# The Poisson Gamma Belief Network

**Mingyuan Zhou[#], Yulai Cong[*], and Bo Chen[*]**

**[#] IROM Department, The University of Texas at Austin, Austin, TX, USA**
**[*] School of Electronic Engineering, Xidian University, Xi'an, Shaanxi, China**

THE UNIVERSITY OF TEXAS AT AUSTIN
McCOMBS SCHOOL OF BUSINESS

西安电子科技大学
XIDIAN UNIVERSITY

## Introduction

The Poisson gamma belief network (PGBN) is proposed to infer a multilayer representation of high-dimensional count vectors.

➤ The PGBN factorizes each of its layers into the product of a connection weight matrix and the nonnegative real hidden units of the next layer.

➤ The PGBN's hidden layers are jointly trained with an upward-downward Gibbs sampler.

➤ The gamma-negative binomial process combined with a layer-wise training strategy allows the PGBN to infer the width of each layer given a fixed budget on the width of the first layer.

➤ Example results illustrate interesting relationships between the width of the first layer and the inferred network structure, and demonstrate that the PGBN can add more layers to increase its performance gains over Poisson factor analysis.

## Hierarchical Model and Properties

### ❑ Poisson Gamma Belief Network (PGBN)

$$\boldsymbol{\theta}_j^{(T)} \sim \text{Gam}\left(\boldsymbol{r}, 1/c_j^{(T+1)}\right),$$
$$\cdots$$
$$\boldsymbol{\theta}_j^{(t)} \sim \text{Gam}\left(\boldsymbol{\Phi}^{(t+1)}\boldsymbol{\theta}_j^{(t+1)}, 1/c_j^{(t+1)}\right),$$
$$\cdots$$
$$\boldsymbol{x}_j^{(1)} \sim \text{Pois}\left(\boldsymbol{\Phi}^{(1)}\boldsymbol{\theta}_j^{(1)}\right), \quad \boldsymbol{\theta}_j^{(1)} \sim \text{Gam}\left(\boldsymbol{\Phi}^{(2)}\boldsymbol{\theta}_j^{(2)}, p_j^{(2)}/\left(1-p_j^{(2)}\right)\right).$$

### ❑ PGBN vs Sigmoid belief network (SBN)

$$P\left(\boldsymbol{x}_j^{(1)}, \{\boldsymbol{\theta}_j^{(t)}\}_t \,\big|\, \{\boldsymbol{\Phi}^{(t)}\}_t\right) = P\left(\boldsymbol{x}_j^{(1)} \,\big|\, \boldsymbol{\Phi}^{(1)}, \boldsymbol{\theta}_j^{(1)}\right)\left[\prod_{t=1}^{T-1} P\left(\boldsymbol{\theta}_j^{(t)} \,\big|\, \boldsymbol{\Phi}^{(t+1)}, \boldsymbol{\theta}_j^{(t+1)}\right)\right] P\left(\boldsymbol{\theta}_j^{(T)}\right)$$

**PGBN:** $P\left(\theta_{vj}^{(t)} \,\big|\, \phi_{v:}^{(t+1)}, \boldsymbol{\theta}_j^{(t+1)}, c_{j+1}^{(t+1)}\right) = \frac{\left(c_{j+1}^{(t+1)}\right)^{\phi_{v:}^{(t+1)}\boldsymbol{\theta}_j^{(t+1)}}}{\Gamma\left(\phi_{v:}^{(t+1)}\boldsymbol{\theta}_j^{(t+1)}\right)}\left(\theta_{vj}^{(t)}\right)^{\phi_{v:}^{(t+1)}\boldsymbol{\theta}_j^{(t+1)}-1} e^{-c_{j+1}^{(t+1)}\theta_{vj}^{(t)}}$

**SBN:** $P\left(\theta_{vj}^{(t)}=1 \,\big|\, \phi_{v:}^{(t+1)}, \boldsymbol{\theta}_j^{(t+1)}, b_v^{(t+1)}\right) = \sigma\left(b_v^{(t+1)} + \phi_{v:}^{(t+1)}\boldsymbol{\theta}_j^{(t+1)}\right)$

### ❑ Properties of PGBN:

**Lemma 1** (Augment-and-conquer the PGBN). *With* $p_j^{(1)} := 1 - e^{-1}$ *and*
$$p_j^{(t+1)} := -\ln(1-p_j^{(t)}) / \left[c_j^{(t+1)} - \ln(1-p_j^{(t)})\right]$$

*for* $t = 1, \ldots, T$, *one may connect the observed (if $t = 1$) or some latent (if $t \geq 2$) counts* $\boldsymbol{x}_j^{(t)} \in \mathbb{Z}^{K_{t-1}}$ *to the product* $\boldsymbol{\Phi}^{(t)}\boldsymbol{\theta}_j^{(t)}$ *at layer $t$ under the Poisson likelihood as*
$$\boldsymbol{x}_j^{(t)} \sim \text{Pois}\left[-\boldsymbol{\Phi}^{(t)}\boldsymbol{\theta}_j^{(t)} \ln\left(1-p_j^{(t)}\right)\right].$$

**Corollary 2.** *With* $m_{kj}^{(t)(t+1)} := x_{\cdot jk}^{(t)} := \sum_{v=1}^{K_{t-1}} x_{vjk}^{(t)}$, *we can propagate the latent counts* $x_{vj}^{(t)}$ *of layer $t$ upward to layer $t + 1$ as*
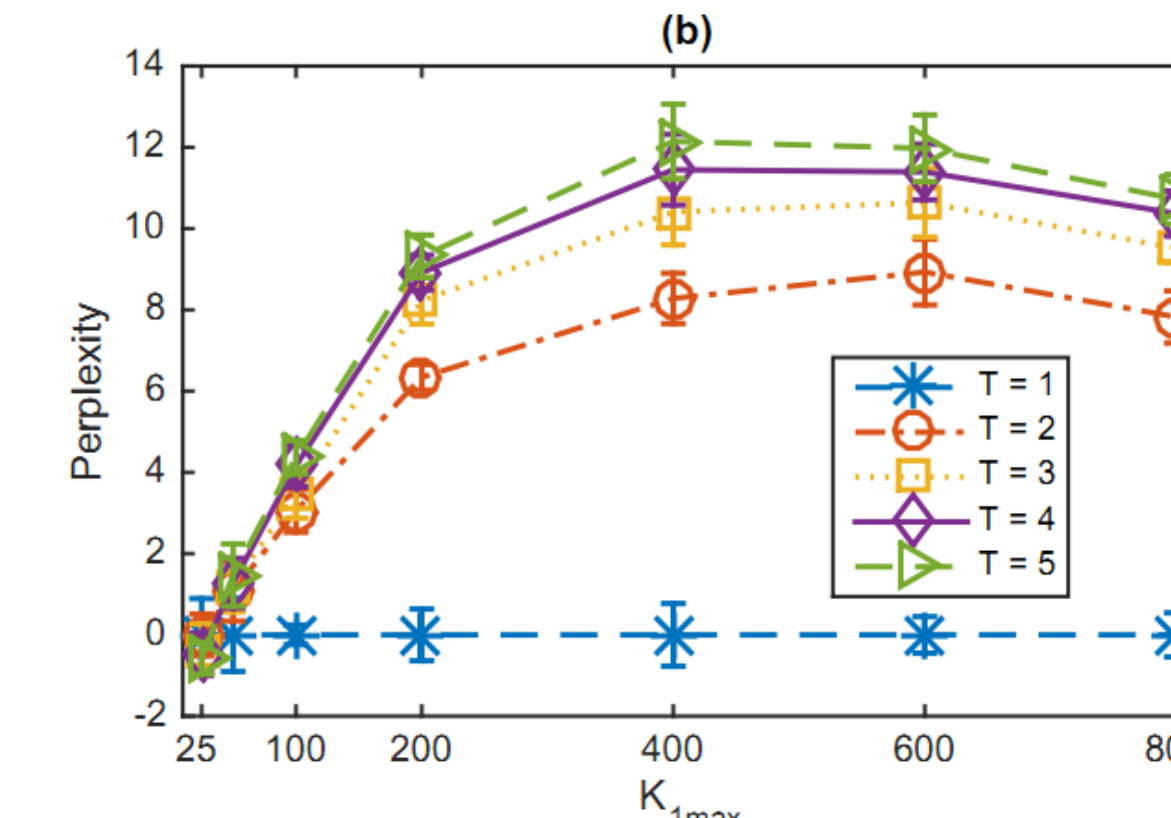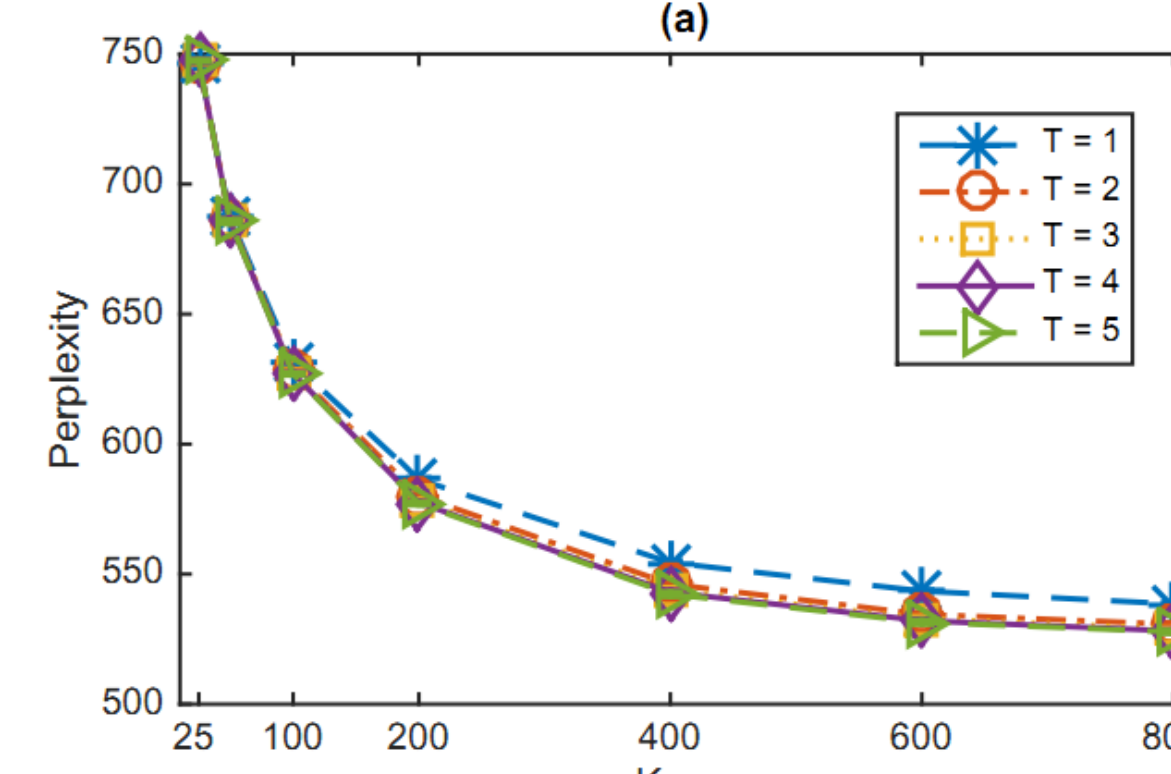$$\left\{\left(x_{vj1}^{(t)}, \ldots, x_{vjK_t}^{(t)}\right) \,\big|\, x_{vj}^{(t)}, \phi_{v:}^{(t)}, \boldsymbol{\theta}_j^{(t)}\right\} \sim \text{Mult}\left(x_{vj}^{(t)}, \frac{\phi_{v1}^{(t)}\theta_{1j}^{(t)}}{\sum_{k=1}^{K_t}\phi_{vk}^{(t)}\theta_{kj}^{(t)}}, \ldots, \frac{\phi_{vK_t}^{(t)}\theta_{K_tj}^{(t)}}{\sum_{k=1}^{K_t}\phi_{vk}^{(t)}\theta_{kj}^{(t)}}\right)$$
$$\left(x_{kj}^{(t+1)} \,\big|\, m_{kj}^{(t)(t+1)}, \phi_{k:}^{(t+1)}, \boldsymbol{\theta}_j^{(t+1)}\right) \sim \text{CRT}\left(m_{kj}^{(t)(t+1)}, \phi_{k:}^{(t+1)}\boldsymbol{\theta}_j^{(t+1)}\right)$$

## ❑ Projecting and ranking latent factors:

$$\mathbb{E}\left[\boldsymbol{x}_j^{(1)} \,\big|\, \boldsymbol{\theta}_j^{(t)}, \{\boldsymbol{\Phi}^{(\ell)}, c_j^{(\ell)}\}_{1,t}\right] = \left[\prod_{\ell=1}^t \boldsymbol{\Phi}^{(\ell)}\right]\frac{\boldsymbol{\theta}_j^{(t)}}{\prod_{\ell=2}^t c_j^{(\ell)}}$$

$$\mathbb{E}\left[\boldsymbol{\theta}_j^{(t)} \,\big|\, \{\boldsymbol{\Phi}^{(\ell)}, c_j^{(\ell)}\}_{t+1,T}, \boldsymbol{r}\right] = \left[\prod_{\ell=t+1}^T \boldsymbol{\Phi}^{(\ell)}\right]\frac{\boldsymbol{r}}{\prod_{\ell=t+1}^{T+1} c_j^{(\ell)}}$$

**Projection of the factors of layer $t$:**
$$\prod_{\ell=1}^t \boldsymbol{\Phi}^{(\ell)}$$

**Weights of the factors of layer $t$:**
$$\boldsymbol{r}^{(t)} := \left[\prod_{\ell=t+1}^T \boldsymbol{\Phi}^{(\ell)}\right]\boldsymbol{r}$$

1: **for** $T = 1, 2, \ldots, T_{\max}$ **do** Jointly train all the $T$ layers of the network
2:     Set $K_{T-1}$, the inferred width of layer $T-1$, as $K_{T\max}$, the upper bound of layer $T$'s width.
3:     **for** $iter = 1 : B_T + C_T$ **do** Upward-downward Gibbs sampling
4:         Sample $\{z_{ji}\}_{j,i}$ using collapsed inference; Calculate $\{x_{vjk}^{(1)}\}_{v,k,j}$; Sample $\{x_{vj}^{(2)}\}_{v,j}$;
5:         **for** $t = 2, 3, \ldots, T$ **do**
6:             Sample $\{x_{vjk}^{(t)}\}_{v,j,k}$; Sample $\{\phi_k^{(t)}\}_k$; Sample $\{x_{vj}^{(t+1)}\}_{v,j}$;
7:         **end for**
8:         **for** $t = T, T-1, \ldots, 2$ **do**
9:             Sample $c_j^{(t+1)}$ and calculate $p_j^{(t+1)}$; Sample $\boldsymbol{r}$ if $t = T$; Sample $\{\boldsymbol{\theta}_j^{(t)}\}_j$;
10:        **end for**
11:        Sample $p_j^{(2)}$ and Calculate $c_j^{(2)}$;
12:        **if** $iter = B_T$ **then**
13:            Prune layer $T$'s inactive factors $\{\phi_k^{(T)}\}_{k:x_{\cdot\cdot k}^{(T)}=0}$, let $K_T = \sum_k \delta(x_{\cdot\cdot k}^{(T)} > 0)$, and update $\boldsymbol{r}$;
14:        **end if**
15:    **end for**
16:    Output the posterior means (according to the last MCMC sample) of all remaining factors $\{\phi_k^{(t)}\}_{k,t}$ as the inferred network of $T$ layers, and $\{r_k\}_{k=1}^{K_T}$ as the gamma shape parameters of layer $T$'s hidden units.
17: **end for**

## ❑ Multi-class Classification ❑ Perplexities on NIPS12 Corpus



## Example Results

### ❑ A Tree on "Religion"



### ❑ Topics of Layer One on 20newsgroups



### ❑ Topics of Layer Three



### ❑ Topics of Layer Five



### ❑ A Tree related to "Turkey"



# NIPS 2015