### Bayesian Deep Learning for Integrated Intelligence: Bridging the Gap between Perception and Inference

### Hao Wang Department of Computer Science and Engineering

Joint work with Naiyan Wang, Xingjian Shi, and Dit-Yan Yeung



# **Perception and Inference**

See (visual object recognition)
Read (text understanding)
Hear (speech recognition)



Think (inference and reasoning)

# Bayesian Deep Learning (BDL)



3

Bayesian deep learning

Deep Learning & Graphical Models

# **Perception and Inference**



#### **Perception component**

Content understanding





### **Task-Specific component**

Target task

### Bayesian deep learning (BDL)

Maximum a posteriori (MAP)
Markov chain Monte Carlo (MCMC)
Variational inference (VI)

# **Example: Medical Diagnosis**



**Bayesian deep learning (BDL)** 

# **Example: Movie Recommender Systems**



#### **Perception component**

### Content understanding





### Task-Specific component

Similarity, preferences Recommendation

Bayesian deep learning (BDL)

6

# **A Principled Probabilistic Framework**



[Wang et al. 2016]

# **BDL Models for Different Applications**

Applications	Models	Hinge Variables	Learning
Recommender Systems	CDL [Wang et al.]	$\{\mathbf{V}\}$	MAP
	Bayesian CDL [Wang et al.]	$\{\mathbf{V}\}$	Gibbs Sampling
	Marginalized CDL [Li et al.]	$\{\mathbf{V}\}$	MAP
	Symmetric CDL [Li et al.]	$\{\mathbf{V},\mathbf{U}\}$	MAP
	Collaborative Deep Ranking [Ying et al.]	$\{\mathbf{V}\}$	MAP
Topic Models	Relational SDAE [Wang et al.]	$\{\mathbf{S}\}$	MAP
	DPFA-SBN [Gan et al.]	$\{\mathbf{X}\}$	Hybrid MC
	DPFA-RBM [Gan et al.]	$\{\mathbf{X}\}$	Hybrid MC
Control	Embed to Control [Watter et al.]	$\{\mathbf{z}_t, \mathbf{z}_{t+1}\}$	Variational Inference

[Wang et al. 2016]

# Bayesian Deep Learning: Under a Principled Framework





[Wang et al. 2015 (KDD)]

# **Recommender Systems**



 $\checkmark$ 

7

# **Recommender Systems with Content**



Content information: Plots, directors, actors, etc.



# **Modeling the Content Information**



Handcrafted features





Automatically learn features and adapt for ratings

# **Modeling the Content Information**

### **1. Powerful features for content information**



# 2. Feedback from rating information Non-i.i.d.

**Collaborative deep learning** 

# **Deep Learning**



Stacked denoising autoencoders

Convolutional neural networks

y i

 $\hat{y}_t \sim P(\mathbf{y}_t \mid \mathbf{h}_t)$ 

 $P(\mathbf{y}_t \mid \boldsymbol{h}_t)$ 

Recurrent neural networks

# **Typically for i.i.d. data**

**Modeling the Content Information** 

### **1. Powerful features for content information**

### **Deep learning**

## 2. Feedback from rating information Non-i.i.d.

### **Collaborative deep learning (CDL)**

# Contribution

•Collaborative deep learning:

- \* deep learning for non-i.i.d. data
- \* joint representation learning and collaborative filtering

# Contribution

Collaborative deep learning

### •Complex target:

\* beyond targets like classification and regression

\* to complete a low-rank matrix

# Contribution

- •Collaborative deep learning
- Complex target
- First hierarchical Bayesian models for deep hybrid recommender system

# **Stacked Denoising Autoencoders (SDAE)**



Corrupted input

**Clean input** 

SDAE solves the following optimization problem:

$$\min_{\{\mathbf{W}_l\},\{\mathbf{b}_l\}} \|\mathbf{X}_c - \mathbf{X}_L\|_F^2 + \lambda \sum_l \|\mathbf{W}_l\|_F^2,$$

where  $\lambda$  is a regularization parameter and  $\|\cdot\|_F$  denotes the Frobenius norm. [Vincent et al. 2010]

# **Probabilistic Matrix Factorization (PMF)**

**Graphical model:** 



### Notation:

- $\left( \mathbf{v}_{\mathbf{j}}
  ight)$  latent vector of item j
- $\overline{\mathbf{U}_{i}}$  latent vector of user i
- **R**<sub>ij</sub> rating of item j from user i

#### **Generative process:**



$$= \prod_{i=1}^{N} \mathcal{N}(U_i|0, \sigma_U^2 \mathbf{I}) \qquad p(V|\sigma_V^2) = \prod_{j=1}^{M} \mathcal{N}(V_j|0, \sigma_V^2 \mathbf{I})$$
$$p(R|U, V, \sigma^2) = \prod_{i=1}^{N} \prod_{j=1}^{M} \left[ \mathcal{N}(R_{ij}|U_i^T V_j, \sigma^2) \right]^{I_{ij}}$$

**Objective function if using MAP:** 

 $p(U|\sigma_U^2) =$ 

$$E = \frac{1}{2} \sum_{i=1}^{N} \sum_{j=1}^{M} I_{ij} \left( R_{ij} - U_i^T V_j \right)^2 + \frac{\lambda_U}{2} \sum_{i=1}^{N} \parallel U_i \parallel_{Fro}^2 + \frac{\lambda_V}{2} \sum_{j=1}^{M} \parallel V_j \parallel_{Fro}^2$$
[Salakhutdinov et al. 2008]

# **Probabilistic SDAE**

#### **Graphical model:**



#### **Generative process:**



# **Collaborative Deep Learning (CDL)**

**Graphical model:** 



### **Collaborative deep learning**

SDAE

Two-way interaction

More powerful representation
Infer missing ratings from content
Infer missing content from ratings



### A Principled Probabilistic Framework (Recap)



[Wang et al. 2016]

# **CDL with Two Components**

### **Graphical model:**





Neural network representation for degenerated CDL

corrupted



26







maximizing the posterior probability is equivalent to maximizing the joint log-likelihood

$$\mathcal{L} = -\frac{\lambda_u}{2} \sum_i \|\mathbf{u}_i\|_2^2 - \frac{\lambda_w}{2} \sum_l (\|\mathbf{W}_l\|_F^2 + \|\mathbf{b}_l\|_2^2)$$
$$-\frac{\lambda_v}{2} \sum_j \|\mathbf{v}_j - \mathbf{X}_{\frac{L}{2},j*}^T\|_2^2 - \frac{\lambda_n}{2} \sum_j \|\mathbf{X}_{L,j*} - \mathbf{X}_{c,j*}\|_2^2$$
$$-\frac{\lambda_s}{2} \sum_l \sum_j \|\sigma(\mathbf{X}_{l-1,j*}\mathbf{W}_l + \mathbf{b}_l) - \mathbf{X}_{l,j*}\|_2^2$$
$$-\sum_{i,j} \frac{\mathbf{C}_{ij}}{2} (\mathbf{R}_{ij} - \mathbf{u}_i^T \mathbf{v}_j)^2.$$

 $ackslash \Lambda_u$  /

Prior (regularization) for user latent vectors, weights, and biases

$$\mathcal{L} = -\frac{\lambda_u}{2} \sum_i \|\mathbf{u}_i\|_2^2 - \frac{\lambda_w}{2} \sum_l (\|\mathbf{W}_l\|_F^2 + \|\mathbf{b}_l\|_2^2)$$
$$-\frac{\lambda_v}{2} \sum_j \|\mathbf{v}_j - \mathbf{X}_{\frac{L}{2},j*}^T\|_2^2 - \frac{\lambda_n}{2} \sum_j \|\mathbf{X}_{L,j*} - \mathbf{X}_{c,j*}\|_2^2$$
$$-\frac{\lambda_s}{2} \sum_l \sum_j \|\sigma(\mathbf{X}_{l-1,j*}\mathbf{W}_l + \mathbf{b}_l) - \mathbf{X}_{l,j*}\|_2^2$$
$$-\sum_{i,j} \frac{\mathbf{C}_{ij}}{2} (\mathbf{R}_{ij} - \mathbf{u}_i^T \mathbf{v}_j)^2.$$

Generating item latent vectors from content representation with Gaussian offset  $-\frac{\lambda_u}{2}\sum_{i} \|\mathbf{u}_i\|_2^2 - \frac{\lambda_w}{2}\sum_{i} (\|\mathbf{W}_l\|_F^2 + \|\mathbf{b}_l\|_2^2)$  $\frac{\lambda_{v}}{2} \sum_{i} \|\mathbf{v}_{j} - \mathbf{X}_{\frac{L}{2}, j*}^{T}\|_{2}^{2} - \frac{\lambda_{n}}{2} \sum_{i} \|\mathbf{X}_{L, j*} - \mathbf{X}_{c, j*}\|_{2}^{2}$  $-\frac{\lambda_s}{2}\sum_l \sum_l \|\sigma(\mathbf{X}_{l-1,j*}\mathbf{W}_l + \mathbf{b}_l) - \mathbf{X}_{l,j*}\|_2^2$  $-\sum_{i=1}^{T} \frac{\mathbf{C}_{ij}}{2} (\mathbf{R}_{ij} - \mathbf{u}_i^T \mathbf{v}_j)^2.$ 

'Generating' clean input from the output of probabilistic SDAE with Gaussian offset

$$\mathcal{L} = -\frac{\lambda_u}{2} \sum_i \|\mathbf{u}_i\|_2^2 - \frac{\lambda_w}{2} \sum_l (\|\mathbf{W}_l\|_F^2 + \|\mathbf{b}_l\|_2^2)$$
$$-\frac{\lambda_v}{2} \sum_j \|\mathbf{v}_j - \mathbf{X}_{\frac{L}{2},j*}^T\|_2^2 - \frac{\lambda_n}{2} \sum_j \|\mathbf{X}_{L,j*} - \mathbf{X}_{c,j*}\|_2^2$$
$$-\frac{\lambda_s}{2} \sum_l \sum_j \|\sigma(\mathbf{X}_{l-1,j*}\mathbf{W}_l + \mathbf{b}_l) - \mathbf{X}_{l,j*}\|_2^2$$
$$-\sum_{i,j} \frac{\mathbf{C}_{ij}}{2} (\mathbf{R}_{ij} - \mathbf{u}_i^T \mathbf{v}_j)^2.$$

Generating the input of Layer I from the output of Layer I-1 with Gaussian offset

$$\mathcal{L} = -\frac{\lambda_u}{2} \sum_i \|\mathbf{u}_i\|_2^2 - \frac{\lambda_w}{2} \sum_l (\|\mathbf{W}_l\|_F^2 + \|\mathbf{b}_l\|_2^2)$$
$$-\frac{\lambda_v}{2} \sum_j \|\mathbf{v}_j - \mathbf{X}_{\frac{L}{2},j*}^T\|_2^2 - \frac{\lambda_n}{2} \sum_j \|\mathbf{X}_{L,j*} - \mathbf{X}_{c,j*}\|_2^2$$
$$-\frac{\lambda_s}{2} \sum_l \sum_j \|\sigma(\mathbf{X}_{l-1,j*}\mathbf{W}_l + \mathbf{b}_l) - \mathbf{X}_{l,j*}\|_2^2$$
$$-\sum_{i,j} \frac{\mathbf{C}_{ij}}{2} (\mathbf{R}_{ij} - \mathbf{u}_i^T \mathbf{v}_j)^2.$$

measures the error of predicted ratings

$$\mathscr{L} = -\frac{\lambda_u}{2} \sum_i \|\mathbf{u}_i\|_2^2 - \frac{\lambda_w}{2} \sum_l (\|\mathbf{W}_l\|_F^2 + \|\mathbf{b}_l\|_2^2)$$
$$-\frac{\lambda_v}{2} \sum_j \|\mathbf{v}_j - \mathbf{X}_{\frac{L}{2},j*}^T\|_2^2 - \frac{\lambda_n}{2} \sum_j \|\mathbf{X}_{L,j*} - \mathbf{X}_{c,j*}\|_2^2$$
$$-\frac{\lambda_s}{2} \sum_l \sum_j \|\sigma(\mathbf{X}_{l-1,j*}\mathbf{W}_l + \mathbf{b}_l) - \mathbf{X}_{l,j*}\|_2^2$$
$$-\sum_{i,j} \frac{\mathbf{C}_{ij}}{2} (\mathbf{R}_{ij} - \mathbf{u}_i^T \mathbf{v}_j)^2.$$

 $(\lambda_u)$ 

**→**(**u**)

If  $\lambda_s$  goes to infinity, the likelihood simplifies to

$$\begin{aligned} \mathscr{L} &= -\frac{\lambda_u}{2} \sum_i \|\mathbf{u}_i\|_2^2 - \frac{\lambda_w}{2} \sum_l (\|\mathbf{W}_l\|_F^2 + \|\mathbf{b}_l\|_2^2) \\ &- \frac{\lambda_v}{2} \sum_j \|\mathbf{v}_j - f_e(\mathbf{X}_{0,j*}, \mathbf{W}^+)^T\|_2^2 \\ &- \frac{\lambda_n}{2} \sum_j \|f_r(\mathbf{X}_{0,j*}, \mathbf{W}^+) - \mathbf{X}_{c,j*}\|_2^2 \\ &- \sum_{i,j} \frac{\mathbf{C}_{ij}}{2} (\mathbf{R}_{ij} - \mathbf{u}_i^T \mathbf{v}_j)^2, \end{aligned}$$


## **Update Rules**

For U and V, use block coordinate descent:

$$\mathbf{u}_{i} \leftarrow (\mathbf{V}\mathbf{C}_{i}\mathbf{V}^{T} + \lambda_{u}\mathbf{I}_{K})^{-1}\mathbf{V}\mathbf{C}_{i}\mathbf{R}_{i}$$

$$\mathbf{v}_{j} \leftarrow (\mathbf{U}\mathbf{C}_{i}\mathbf{U}^{T} + \lambda_{v}\mathbf{I}_{K})^{-1}(\mathbf{U}\mathbf{C}_{j}\mathbf{R}_{j} + \lambda_{v}f_{e}(\mathbf{X}_{0,j*}, \mathbf{W}^{+})^{T})$$
For W and b, use a modified version of backpropagation
$$\nabla_{\mathbf{W}_{l}}\mathscr{L} = -\lambda_{w}\mathbf{W}_{l}$$

$$-\lambda_{v}\sum_{j}\nabla_{\mathbf{W}_{l}}f_{e}(\mathbf{X}_{0,j*}, \mathbf{W}^{+})^{T}(f_{e}(\mathbf{X}_{0,j*}, \mathbf{W}^{+})^{T} - \mathbf{v}_{j})$$

$$-\lambda_{n}\sum_{j}\nabla_{\mathbf{W}_{l}}f_{r}(\mathbf{X}_{0,j*}, \mathbf{W}^{+})(f_{r}(\mathbf{X}_{0,j*}, \mathbf{W}^{+}) - \mathbf{X}_{c,j*})$$

$$\nabla_{\mathbf{b}_{l}}\mathscr{L} = -\lambda_{w}\mathbf{b}_{l}$$

$$-\lambda_{v}\sum_{j}\nabla_{\mathbf{b}_{l}}f_{e}(\mathbf{X}_{0,j*}, \mathbf{W}^{+})^{T}(f_{e}(\mathbf{X}_{0,j*}, \mathbf{W}^{+})^{T} - \mathbf{v}_{j})$$

$$-\lambda_n \sum_j \nabla_{\mathbf{b}_l} f_r(\mathbf{X}_{0,j*}, \mathbf{W}^+) (f_r(\mathbf{X}_{0,j*}, \mathbf{W}^+) - \mathbf{X}_{c,j*})$$

### Datasets

	citeulike-a	citeulike-t	Netflix
#users	5551	7947	407261
#items	16980	25975	9228
#ratings	204987	134860	15348808

#### Collaborative Deep Learning for Recommender Systems ABSTRACT

ABSTRUCE Collaboration Billing (CD) is a successful approach com-financy and by many mecommerler systems. Conventional CP lossed nucleosis with the string stars to trains by aver-mentation. However, the ratings are afree very papers in space systems, and the string stars of the string star-sen systems and the string stars of the string star-tion of the space systems and the string stars of the string with the space systems and the string stars of the string string string strings and strings string strings and string the space systems and the string string strings and string strings and strings and strings and strings and string strings and strings and strings and strings and strings part and propose in the paper a hereached and strings. The string strings and collaborative Bitering for the rating (reduction string). The strings and strings and strings and strings. Extension sequences and strings and strings and strings. The strings are a string string string strings and strings. The strings are a string string strings and strings and collaborative Bitering for the rating (reduction) string. Extension sequences and strings and strings and the string of the string (Tota) and strings and strings and the strings of the string (Tota) and strings and strings and the string of the string (Tota) and strings and strings and the strings of the string (Tota) and strings and strings and the strings of the string (Tota) and strings and strings

#### Collaborative Deep Learning for Recommender Systems ABSTRACT

ABTRACT Calabranies' length (G7) is is mercandid approach com-mody and by many recommender systems. Conventional G7-based methods use the rating gives to items by users at the sole source of information for learning to make rec-mmany applications, consiting G7-based methods in degrade applications of the source of the source of the system of the sole source of information and the sole of the proper preparison (CTR) is an appealing record method is degrade to the spectra of the source of information and and the spectra of the source of information in and a proper preparison (CTR) is an appealing record method taking this approach which the acalitary information is usery parse-tion address this problem, we generative recent solutions in address the problem, we generative recent solutions in address the problem, we generative recent solutions in address the source of the source of the source in degrade source of the source of the source of the source of the solution of the source o



universe which alters their physical form in shocking ways. The four must learn to harness their new abilities and work together to save Earth from a former friend turned enemy.

Titles and abstracts Titles and abstracts

#### Movie plots

[Wang et al. 2011] [Wang et al. 2013]

### Content information

## **Evaluation Metrics**

### **Recall:**

recall@ $M = \frac{\text{number of items that the user likes among the top } M}{\text{total number of items that the user likes}}$ 

Mean Average Precision (mAP):

$$mAP = \frac{\sum_{q=1}^{Q} AveP(q)}{Q}$$
$$AveP = \frac{\sum_{k=1}^{n} (P(k) \times rel(k))}{\text{number of relevant items}}$$

Higher recall and mAP indicate better recommendation performance

## **Comparing Methods**

- **CMF**: Collective Matrix Factorization (Singh et al. 2008) is a model incorporating different sources of information by simultaneously factorizing multiple matrices.
- **SVDFeature**: SVDFeature (Chen et al. 2012) is a model for feature-based collaborative filtering.
- **DeepMusic**: DeepMusic (Oord et al. 2013) is a model for music recommendation.
- **CTR**: Collaborative Topic Regression (Wang et al. 2011) is a model performing topic modeling and collaborative filtering simultaneously.

Hybrid methods using BOW and ratings

Loosely coupled; interaction is not two-way





citeulike-t, dense setting

Netflix, dense setting

## Mean Average Precision (mAP)

	citeulike-a	citeulike-t	Netflix
CDL	0.0514	0.0453	0.0312
CTR	0.0236	0.0175	0.0223
DeepMusic	0.0159	0.0118	0.0167
CMF	0.0164	0.0104	0.0158
SVDFeature	0.0152	0.0103	0.0187

Exactly the same as Oord et al. 2013, we set the cutoff point at 500 for each user.

A relative performance boost of about 50%

## **Number of Layers**

### **Sparse Setting**

#layers	1	2	3
citeulike-a	27.89	31.06	30.70
citeulike-t	32.58	34.67	35.48
Netflix	29.20	30.50	31.01

### **Dense Setting**

#layers	1	2	3
citeulike- $a$	58.35	<b>59.43</b>	59.31
citeulike-t	52.68	53.81	<b>54.48</b>
Netflix	69.26	70.40	70.42

The best performance is achieved when the number of layers is **2 or 3** (**4 or 6** layers of generalized neural networks).

## **Example User**



**Romance** Moonstruck

### Movies





**True Romance** 

# training samples	2
	Swordfish
	A Fish Called Wanda
	Terminator 2
	A Clockwork Orange
Top 10 recommended	Sling Blade
movies by $\mathbf{CTR}$	Bridget Jones's Diary
	Raising Arizona
	A Streetcar Named Desire
	The Untouchables
	The Full Monty
	rno ran money
# training samples	2
# training samples	2 Snatch
# training samples	2 Snatch The Big Lebowski
# training samples	2 Snatch The Big Lebowski Pulp Fiction
# training samples	2 Snatch <b>The Big Lebowski</b> <b>Pulp Fiction</b> Kill Bill
# training samples Top 10 recommended	2 Snatch The Big Lebowski Pulp Fiction Kill Bill Raising Arizona
# training samples Top 10 recommended movies by <b>CDL</b>	2 Snatch <b>The Big Lebowski</b> <b>Pulp Fiction</b> Kill Bill <b>Raising Arizona</b> The Big Chill
# training samples Top 10 recommended movies by <b>CDL</b>	2 Snatch The Big Lebowski Pulp Fiction Kill Bill Raising Arizona The Big Chill Tootsie
# training samples Top 10 recommended movies by <b>CDL</b>	2 Snatch The Big Lebowski Pulp Fiction Kill Bill Raising Arizona The Big Chill Tootsie Sense and Sensibility
# training samples Top 10 recommended movies by <b>CDL</b>	2 Snatch The Big Lebowski Pulp Fiction Kill Bill Raising Arizona The Big Chill Tootsie Sense and Sensibility Sling Blade

## **Precision: 30% VS 20%**

## **Example User**

# training samples 4 Pulp Fiction A Clockwork Orange Being John Malkovich Raising Arizona Sling Blade Top 10 recommended movies by **CTR** Swordfish A Fish Called Wanda 51/15/15 Saving Grace C D L R.E C T I O N The Graduate Action & Monster's Ball **Johnny English** # training samples 4 Drama Pulp Fiction Movies Snatch The Usual Suspect Kill Bill KEVIN SI Top 10 recommended Momento The Big Lebowski movies by **CDL** One Flew Over the Cuckoo's Nest As Good as It Gets Goodfellas AMERICAN The Matrix

**American Beauty** 

## **Precision: 50% VS 20%**

## Example User





BEAUT

TOP GUN











# training samples	10
	Best in Snow
	Chocolat
	Good Will Hunting
	Monty Python and the Holy Grail
Top 10 recommended	Being John Malkovich
movies by <b>CTR</b>	Raising Arizona
	The Graduate
	Swordfish
	Tootsie
	Saving Private Ryan
# training samples	10
	Good Will Hunting
	Best in Show
	The Dig Leboughi
	The big Lebowski
	A Few Good Men
Top 10 recommended	A Few Good Men Monty Python and the Holy Grail
Top 10 recommended movies by <b>CDL</b>	A Few Good Men Monty Python and the Holy Grail Pulp Fiction
Top 10 recommended movies by <b>CDL</b>	A Few Good Men Monty Python and the Holy Grail Pulp Fiction The Matrix
Top 10 recommended movies by <b>CDL</b>	A Few Good Men Monty Python and the Holy Grail Pulp Fiction The Matrix Chocolat
Top 10 recommended movies by <b>CDL</b>	A Few Good Men Monty Python and the Holy Grail Pulp Fiction The Matrix Chocolat The Usual Suspect

## **Precision: 90% VS 50%**

## **Summary: Collaborative Deep Learning**



- Non-i.i.d (collaborative) deep learning
- •With a complex target
- First hierarchical Bayesian models for

hybrid deep recommender system

Significantly advance the state of the art

## **Marginalized CDL**

CDL:

#### Transformation to latent factors

$$\mathscr{L} = -\frac{\lambda_u}{2} \sum_i \|\mathbf{u}_i\|_2^2 - \frac{\lambda_v}{2} \sum_j \|\mathbf{v}_j - f_e(\mathbf{X}_{0,j*}, \mathbf{W}^+)^T\|_2^2$$
$$-\frac{\lambda_n}{2} \sum_j \|f_r(\mathbf{X}_{0,j*}, \mathbf{W}^+) - \mathbf{X}_{c,j*}\|_2^2 - \sum_{i,j} \frac{\mathbf{C}_{ij}}{2} (\mathbf{R}_{ij} - \mathbf{u}_i^T \mathbf{v}_j)^2$$
$$\mathbf{Reconstruction \ error}$$

Transformation to latent factors

$$\begin{aligned} \mathscr{L} &= -\frac{\lambda_u}{2} \sum_i \|\mathbf{u}_i\|_2^2 - \frac{\lambda_v}{2} \sum_j \|\mathbf{v}_j \mathbf{P}_1 - \mathbf{X}_{0,j*} \mathbf{W}_1\|_2^2 \\ \text{Marginalized CDL:} \quad -\sum_j \|\widetilde{\mathbf{X}}_{0,j*} \mathbf{W}_1 - \overline{\mathbf{X}}_{c,j*}\|_2^2 - \sum_{i,j} \frac{\mathbf{C}_{ij}}{2} (\mathbf{R}_{ij} - \mathbf{u}_i^T \mathbf{v}_j)^2 \end{aligned}$$

**Reconstruction error** 

[Li et al., CIKM 2015]

## **Collaborative Deep Ranking**



Fig. 1. The graphic model of CDR. SDAE with L = 4 is presented inside the dashed rectangle. Note that  $W^+$  denotes the set of weight matrices and bias vectors of all layers.

[Ying et al., PAKDD 2016]

### **Generative Process: Collaborative Deep Ranking**

- 1. For each layer l of the SDAE network,
  - (a) For each column q, draw the weight matrix and bias vector  $W_l^+$ , draw  $W_{l,*q}^+ \sim \mathcal{N}(0, \lambda_w^{-1} I_{K_l}).$
  - (b) For each row j of  $X_l$ , draw  $X_{l,j*} \sim \mathcal{N}(\sigma(X_{l-1,j*}W_l + b_l), \lambda_s^{-1}I_{K_l})$
- 2. For each item j,
  - (a) Draw a clean input  $X_{c,j*} \sim \mathcal{N}(X_{L,j*}, \lambda_n^{-1} I_m)$
  - (b) Draw a latent item offset vector  $\epsilon_j \sim \mathcal{N}(0, \lambda_v^{-1}I_K)$  and then set the latent item vector to be:

$$v_j = \epsilon_j + X_{\frac{L}{2},j}^T,$$

- 3. For each user i,
  - (a) Draw user factor vector  $u_i \sim \mathcal{N}(0, \lambda_u^{-1} I_K)$
  - (b) For each pair-wise preference  $(j,k) \in \mathcal{P}_i$ , where  $\mathcal{P}_i = \{(j,k) : r_{ij} r_{ik} > 0\}$ , draw the estimator,

$$\delta_{ijk} \sim \mathcal{N}(u_i^T v_j - u_i^T v_k, c_{ijk}^{-1})$$

## Symmetric CDL



Both item content and user attributesUser attributes: age, gender, occupation, country,city, geolacation, domain, etc[Li et al., CIKM 2015]

## Symmetric CDL



## **Other Extensions of CDL**

Word2vec, tf-idf

- Sampling-based, variational inference
- Tagging information, networks

# Relational Stacked Denoising Autoencoders

[Wang et al. 2015 (AAAI)]

## **BDL for Topic Models and Relational Learning**



#### **BDL-Based Topic Models**

## **Relational SDAE as Relational Topic Models**



**BDL-Based Topic Models** 

[Wang et al. 2015 (AAAI)]

## **Relational SDAE: Motivation**



- Unsupervised representation learning
- Enhance representation power with relational information

## **Probabilistic SDAE**

#### **Graphical model:**



#### **Generative process:**



## **Relational SDAE: Graphical Model**



## **Relational SDAE: Two Components**





## **Relational SDAE: Generative Process**

Oraw the relational latent matrix S from a matrix variate normal distribution:

$$\mathbf{S} \sim \mathcal{N}_{K,J}(0, \mathbf{I}_K \otimes (\lambda_l \mathscr{L}_a)^{-1}).$$

- **2** For layer l of the SDAE where  $l = 1, 2, \ldots, \frac{L}{2} 1$ ,
  - For each column n of the weight matrix  $\mathbf{W}_l$ , draw  $\mathbf{W}_{l,*n} \sim \mathcal{N}(0, \lambda_w^{-1} \mathbf{I}_{K_l})$ .
  - **2** Draw the bias vector  $\mathbf{b}_l \sim \mathcal{N}(0, \lambda_w^{-1} \mathbf{I}_{K_l})$ .
  - For each row j of  $X_l$ , draw

$$\mathbf{X}_{l,j*} \sim \mathcal{N}(\sigma(\mathbf{X}_{l-1,j*}\mathbf{W}_l + \mathbf{b}_l), \lambda_s^{-1}\mathbf{I}_{K_l}).$$

**③** For layer  $\frac{L}{2}$  of the SDAE network, draw the representation vector for item j from the product of two Gaussians (PoG):

$$\mathbf{X}_{\frac{L}{2},j*} \sim \mathsf{PoG}(\sigma(\mathbf{X}_{\frac{L}{2}-1,j*}\mathbf{W}_l + \mathbf{b}_l), \mathbf{s}_j^T, \lambda_s^{-1}\mathbf{I}_K, \lambda_r^{-1}\mathbf{I}_K).$$

## **Relational SDAE: Generative Process**

- For layer l of the SDAE network where  $l = \frac{L}{2} + 1, \frac{L}{2} + 2, \dots, L$ ,
  - For each column n of the weight matrix  $\mathbf{W}_l$ , draw  $\mathbf{W}_{l,*n} \sim \mathcal{N}(0, \lambda_w^{-1} \mathbf{I}_{K_l}).$
  - Draw the bias vector  $\mathbf{b}_l \sim \mathcal{N}(0, \lambda_w^{-1} \mathbf{I}_{K_l})$ .
  - **3** For each row j of  $\mathbf{X}_l$ , draw

$$\mathbf{X}_{l,j*} \sim \mathcal{N}(\sigma(\mathbf{X}_{l-1,j*}\mathbf{W}_l + \mathbf{b}_l), \lambda_s^{-1}\mathbf{I}_{K_l}).$$

**2** For each item j, draw a clean input

$$\mathbf{X}_{c,j*} \sim \mathcal{N}(\mathbf{X}_{L,j*}, \lambda_n^{-1} \mathbf{I}_B).$$

## **Multi-Relational SDAE: Graphical Model**



#### Product of Q+1 Gaussians

#### Multiple networks:

citation networks co-author networks



## **Relational SDAE: Objective Function**

$$\mathscr{L} = -\frac{\lambda_l}{2} \operatorname{tr}(\mathbf{S}\mathscr{L}_a \mathbf{S}^T) - \frac{\lambda_r}{2} \sum_j \|(\mathbf{s}_j^T - \mathbf{X}_{\frac{L}{2}, j*})\|_2^2 - \frac{\lambda_w}{2} \sum_l (\|\mathbf{W}_l\|_F^2 + \|\mathbf{b}_l\|_2^2) - \frac{\lambda_w}{2} \sum_l \|\mathbf{X}_{L, j*} - \mathbf{X}_{c, j*}\|_2^2 - \frac{\lambda_s}{2} \sum_l \sum_j \|\sigma(\mathbf{X}_{l-1, j*} \mathbf{W}_l + \mathbf{b}_l) - \mathbf{X}_{l, j*}\|_2^2$$

Similar to the generalized SDAE, taking  $\lambda_s$  to infinity, the joint log-likelihood becomes:

$$\mathcal{L} = -\frac{\lambda_l}{2} \operatorname{tr}(\mathbf{S}\mathcal{L}_a \mathbf{S}^T) - \frac{\lambda_r}{2} \sum_j \|(\mathbf{s}_j^T - \mathbf{X}_{\frac{L}{2}, j*})\|_2^2 - \frac{\lambda_w}{2} \sum_l (\|\mathbf{W}_l\|_F^2 + \|\mathbf{b}_l\|_2^2) - \frac{\lambda_n}{2} \sum_j \|\mathbf{X}_{L, j*} - \mathbf{X}_{c, j*}\|_2^2,$$

Network  $A \rightarrow$  Relational Matrix S



Relational Matrix  $S \rightarrow$  Middle-Layer Representations

## **Update Rules**

For S:

$$\mathbf{S}_{k*}(t+1) \leftarrow \mathbf{S}_{k*}(t) + \delta(t)r(t)$$
$$r(t) \leftarrow \lambda_r \mathbf{X}_{\frac{L}{2},*k}^T - (\lambda_l \mathscr{L}_a + \lambda_r \mathbf{I}_J)\mathbf{S}_{k*}(t)$$
$$\delta(t) \leftarrow \frac{r(t)^T r(t)}{r(t)^T (\lambda_l \mathscr{L}_a + \lambda_r \mathbf{I}_J)r(t)}.$$

For X, W, and b: Use Back Propagation.

### From Representation to Tag Recommendation

Objective function:

$$\mathscr{L} = -\frac{\lambda_u}{2} \sum_i \|\mathbf{u}_i\|_2^2 - \frac{\lambda_v}{2} \sum_j \|\mathbf{v}_j - \mathbf{X}_{\frac{L}{2},j*}^T\|_2^2$$
$$-\sum_{i,j} \frac{c_{ij}}{2} (\mathbf{R}_{ij} - \mathbf{u}_i^T \mathbf{v}_j)^2,$$

where  $\lambda_u$  and  $\lambda_v$  are hyperparameters.  $c_{ij}$  is set to 1 for the existing ratings and 0.01 for the missing entries.

# Algorithm

1. Learning representation:

repeat Update S using the updating rules Update X, W, and b until convergence Get resulting representation  $X_{\frac{L}{2},j*}$ 

**2.** Learning  $u_i$  and  $v_j$ :

Optimize the objective function  $\mathscr{L}$ 

3. Recommend tags to items according to the predicted  $\mathbf{R}_{ij}$ :  $\mathbf{R}_{ij} = \mathbf{u}_i^T \mathbf{v}_j$ Rank  $\mathbf{R}_{1j}, \mathbf{R}_{2j}, \dots, \mathbf{R}_{Ij}$ Recommend tags with largest  $\mathbf{R}_{ij}$  to item j

### Datasets

### Description of datasets

	citeulike-a	citeulike-t	movielens-plot
#items	16980	25975	7261
#tags	7386	8311	2988
#tag-item paris	204987	134860	51301
#relations	44709	32665	543621

## Sparse Setting, citeulike-a



## Dense Setting, citeulike-a



## Sparse Setting, movielens-plot



## **Dense Setting, movielens-plot**


## Case Study 1: Tagging Scientific Articles

#### An example article with recommended tags

Example Article	Title: Mining the Peanut Gallery: Opinion Extraction and Semantic				
	Top topic 1: language, text, mining, representation, semantic, concepts				
	words, relations, processing, categories				
Top 10 tags	SDAE	True?	RSDAE	True?	
	1. instance	no	<ol> <li>sentiment_analysis</li> </ol>	no	
	2. consumer	yes	2. instance	no	
	<ol><li>sentiment_analysis</li></ol>	no	3. consumer	yes	
	4. summary	no	4. summary	no	
	5. 31july09	no	5. sentiment	yes	
	6. medline	no	6. product_review_mining	yes	
	7. eit2	no	7. sentiment_classification	yes	
	8. l2r	no	8. 31july09	no	
	9. exploration	no	9. opinion_mining	yes	
	10. biomedical	no	10. product	yes	

#### **Precision: 10% VS 60%**

## Case Study 2: Tagging Movies (SDAE)

An example movie with recommended tags

	Title: E.T. the Extra-Terrestrial			
Example Movie	Top topic 1: crew, must, on, earth, human, save, ship, rescue,			
	by, find, scientist, planet			
	SDAE	True tag?		
	1. Saturn Award (Best Special Effects)	yes		
	2. Want	no		
	3. Saturn Award (Best Fantasy Film)	no		
	4. Saturn Award (Best Writing)	yes		
Top 10 recommended tags	5. Cool but freaky	no		
	6. Saturn Award (Best Director)	no		
	7. Oscar (Best Editing)	no		
	8. almost favorite	no		
	9. Steven Spielberg	yes		
	10. sequel better than original	no		

#### **Precision: 30% VS 60%**

## Case Study 2: Tagging Movies (RSDAE)

An example movie with recommended tags

	Title: E.T. the Extra-Terrestrial			
Example Movie	Top topic 1: crew, must, on, earth, human, save, ship, rescue,			
	by, find, scientist, planet			
	RSDAE	True tag?		
	1. Steven Spielberg	yes		
	2. Saturn Award (Best Special Effects)	yes		
	3. Saturn Award (Best Writing)	yes		
	4. Oscar (Best Editing)	no		
Top 10 recommended tags	5. Want	no		
	6. Liam Neeson	no		
	7. AFI 100 (Cheers)	yes		
	8. Oscar (Best Sound)	yes		
	9. Saturn Award (Best Director)	no		
	10. Oscar (Best Music - Original Score)	yes		

Does not appear in the tag lists of movies linked to 'E.T. the Extra-Terrestrial'

Very difficult to discover this tag

#### **Relational SDAE as Deep Relational Topic Models**



**BDL-Based Topic Models** 

Unified into a probabilistic relational model

for relational deep learning

[Wang et al. 2015 (AAAI)]

### Applications of Bayesian Deep Learning: Under a Principled Framework



## **Take-home Messages**

- Probabilistic graphical models for formulating both representation learning and inference/reasoning components
- Learnable representation serving as a bridge
- Tight, two-way interaction is crucial

### **Future Goals**



#### General Framework:

- 1. Ability of understanding text, images, and videos
- 2. Ability of inference and planning under uncertainty
- 3. Close the gap between human intelligence and artificial intelligence





# Thanks! Q&A