

# Learning the structure of deep sparse directed graphical models

Ryan P. Adams<sup>1</sup>, Hanna M. Wallach<sup>2</sup>, and Zoubin Ghahramani<sup>3</sup>

<sup>1</sup>University of Toronto

<sup>2</sup>University of Massachusetts at Amherst

<sup>3</sup>University of Cambridge

Machine Learning Summer School  
August 2009

# Motivation

- ▶ Present some recent research on **graphical model** structure learning...
- ▶ ...related to **deep belief networks**...
- ▶ ...which uses **Markov chain Monte Carlo** inference...
- ▶ ...in a **non-parametric Bayesian model**.

# Deep networks

There is a great deal of interest on “deep belief networks”.

Deep belief nets are probabilistic generative models that are composed of multiple layers of stochastic, latent variables. The latent variables typically have binary values and are often called hidden units or feature detectors. The top two layers have undirected, symmetric connections between them and form an associative memory. The lower layers receive top-down, directed connections from the layer above. The states of the units in the lowest layer represent a data vector.

Geoffrey E. Hinton (2009) Scholarpedia.

# Deep networks

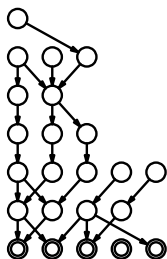
## Questions:

- ▶ How many layers should there be?
- ▶ How wide should each layer be?
- ▶ What sorts of units?

**Goal:** To learn the structure of a deep network.

**Approach:** A *nonparametric Bayesian* method that learns the structure of a layered *directed* deep belief network.

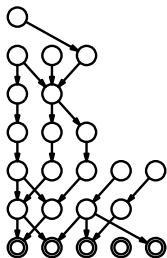
## Layered directed deep networks



$$p(\mathbf{x}) = \prod_{i=1}^K p(x_i | \mathbf{x}_{\pi_i})$$

Where  $\mathbf{x} = (x_1, \dots, x_K)$  and  $\pi_i$  are the parents of node  $i$ .  
aka Bayesian networks, probabilistic directed graphical models.  
Assume a layered graph structure.  
How many layers? How wide should each layer be?

## Priors over graph structures



Let  $z_{ij}^{(m)} = 1$  mean that  $j \in \pi_i$ , that is, node  $j$  is a parent of node  $i$  in layer  $m$ .

If we specify a sequence of matrices  $Z^{(0)}, Z^{(1)}, Z^{(2)}, \dots$  we have defined the layered graph structure.



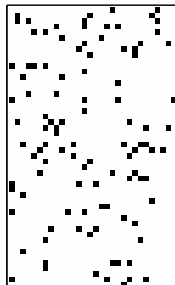
# Preview of the Indian buffet process

From finite to infinite matrices

$z_{ik} = 1$  means object  $i$  has feature  $k$ :

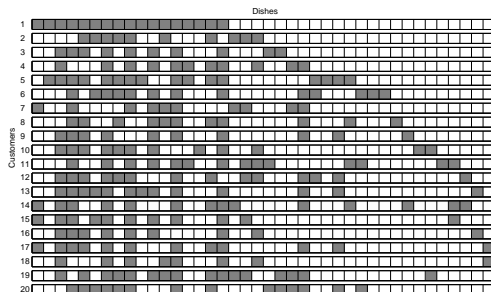
$$z_{ik} \sim \text{Bernoulli}(\theta_k)$$

$$\theta_k \sim \text{Beta}(\alpha/K, 1)$$



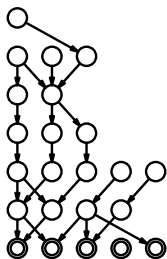
- ▶ Note that  $P(z_{ik} = 1|\alpha) = E(\theta_k) = \frac{\alpha/K}{\alpha/K+1}$ , so as  $K$  grows larger the matrix gets **sparser**.
- ▶ So if  $Z$  is  $N \times K$ , the expected number of nonzero entries is  $N\alpha/(1 + \alpha/K) < N\alpha$ .
- ▶ Even in the  $K \rightarrow \infty$  limit, the matrix is expected to have a finite number of non-zero entries.
- ▶ Two parameter extension  $\theta_k \sim \text{Beta}(\alpha\beta/K, \beta)$

# Indian buffet process



- ▶ First customer starts at the left of the buffet, and takes a serving from each dish, stopping after a  $\text{Poisson}(\alpha)$  number of dishes as his plate becomes overburdened.
- ▶ The  $n$ th customer moves along the buffet, sampling dishes in proportion to their popularity, serving himself dish  $k$  with probability  $m_k/n$ , and trying  $\text{Poisson}(\alpha/n)$  new dishes.
- ▶ The customer-dish matrix is the feature matrix,  $Z$ .

# Cascading Indian buffet process

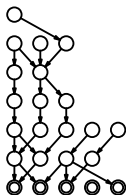


Start with  $K^{(0)}$  rows (visible units)

- ▶  $Z^{(0)} \sim IBP(\alpha, \beta)$  with  $K^{(0)}$  rows and  $K^{(1)}$  non-zero columns
- ▶  $Z^{(1)} \sim IBP(\alpha, \beta)$  with  $K^{(1)}$  rows and  $K^{(2)}$  non-zero columns
- ▶  $Z^{(2)} \sim IBP(\alpha, \beta)$  with  $K^{(2)}$  rows and  $K^{(3)}$  non-zero columns
- ▶ ...

This defines a sequences of infinite sparse binary matrices.

# Properties of the Cascading IBP



$$Z^{(m)} \sim IBP(\alpha, \beta) \quad \text{for } m = 0, 1, 2, \dots$$

- ▶ The expected in-degree of each unit (number of parents) is  $\alpha$ .
- ▶ The expected out-degree of each unit in  $m$  (number of children) is

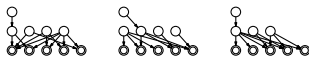
$$c(\beta, m) = 1 + \frac{K^{(m-1)} - 1}{1 + \beta}$$

Note that  $\lim_{\beta \rightarrow 0} c(\beta, m) = K^{(m-1)}$  and  $\lim_{\beta \rightarrow \infty} c(\beta, m) = 1$ .

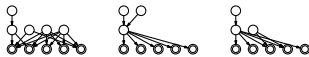
- ▶ Hidden units are exchangeable at each layer.
- ▶ **Theorem:** For  $K^{(m)} \in \mathbb{N}$ ,  $0 < \alpha < \infty$ ,  $0 < \beta < \infty$ , the sequence of  $K^{(m)}$  defined by the CIBP reaches the absorption state 0, with probability one, i.e.  $\lim_{m \rightarrow \infty} p(K^{(m)} = 0) = 1$ .

# Samples from the prior over structures

$$\alpha = 1, \beta = 1$$



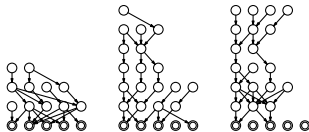
$$\alpha = 1, \beta = \frac{1}{2}$$



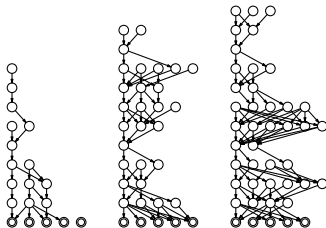
$$\alpha = \frac{1}{2}, \beta = 1$$



$$\alpha = 1, \beta = 2$$



$$\alpha = \frac{3}{2}, \beta = 1$$



Samples from the CIBP prior starting from five visible units.

## What kinds of units?

We want a model that is flexible enough to learn what types of unit it needs, ranging from binary to linear-Gaussian.

This idea was explored in [Nonlinear Gaussian belief networks](#) (NLGBNs) by [\(Frey and Hinton, 1999\)](#).

Let  $\mathbf{u}^{(m)}$  be the activity of units in layer  $m$ .

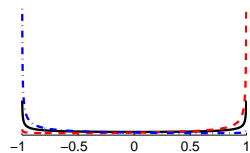
$$\mathbf{y}^{(m)} = (W^{(m+1)} \odot Z^{(m+1)})\mathbf{u}^{(m+1)} + \boldsymbol{\gamma}^{(m)}$$

where  $W$  is a weight matrix,  $\boldsymbol{\gamma}$  is a bias vector and  $\odot$  is Hadamard (elementwise) product.

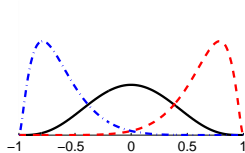
$$u_k^{(m)} = \sigma(y_k^{(m)} + \epsilon_k^{(m)})$$

$\sigma$  is a sigmoid function and noise  $\epsilon_k^{(m)} \sim \mathcal{N}(0, \frac{1}{\nu_k^{(m)}})$  has precision  $\nu_k^{(m)}$ .

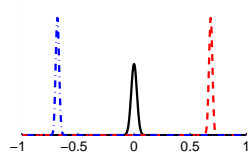
# NLGBN units



(a)  $\nu = \frac{1}{2}$



(b)  $\nu = 5$



(c)  $\nu = 1000$

Three modes of operation for the NLGBN unit. The black solid line shows the zero mean distribution, the red dashed line shows a pre-sigmoid mean of +1 and the blue dash-dot line shows a pre-sigmoid mean of -1.

(a) Binary behavior from small precision.

(b) Roughly Gaussian behavior from medium precision.

(c) Deterministic behavior from large precision.

# Inference

using Markov chain Monte Carlo

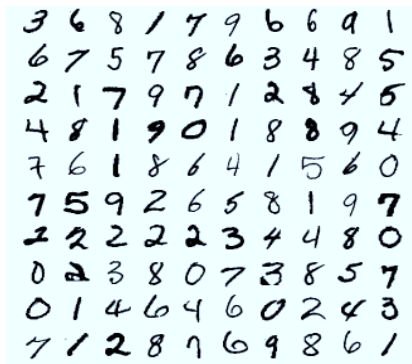
$$W \sim \mathcal{N} \quad \gamma \sim \mathcal{N} \quad \alpha \sim \mathcal{G} \quad \beta \sim \mathcal{G} \quad \nu \sim \mathcal{G}$$

We design an MCMC scheme to sample from the posterior:

$$p(\{Z^{(m)}, W^{(m)}\}_{m=1}^{\infty}, \{\gamma^{(m)}, \nu^{(m)}\}_{m=0}^{\infty}, \{\{\mathbf{u}_n^{(m)}\}_{m=1}^{\infty}\}_{n=1}^N | \{\mathbf{x}_n\}_{n=1}^N)$$

- ▶  $\mathbf{u}$  - slice sample
- ▶  $W$  and  $\gamma$  - Gibbs
- ▶  $\nu$  - Gibbs
- ▶  $Z$  - Gibbs (cf Algorithm 8 of CRPs)

## Experiments on MNIST data

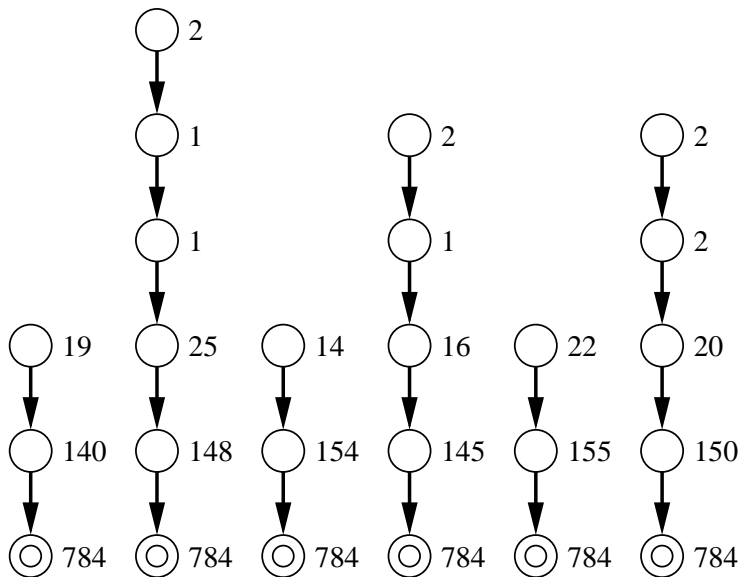


Small subset

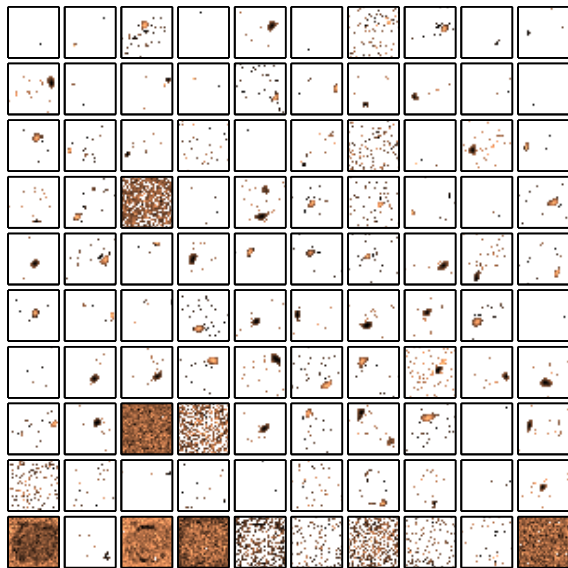
$28 \times 28$  pixels

100 images (10 from each class)

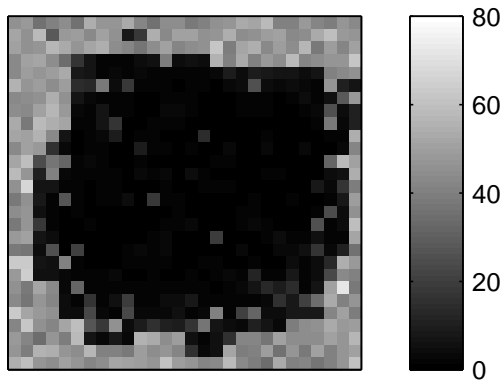
## Samples from Posterior over Structures



# First-Layer Features



# Visible Unit Precisions



# Summary

This work provides an initial attempt at addressing three issues with layered belief networks.

- ▶ It provides a way to learn belief networks that contain an arbitrary number of hidden units with nontrivial joint distributions due to a deep structure.
- ▶ It allows the units to have different operating regimes and infer appropriate local representations ranging from discrete binary to nonlinear continuous behavior.
- ▶ It provides a way to infer the appropriate directed graph structure of a layered network.

Initial work... many open questions!

Ryan P. Adams and Hanna M. Wallach

# Overall Summary

- ▶ Graphical models provide a powerful and intuitive framework for modelling and inference.
- ▶ Directed, undirected and factor graphs.
- ▶ Inference by message passing.
- ▶ Parameter and structure learning.
- ▶ A recent bit of research on structure learning.

Thanks!

# Questions