

CoNet: Collaborative Cross Networks for Cross-Domain Recommendation

Guangneng Hu*, Yu Zhang, and Qiang Yang

CIKM 2018
Oct 22-26 (Mo-Fr),
Turin, Italy



香港科技大學
THE HONG KONG
UNIVERSITY OF SCIENCE
AND TECHNOLOGY

Recommendations Are Ubiquitous: Products, Medias, Entertainment...

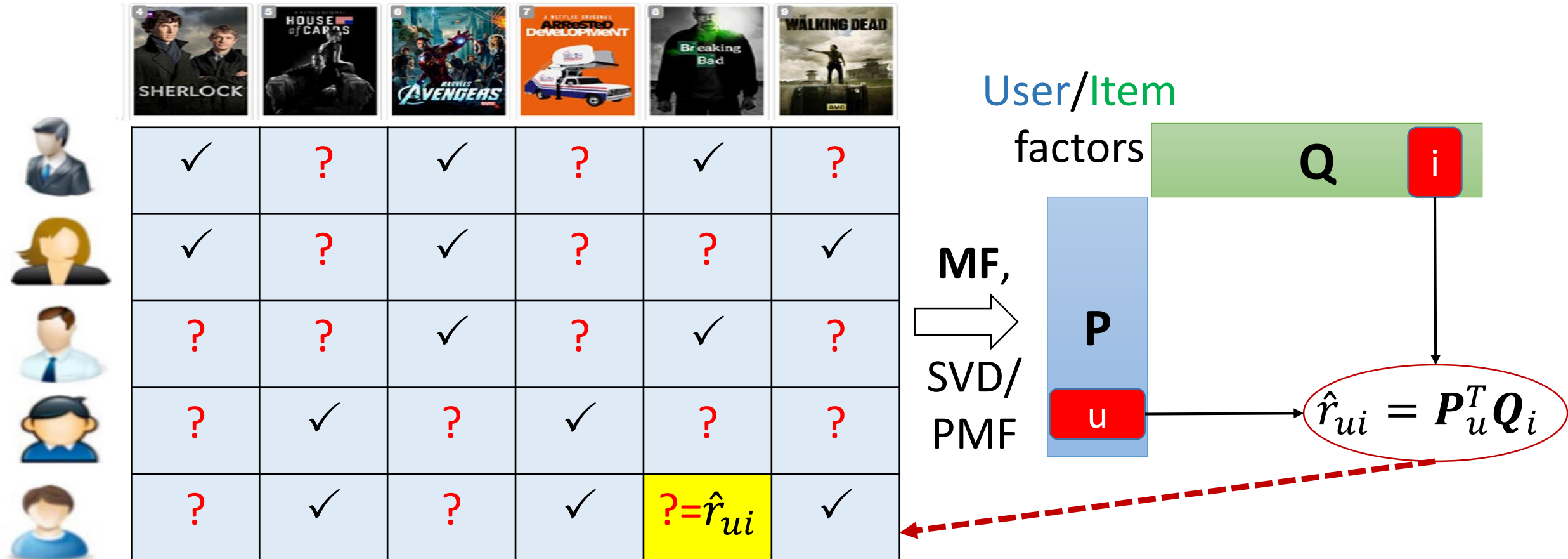
- Amazon
 - 300 million customers
 - 564 million products
- Netflix
 - 480,189 users
 - 17,770 movies
- Spotify
 - 40 million songs
- OkCupid
 - 10 million members

The image displays four distinct examples of recommendation systems:

- Amazon:** A screenshot of the 'Recommended for You' section for the book 'Applied Predictive Modeling' by Max Kuhn. It shows the book cover, a 'LOOK INSIDE!' button, and pricing information (List Price: \$89.95, Price: \$65.81).
- Netflix:** A screenshot of the Netflix homepage featuring a 'Netflix Prize' banner and a 'Movies For You' section with personalized movie recommendations.
- OkCupid:** A screenshot of the 'Today's Most Popular!' section, showing a grid of user profile pictures.
- News:** A screenshot of a news article recommendation for 'Glenn Frey, a Founding Member of the Eagles, Dies at 67' from the New York Times.

Typical Methods: Matrix Factorization

(Koren KDD'08, KDD 2018 TEST OF TIME award)

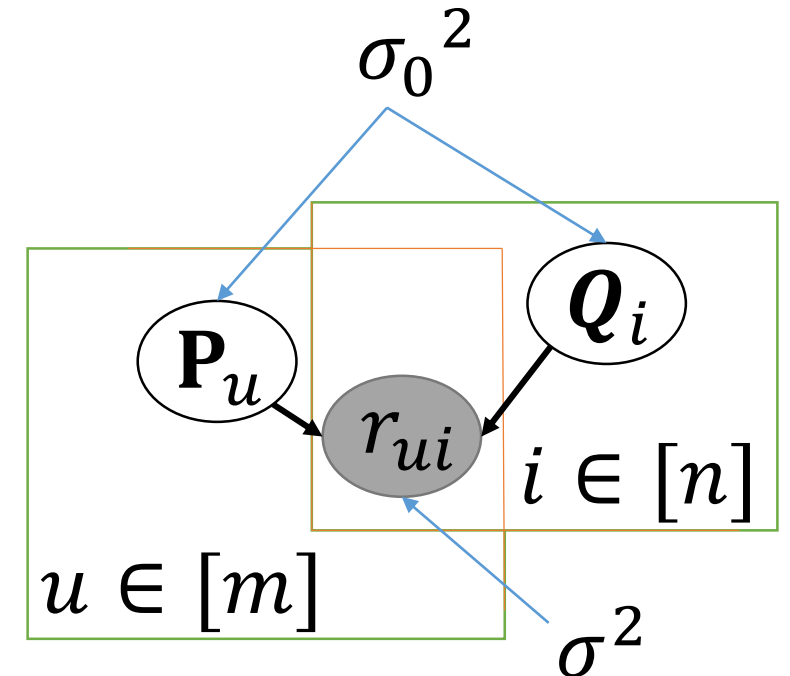


Probabilistic Interpretations: PMF

- The objective of matrix factorization

$$\min_{\mathbf{P}, \mathbf{Q}} \sum_{r_{ui} \neq 0} (r_{ui} - \hat{r}_{ui})^2 + \lambda (\|\mathbf{P}\|_{Frob}^2 + \|\mathbf{Q}\|_{Frob}^2)$$

- Probabilistic interpretations (PMF)
 - Gaussian observations & priors
- Log posterior distribution



$$\ln p(\Theta | \mathbf{R}, \Phi) = -\frac{1}{2\sigma^2} \sum_{u,i} \delta(r_{ui}) (r_{ui} - \mathbf{P}_u^T \mathbf{Q}_i)^2 - \frac{1}{2\sigma_0^2} (\|\mathbf{P}\|_{Frob}^2 + \|\mathbf{Q}\|_{Frob}^2)$$

- Maximum a posteriori (**MAP**) estimation \leftrightarrow Minimizing sum-of-squared-errors with quadratic regularization (**Loss + Regu**)

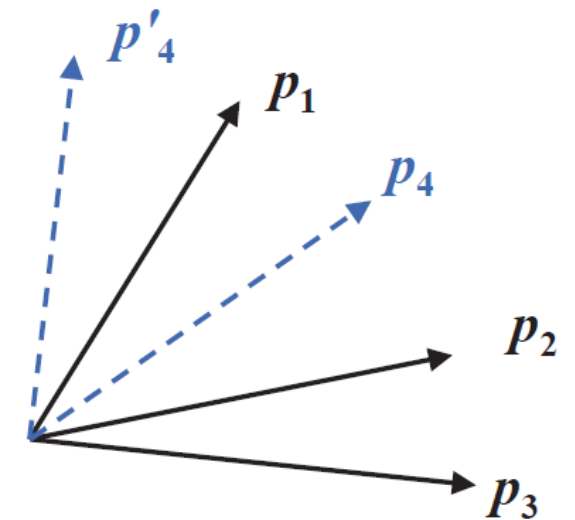
Limited Expressiveness of MF: Example I

- Similarity of user u_4 :
 - Given: $\text{Sim}(u_4, u_1) > \text{Sim}(u_4, u_3) > \text{Sim}(u_4, u_2)$
 - Q: Where to put the latent factor vector p_4 ?
- MF can not capture highly nonlinear
 - Deep learning, nonlinearity

	i_1	i_2	i_3	i_4	i_5
u_1	1	1	1	0	1
u_2	0	1	1	0	0
u_3	0	1	1	1	0
u_4	1	0	1	1	1

← items →

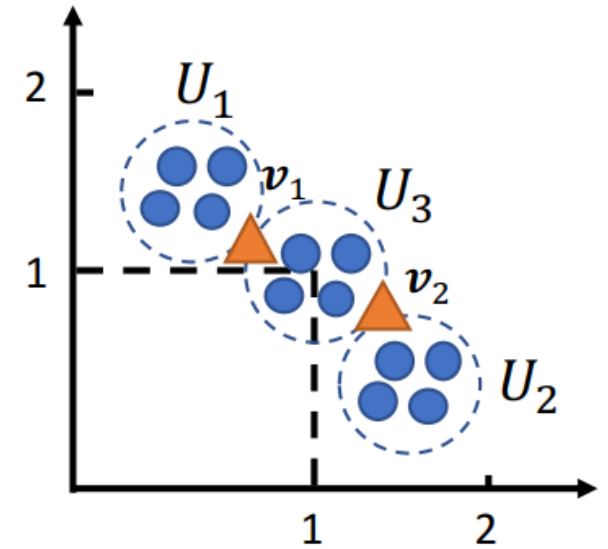
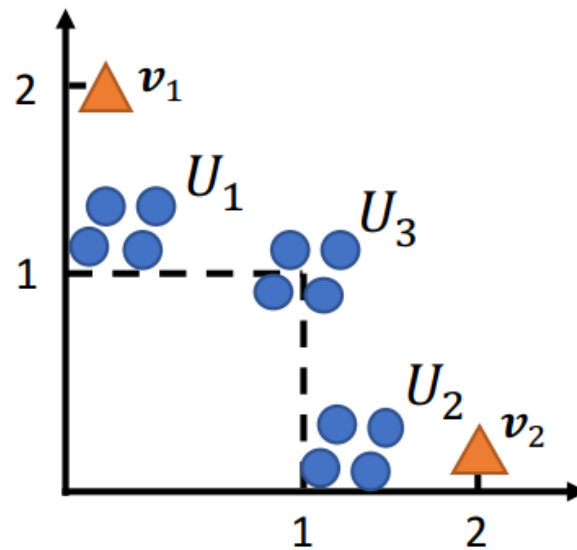
↑ users ↓



Xiangnan He et al. Neural collaborative filtering. WWW'17

Limited Expressiveness of MF: Example II

- Transitivity of user U_3 :
 - Given: U_3 close to item v_1 and v_2
 - Q: Where v_1 and v_2 should be?
- MF can not capture transitivity
 - Metric learning, triangle inequality



Cheng-Kang Hsieh et al. Collaborative metric learning. WWW'17

Modelling Nonlinearity: Generalized Matrix Factorization

- Matrix factorization as a single layer **linear** neural network
 - Input: one-hot encodings of the user and item indices (u, i)
 - Embedding: embedding matrices (P, Q)
 - Output: **Hadamard product** between embeddings with an **identity activation** and a fixed **all-one vector h**
- Generalized Matrix Factorization
 - Learning weights \mathbf{h} instead of fixing it
 - Using non-linear activation (e.g., sigmoid) instead of identity

The diagram shows the equation $\hat{r}_{u,i} = \sigma \left(\mathbf{h}^T (P_u \odot Q_i) \right)$. Annotations include: a blue arrow pointing to the Hadamard product symbol \odot with the label "Hadamard product"; a green arrow pointing to the sigmoid function σ with the label "identity activation"; and a red arrow pointing to the vector \mathbf{h} with the label "all-one vector".

$$\hat{r}_{u,i} = \sigma \left(\mathbf{h}^T (P_u \odot Q_i) \right)$$

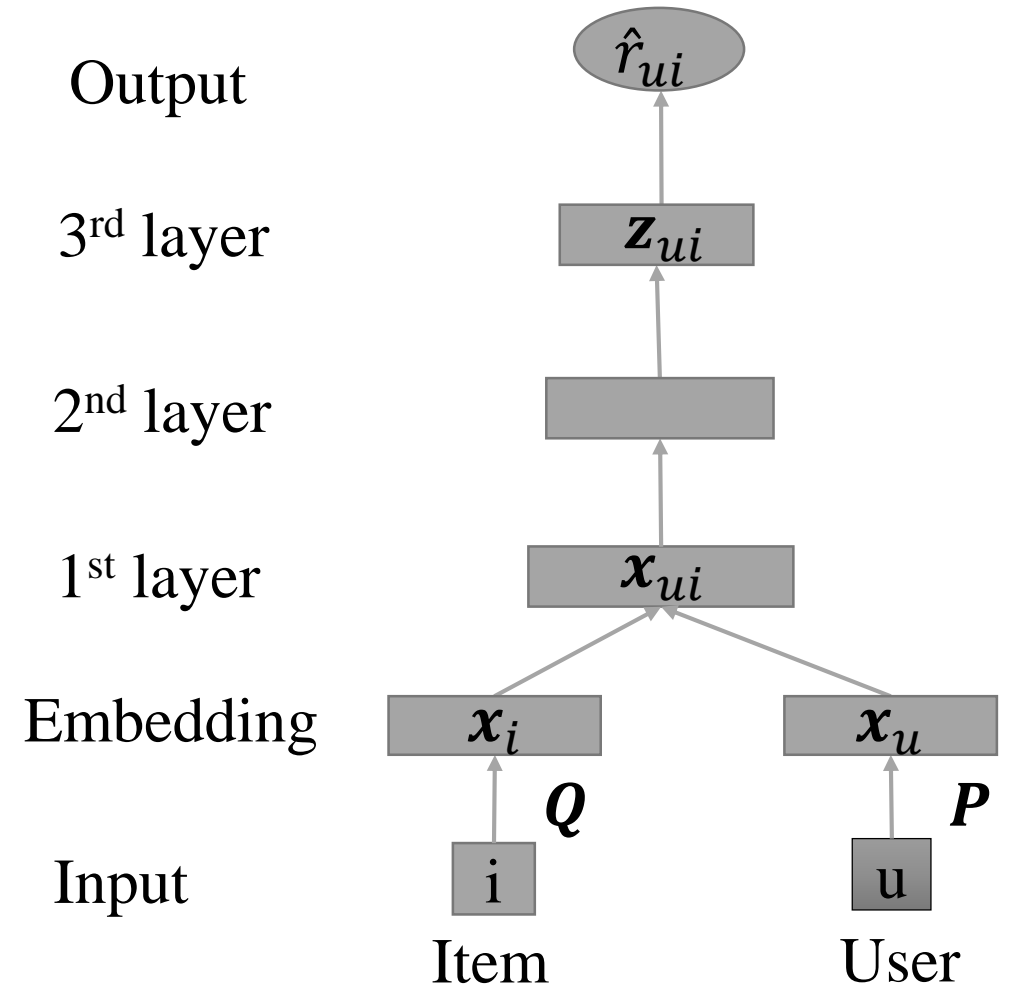
identity activation all-one vector

Go Deeper: Neural Collaborative Filtering

- Stack multilayer feedforward NNs to learn highly non-linear representations


$$f(\mathbf{x}_{ui} | \mathbf{P}, \mathbf{Q}, \theta_f) = \phi_o(\phi_L(\dots(\phi_1(\mathbf{x}_{ui}))\dots))$$

- Capture the complex user-item interaction relationships via the expressiveness of multilayer NNs

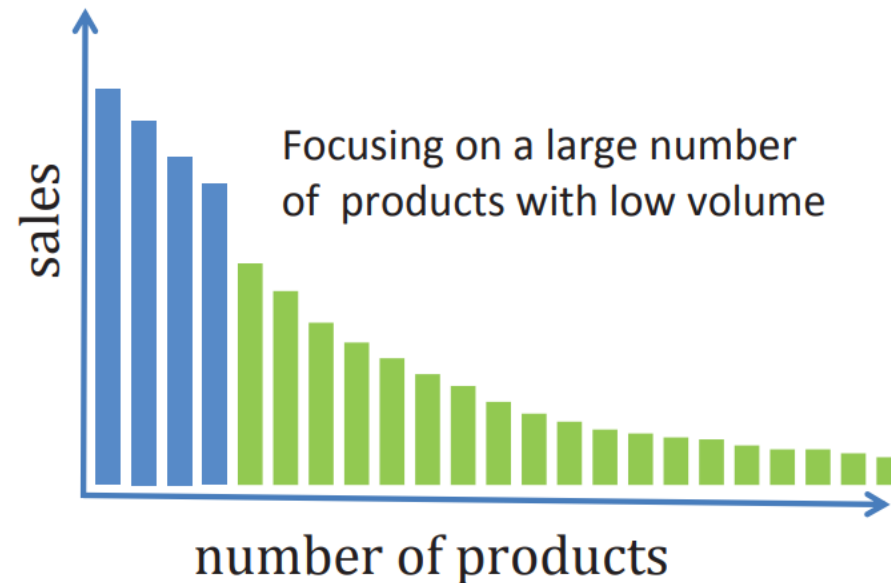


Collaborative Filtering Faces Challenges: Data Sparsity and Long Tail

- Data sparsity
 - Netflix
 - **1.225%**
 - Amazon
 - **0.017%**
- Long tail
 - Pareto principle (80/20 rule):
 - A small proportion (e.g., 20%) of products generate a large proportion (e.g., 80%) of sales



	SHERLOCK	HOUSE OF CARDS	AVENGERS	AERONAUT DEVELOPMENT	Breaking Bad	WALKING DEAD
User 1	2	?	2	?	5	?
User 2	5	?	4	?	?	1
User 3	?	?	5	?	2	?
User 4	?	1	?	5	?	?
User 5	?	5	?	1	?	4



A Solution: Cross-Domain Recommendation

- Two domains
 - A target domain (e.g., Books domain) $\mathbf{R}=\{(u,i)\}$,
 - A related source domain (e.g., Movies domain) $\{(u,j)\}$
- Probability of a user prefers an item by two factors
 - His/her individual preferences (in the target domain), and
 - His/her behavior in a related source domain

	The Name of the Wind	American Gods	The Lord of the Rings (2001)	The Matrix (1999)	Star Wars (1977)
Alice	5		?		
Bob	4		5		
Carol		4			3
Dave			5	5	

u_A (Alice, Bob) and u_B (Carol, Dave) are indicated by brackets on the left.

 I_A (Books) and I_B (Movies) are indicated by brackets at the bottom with corresponding icons.

$$\hat{r}_{ui} \triangleq p(r_{ui} = 1 | u, [j]^u)$$

Typical Methods: Collective Matrix Factorization (Singh & Gordon, KDD'08)

- User-Item interaction matrix **R**
- Relational domain: Item-Genre content matrix **Y**
- Sharing the **item-specific** latent feature matrix **Q**

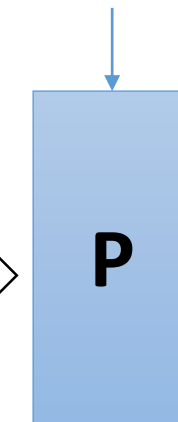
movie	budget	gross	genre	year
Goodfellas	25M	47M	crime	1990
My Cousin Vinny	11M	64M	comedy	1992
...
Clue	15M	15M	comedy	1985

User x Movie

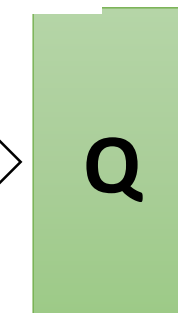
Movie x Genre

$$R \approx PQ^T, Y \approx QW^T$$

User factors



Shared item factors



Genre factors

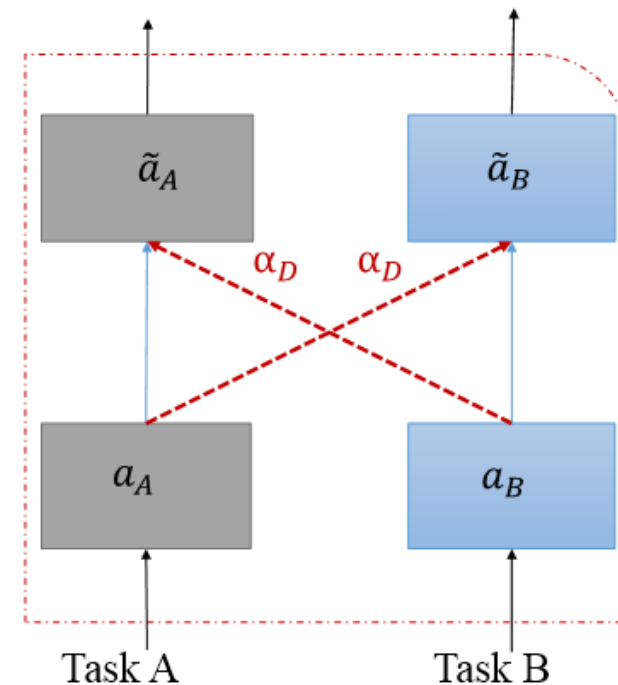


Deep Methods: Cross-Stitch Networks (CSN)

- Linear combination of activation maps from two tasks

$$\tilde{a}_A^{ij} = \alpha_S a_A^{ij} + \alpha_D a_B^{ij}, \quad \tilde{a}_B^{ij} = \alpha_S a_B^{ij} + \alpha_D a_A^{ij},$$

- Strong assumptions (SA)
 - SA 1: Representations from other network are **equally important** with weights being all the same scalar
 - SA 2: Representations from other network are **all useful** since it transfers activations from every location in a dense way

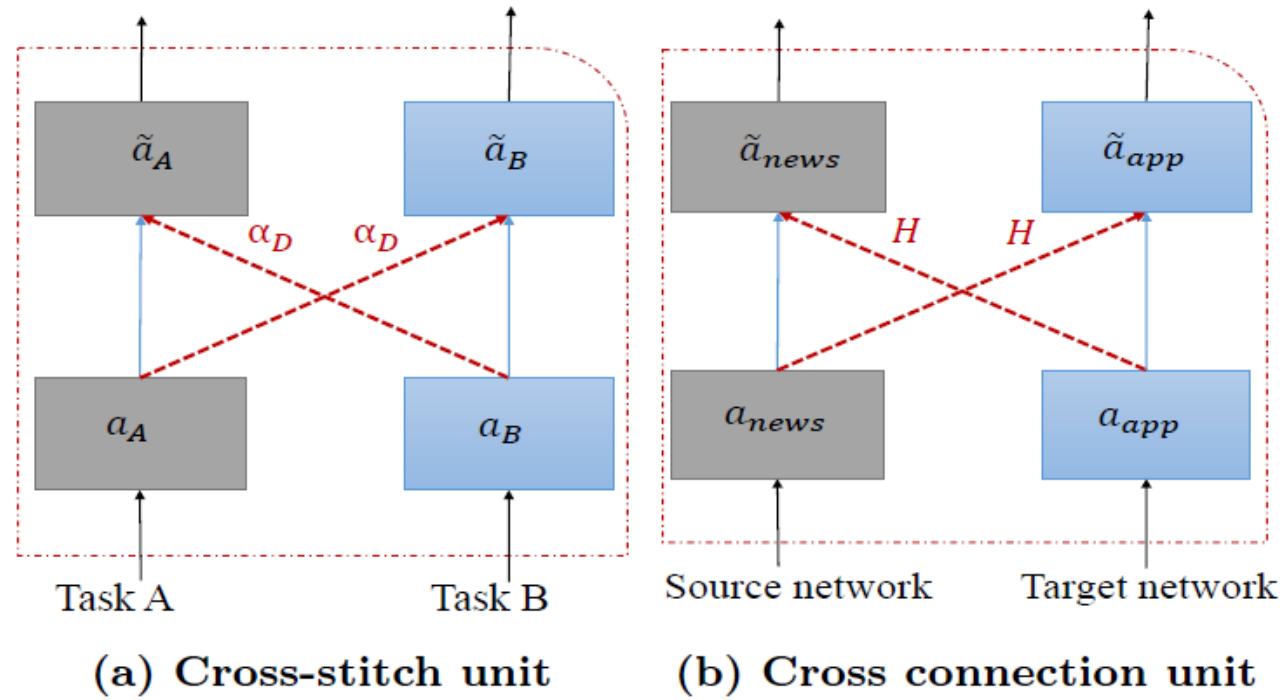


The Proposed Collaborative Cross Networks

- We propose a novel deep transfer learning method, Collaborative Cross Networks, to
 - Alleviate the data sparsity issue faced by the deep collaborative filtering
 - By transferring knowledge from a related source domain
 - Relax the strong assumptions faced by the existing cross-domain recommendation
 - By transferring knowledge via a matrix and enforcing sparsity-induced regularization

Idea 1: Using a matrix rather than a scalar (used in cross-stitch networks) to transfer

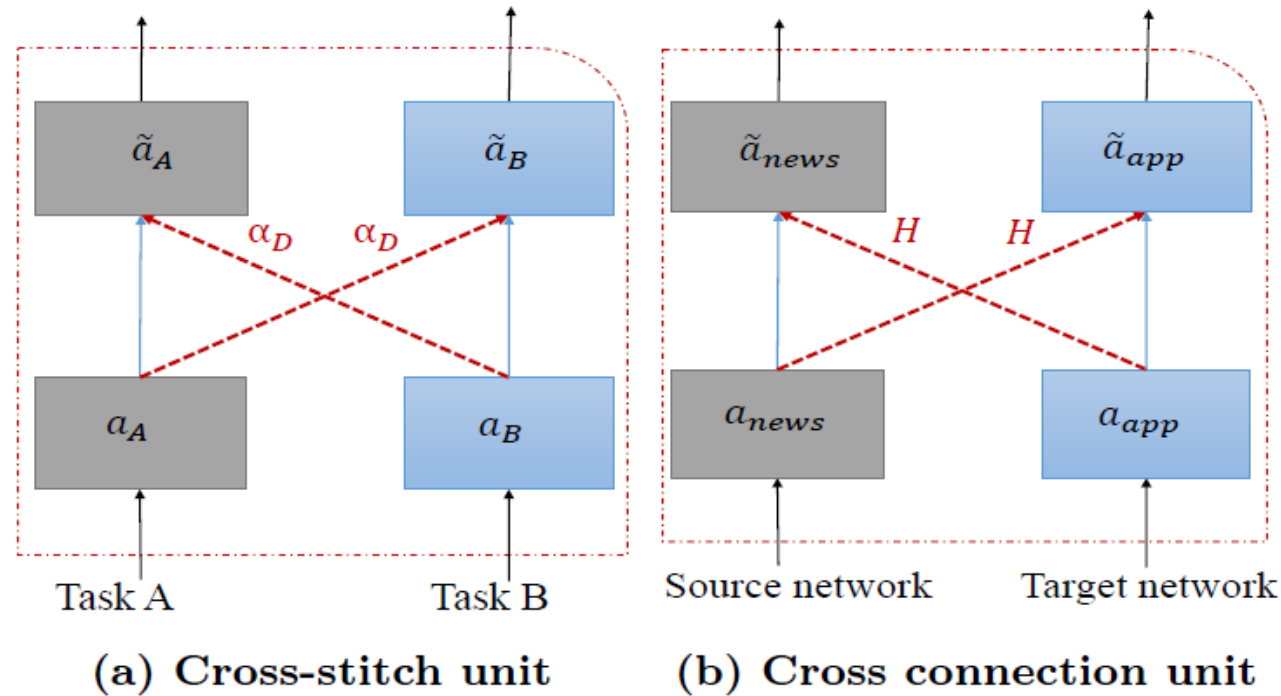
- We can relax the SA 1 assumption (equally important)



$$\mathbf{a}_{app}^{l+1} = \sigma(\mathbf{W}_{app}^l \mathbf{a}_{app}^l + \mathbf{H}^l \mathbf{a}_{news}^l),$$
$$\mathbf{a}_{news}^{l+1} = \sigma(\mathbf{W}_{news}^l \mathbf{a}_{news}^l + \mathbf{H}^l \mathbf{a}_{app}^l)$$

Idea 2: Selecting representations via sparsity-induced regularization

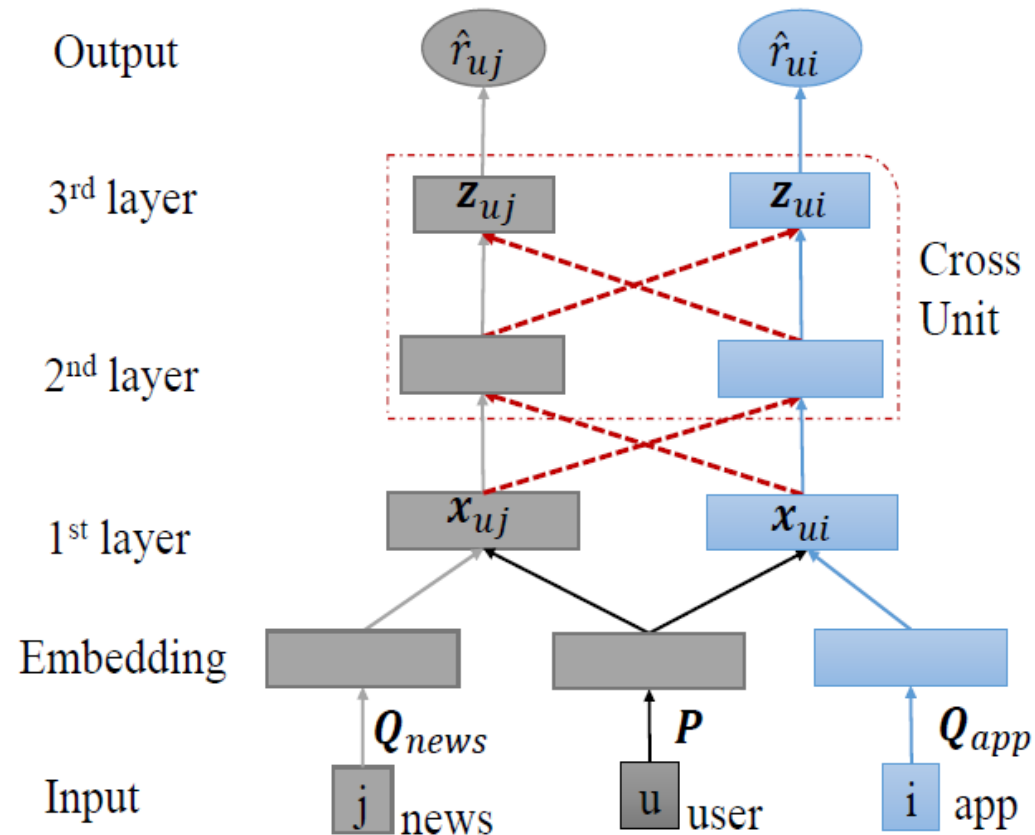
- We can relax the SA 2 assumption (all useful)



$$\Omega(\mathbf{H}^l) = \lambda \sum_{i=1}^r \sum_{j=1}^p |h_{ij}|$$

The Architecture of the CoNet Model

- A version of three hidden layers and two cross units



Model Learning Objective

- The likelihood function (randomly sample negative examples)

$$L(\Theta|\mathcal{S}) = \prod_{(u,i) \in \mathbf{R}_T^+} \hat{r}_{ui} \prod_{(u,i) \in \mathbf{R}_T^-} (1 - \hat{r}_{ui});$$

- The negative logarithm likelihood \leftrightarrow Binary cross-entropy loss

$$\mathcal{L} = - \sum_{(u,i) \in \mathcal{S}} r_{ui} \log \hat{r}_{ui} + (1 - r_{ui}) \log(1 - \hat{r}_{ui});$$

- Stochastic gradient descent (and variants)

$$\Theta^{new} \leftarrow \Theta^{old} - \eta \frac{\partial L(\Theta)}{\partial \Theta}$$

Model Learning Objective (cont')

- Basic model (CoNet)

$$\mathcal{L}(\Theta) = \mathcal{L}_{app}(\Theta_{app}) + \mathcal{L}_{news}(\Theta_{news})$$

- Adaptive model (SCoNet)
 - Added the sparsity-induced penalty term into the basic model
- Typical deep learning library like TensorFlow (<https://www.tensorflow.org>) provides automatic differentiation which can be computed by chain rule in back-propagation.

Complexity Analysis

- Model analysis

The model parameters Θ include $\{P, (H^l)_{l=1}^L\} \cup \{Q_{app}, (W_{app}^l, b_{app}^l)_{l=1}^L, h_{app}\} \cup \{Q_{news}, (W_{news}^l, b_{news}^l)_{l=1}^L, h_{news}\}$,

- Linear with the input size and is close to the size of typical latent factors models and neural CF approaches

- Learning analysis

- Update the target network using the target domain data and update the source network using the source domain data
- The learning procedure is similar to the cross-stitch networks. And the cost of learning each base network is approximately equal to that of running a typical neural CF approach

Dataset and Evaluation Metrics

Dataset	#Users	Target Domain			Source Domain		
		#Items	#Interactions	Density	#Items	#Interactions	Density
Mobile	23,111	14,348	1,164,394	0.351%	29,921	617,146	0.089%
Amazon	80,763	93,799	1,323,101	0.017%	35,896	963,373	0.033%

- Mobile: Apps and News
- Amazon: Books and Movies
- A higher value (HR, NDCG, MRR) with lower cutoff $topK$ indicates better performance

$$HR = \frac{1}{|\mathcal{U}|} \sum_{u \in \mathcal{U}} \delta(p_u \leq topK),$$

$$NDCG = \frac{1}{|\mathcal{U}|} \sum_{u \in \mathcal{U}} \frac{\log 2}{\log(p_u + 1)},$$

$$MRR = \frac{1}{|\mathcal{U}|} \sum_{u \in \mathcal{U}} \frac{1}{p_u}.$$

Baselines

- BPRMF: Bayesian personalized ranking
- MLP: Multilayer perceptron
- MLP++: Combine two MLPs by sharing the user embedding matrix
- CDCF: Cross-domain CF with factorization machines
- CMF: Collective MF
- CSN: The cross-stitch network

Baselines	Shallow method	Deep method
Single-domain	BPRMF [36]	MLP [13]
Cross-domain	CDCF [24], CMF [37]	MLP++, CSN [27]

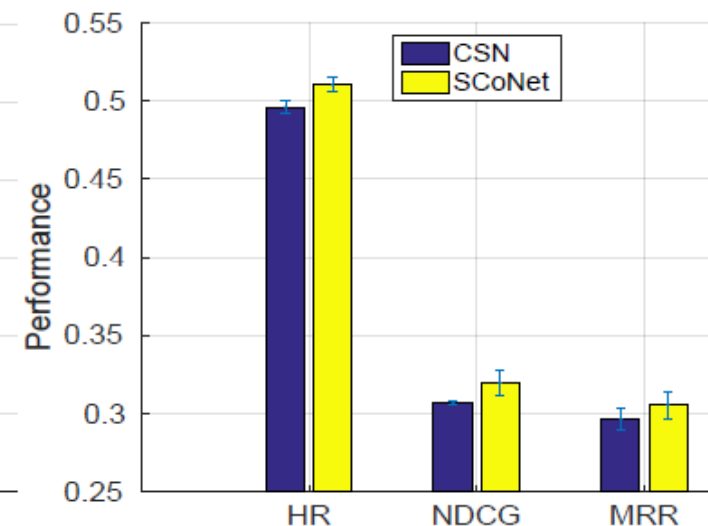
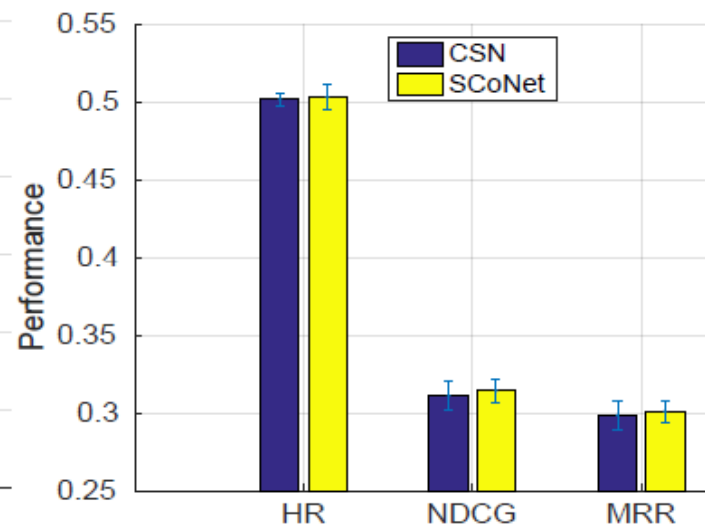
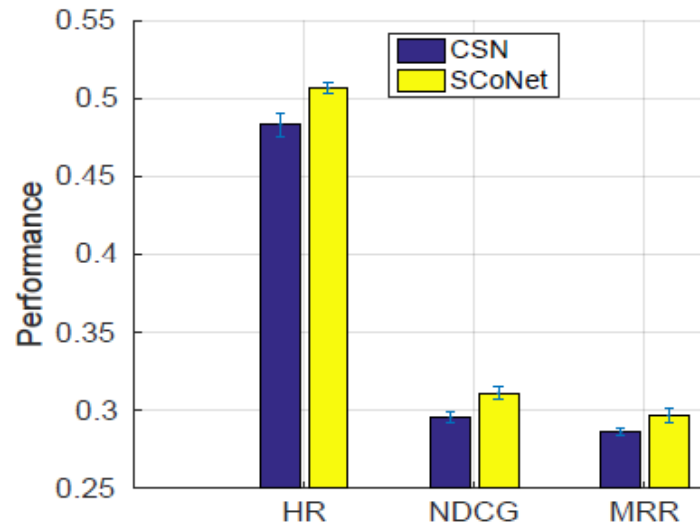
Comparing Different Approaches

- CSN has some difficulty in benefitting from knowledge transfer on the Amazon since it is inferior to the non-transfer base network MLP
- The proposed model outperforms baselines on real-world datasets under three ranking metrics

Dataset	Metric	BPRMF	CMF	CDCF	MLP	MLP++	CSN	CoNet	SCoNet	improve
Mobile	HR	.6175	.7879	.7812	.8405	.8445	.8458*	.8480	.8583	1.47%
	NDCG	.4891	.5740	.5875	.6615	.6683	.6733*	.6754	.6887	2.29%
	MRR	.4489	.5067	.5265	.6210	.6268	.6366*	.6373	.6475	1.71%
Amazon	HR	.4723	.3712	.3685	.5014	.5050*	.4962	.5167	.5338	5.70%
	NDCG	.3016	.2378	.2307	.3143	.3175*	.3068	.3261	.3424	7.84%
	MRR	.2971	.1966	.1884	.3113*	.3053	.2964	.3163	.3351	7.65%

Impact of Selecting Representations

- Configurations are $\{16, 32, 64\} * 4$, on Mobile data
- Naïve transfer learning approach may confront the negative transfer
- We demonstrate the necessity of adaptively selecting representations to transfer



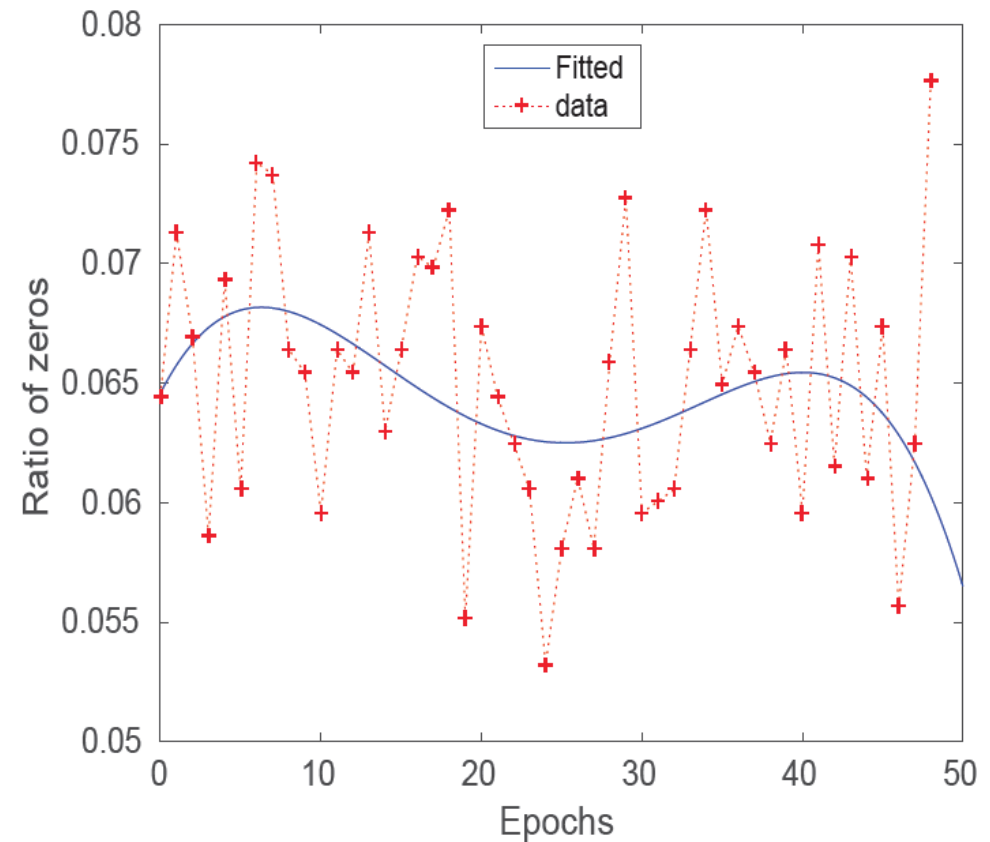
Benefit of Transferring Knowledge

- The more training examples we can reduce, the more benefit we can get from transferring knowledge
- Our model can reduce tens of thousands training examples by comparing with non-transfer methods without performance degradation

Dataset	Method	Reduction		HR	NDCG	MRR
		percent	amount			
Mobile	MLP	0%	0	.8405	.6615	.6210
	SCoNet	0%	0	.8547	.6802	.6431
		2.05%	23,031	.8439	.6640	.6238
		4.06%	45,468	.8347*	.6515*	.6115*
Amazon	MLP	0%	0	.5014	.3143	.3113
	SCoNet	0%	0	.5338	.3424	.3351
		1.11%	12,850	.5110	.3209	.3080*
		2.18%	25,318	.4946*	.3082*	.2968*

Analysis: Ratio of Zeros in Transfer Matrix H

- The percent of zero entries in transfer matrix is 6.5%
- A 4-order polynomial to robustly fit the data
- It may be better to transfer many instead of all representations



Conclusions and Future Works

- In general,
 - Neural/Deep approaches are better than shallow models,
 - Transfer learning approaches are better than non-transfer ones,
 - Shallow models are mainly based on MF techniques,
 - Deep models can be based on various NNs (MLP, CNN, RNN),
- Future works,
 - Data privacy
 - Source domain can not share the raw data, but model parameters
 - Transferable graph convolutional networks

Thanks!

Q & A

Acknowledgment: SIGIR Student Travel Grant