan, M. I. (2003). Latent Dirichlet allocation. *JMLR*.
- Ding, N., Fang, Y., Babbush, R., Chen, C., Skeel, R. D., and Neven, H. (2014). Bayesian sampling using stochastic gradient thermostats. *NIPS*.
- Guhaniyogi, R., Qamar, S., and Dunson, D. B. (2014). Bayesian conditional density filtering. *arXiv:1401.3632*.
- Hinton, G. E. (2002). Training products of experts by minimizing contrastive divergence. *Neural computation*.
- Hinton, G. E., Osindero, S., and Teh, Y. W. (2006). A fast learning algorithm for deep belief nets. *Neural computation*.
- Hinton, G. E. and Salakhutdinov, R. (2011). Discovering binary codes for documents by learning deep generative models. *Topics in Cognitive Science*.
- Neal, R. M. (1992). Connectionist learning of belief networks. *Artificial Intelligence*.
- Paisley, J., Wang, C., Blei, D. M., and Jordan, M. I. (2015). Nested hierarchical Dirichlet processes. *PAMI*.

References II

- Salakhutdinov, R. and Hinton, G. E. (2009a). Deep Boltzmann machines. *AISTATS*.
- Salakhutdinov, R. and Hinton, G. E. (2009b). Replicated softmax: an undirected topic model. *NIPS*.
- Srivastava, N., Salakhutdinov, R., and Hinton, G. E. (2013). Modeling documents with deep Boltzmann machines. *UAI*.
- Teh, Y. W., Jordan, M. I., Beal, M. J., and Blei, D. M. (2006). Hierarchical Dirichlet processes. *JASA*.
- Williamson, S., Wang, C., Heller, K., and Blei, D. M. (2010). The IBP compound Dirichlet process and its application to focused topic modeling. *ICML*.
- Zhou, M. and Carin, L. (2015). Negative binomial process count and mixture modeling. *PAMI*.
- Zhou, M., Hannah, L., Dunson, D., and Carin, L. (2012). Beta-negative binomial process and Poisson factor analysis. *AISTATS*.