

Recommender Systems: Two Threads and Their Meeting

Guangneng Hu

22 Mar 2019 (Fri),
WeChat, Guangzhou



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

Outline

- Introduction
- Collaborative filtering
 - Matrix factorization, metric learning, neural approaches
- Cross-domain recommendation
 - Collective matrix factorization
 - Deep transfer learning
- Hybrid filtering
 - Personalized word embeddings
- Transfer meets hybrid
- Conclusion

Introduction

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

- Amazon
 - 300 million customers
 - 564 million products
- Netflix
 - 480,189 users
 - 17,770 movies
- WeChat
 - 474,726 groups
 - 245,352,140 users
- 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, author information, and pricing details.
- 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 OkCupid website showing a grid of 'Today's Most Popular!' dating profiles.
- 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.

Evaluating Recommender Systems

- Accuracy of predictions

- Root Mean Square Error (RMSE)
 - E.g. Netflix grand prize \$1M
- Mean Absolute Error (MAE)

$$RMSE_{\mathcal{T}} = \sqrt{\sum_{(u_i, v_j) \in \mathcal{T}} (R_{i,j} - \hat{R}_{i,j})^2 / |\mathcal{T}|}$$

- Accuracy of classifications

- Hit Rate/Ratio (HR)
- Precision, Recall, F1, ROC curves

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

- Accuracy of ranks

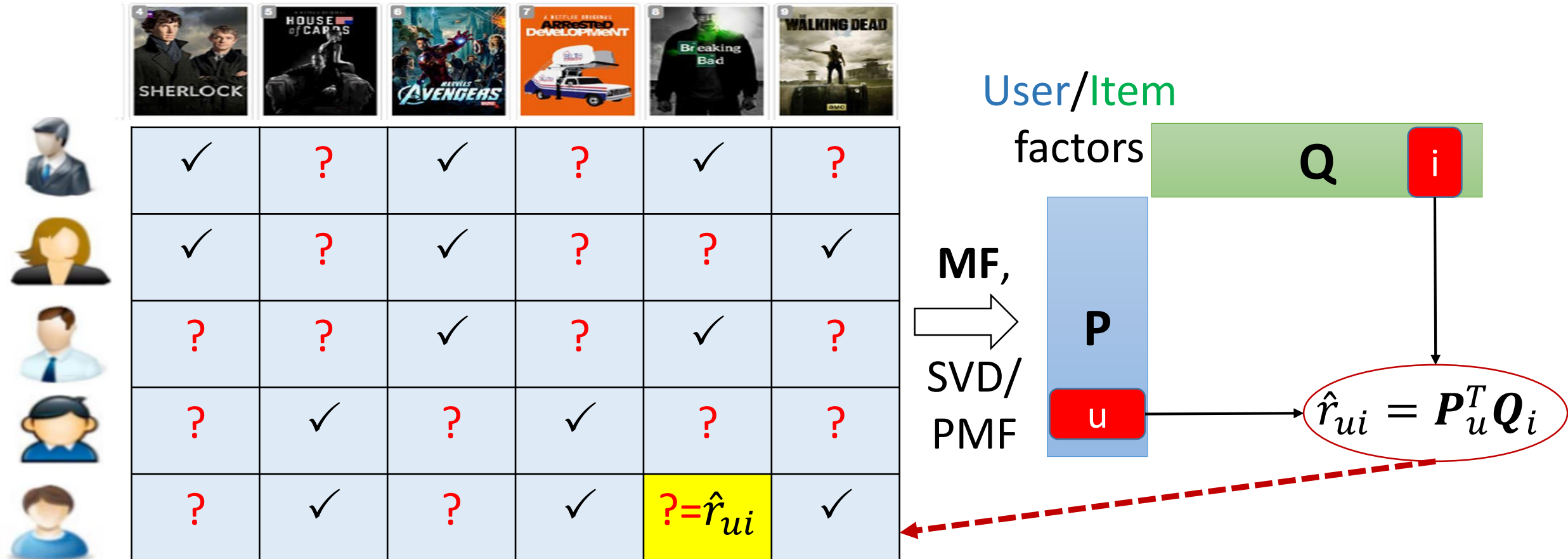
- Mean Reciprocal Rank (MRR)
- Normalized Discounted Cumulative Gain (NDCG)

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

Collaborative filtering

Typical Methods: Matrix Factorization

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

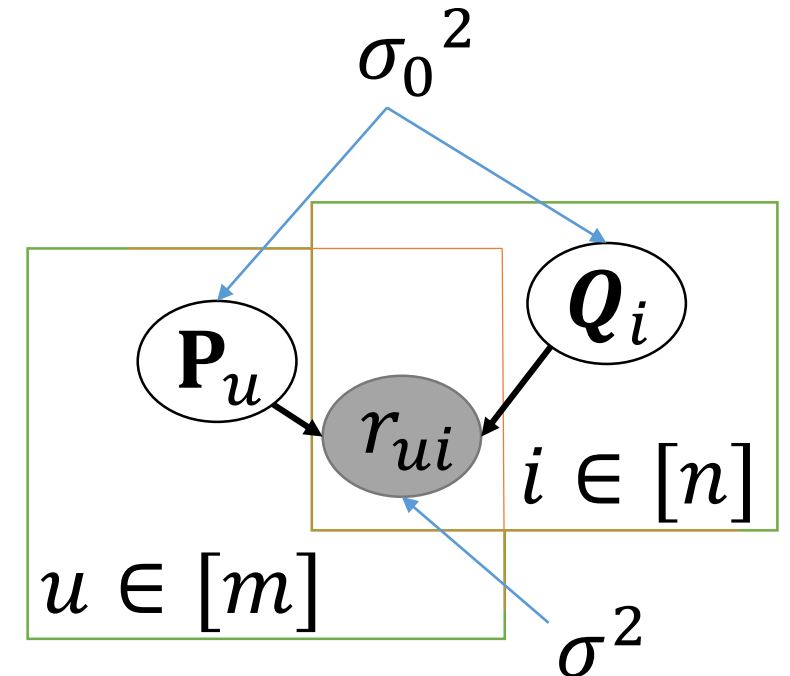


Probabilistic Interpretations: PMF

- 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

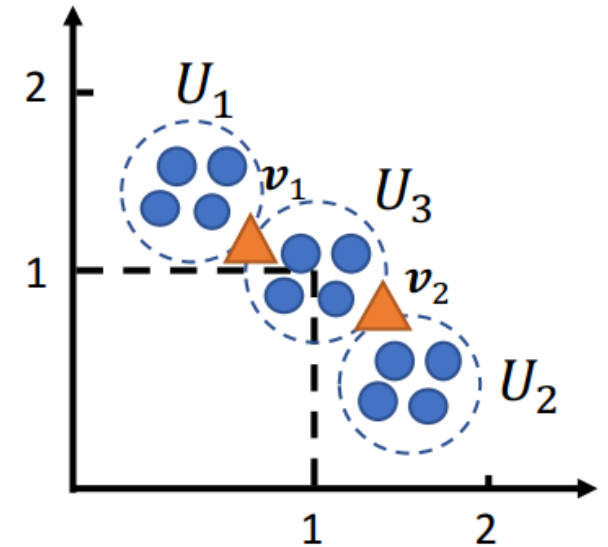
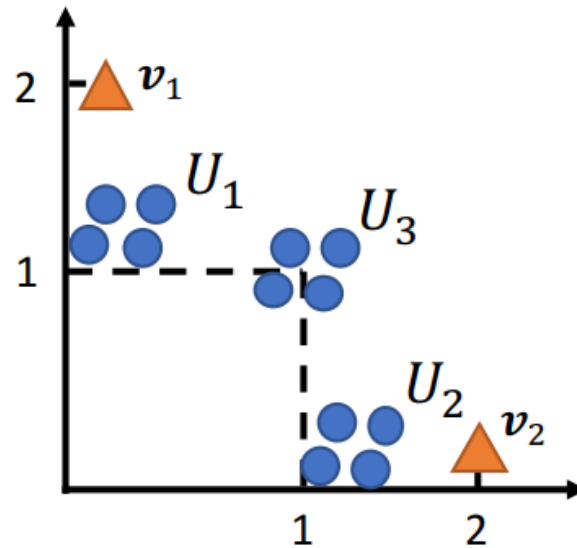


$$\ln p(\Theta | \mathbf{R}, \Phi) = -\frac{1}{2\sigma^2} \sum_{u,i} \delta(r_{ui}) (r_{ui} - \mathbf{P}_u^T \mathbf{Q}_i)^2 - \frac{1}{2\sigma_0^2} (\|\mathbf{P}\|_{Frob}^2 + \|\mathbf{Q}\|_{Frob}^2)$$

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

Limitations of MF: Transitivity

- 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



Metric Learning: Replace Inner Products in MF with (Euclidean) Distances

- An item users liked will be closer to them than other items they did not like

$$d(i, j) = \|\mathbf{u}_i - \mathbf{v}_j\|,$$

- Hinge loss (margin-based)
 - For items user likes, their gradients move inward. For other items, their gradients move outward until they are pushed out by a safe margin

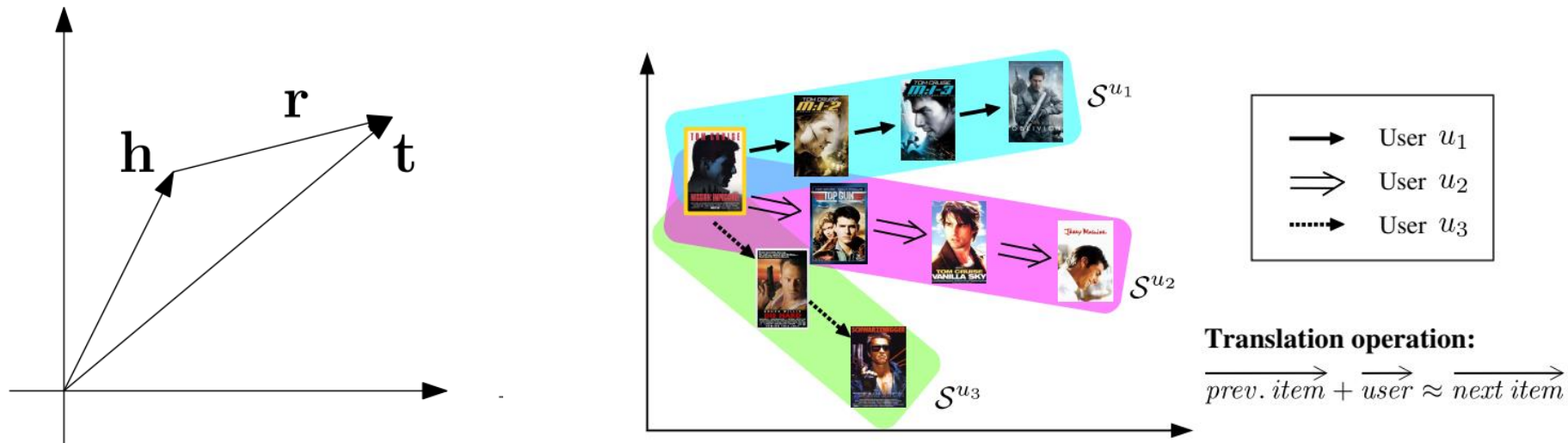
$$\mathcal{L}_m(d) = \sum_{(i,j) \in \mathcal{S}} \sum_{(i,k) \notin \mathcal{S}} w_{ij} [m + d(i, j)^2 - d(i, k)^2]_+,$$

- Rank-based weighting scheme
 - Penalizes a positive item at a lower rank heavily than one at the top

$$w_{ij} = \log(\text{rank}_d(i, j) + 1).$$

Translation-based Recommendation: Capture Sequential Behavior

- Inspired by advances in knowledge graph completion
 - Entities as points and relations as translation vectors
- Items as “entities”, users as “relations” from one item to another



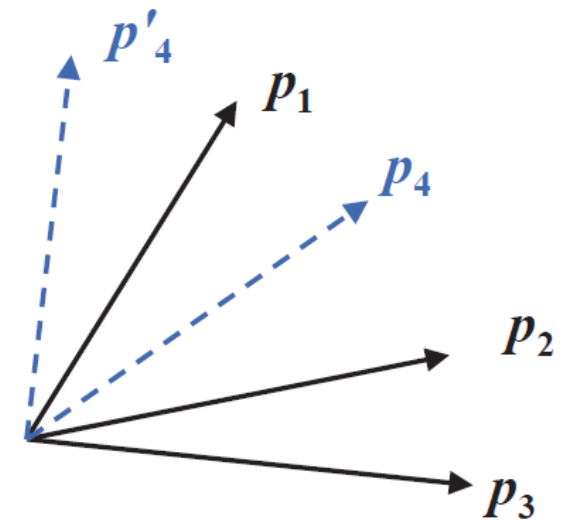
Limited Expressiveness of MF: Nonlinearity

- Similarity of given 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

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

$$\hat{r}_{u,i} = \sigma \left(\mathbf{h}^T (\mathbf{P}_u \odot \mathbf{Q}_i) \right)$$

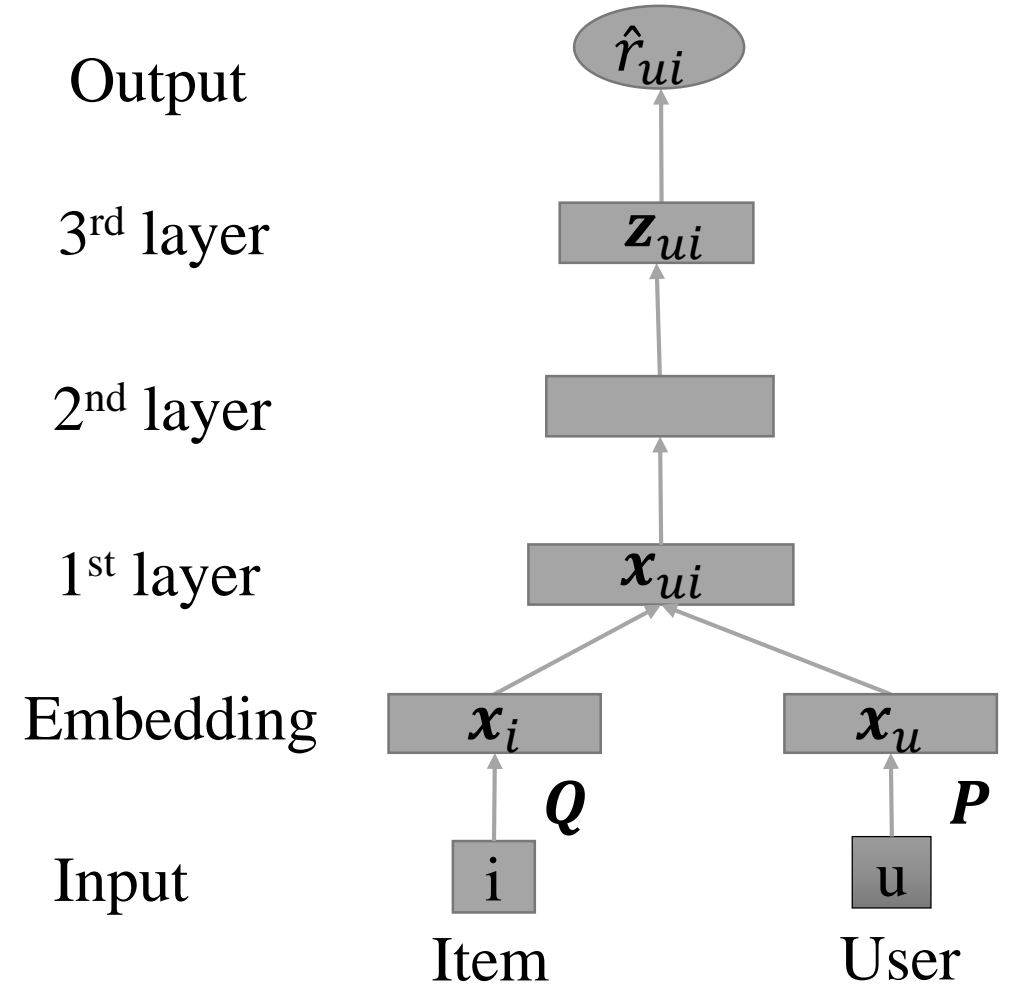
identity activation all-one vector

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

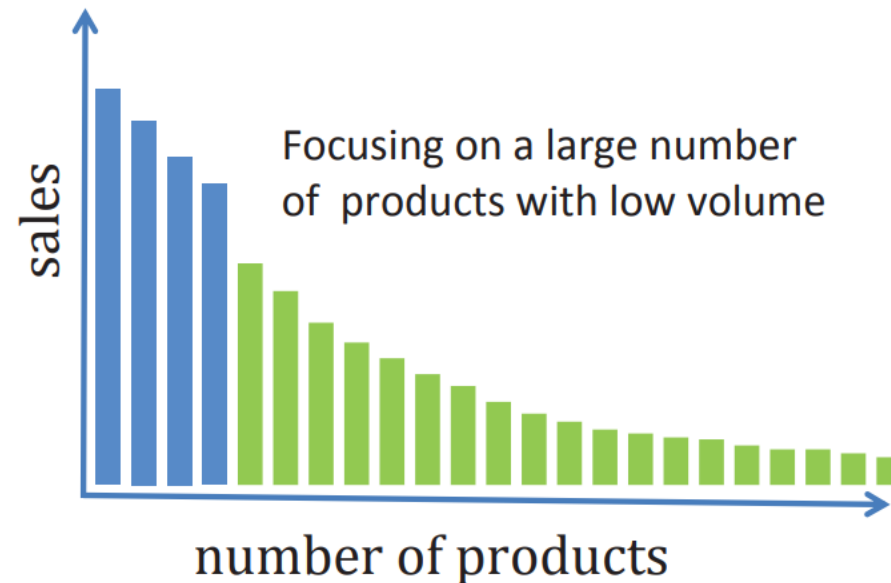


Collaborative Filtering Faces Challenges: Data Sparsity and Long Tail

- 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



Cross-domain recommendation

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 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**

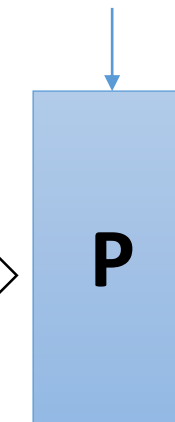
movie	budget	gross	genre	year
Goodfellas	25M	47M	crime	1990
My Cousin Vinny	11M	64M	comedy	1992
...
Clue	15M	15M	comedy	1985

User x Movie

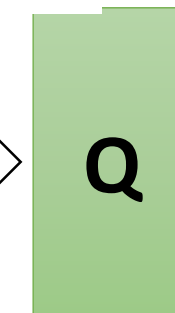
Movie x Genre

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

User factors



Shared item factors



Genre factors

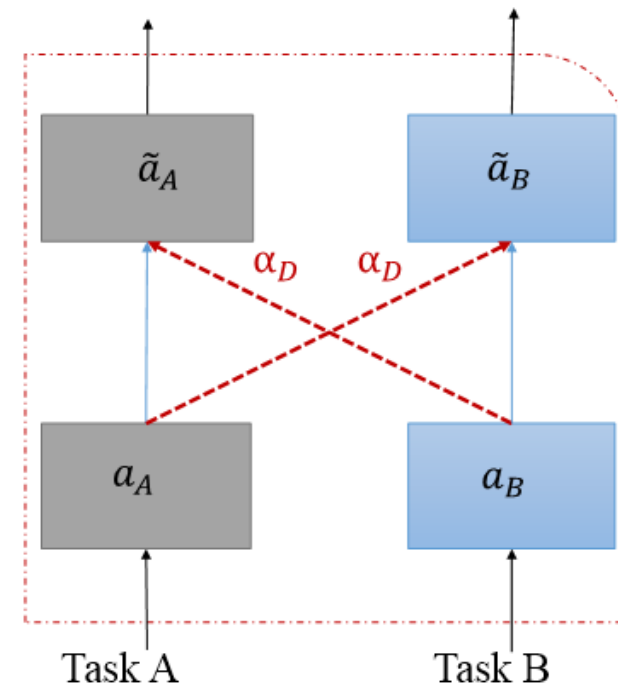


Deep Methods: Cross-Stitch Networks (CSN)

- 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

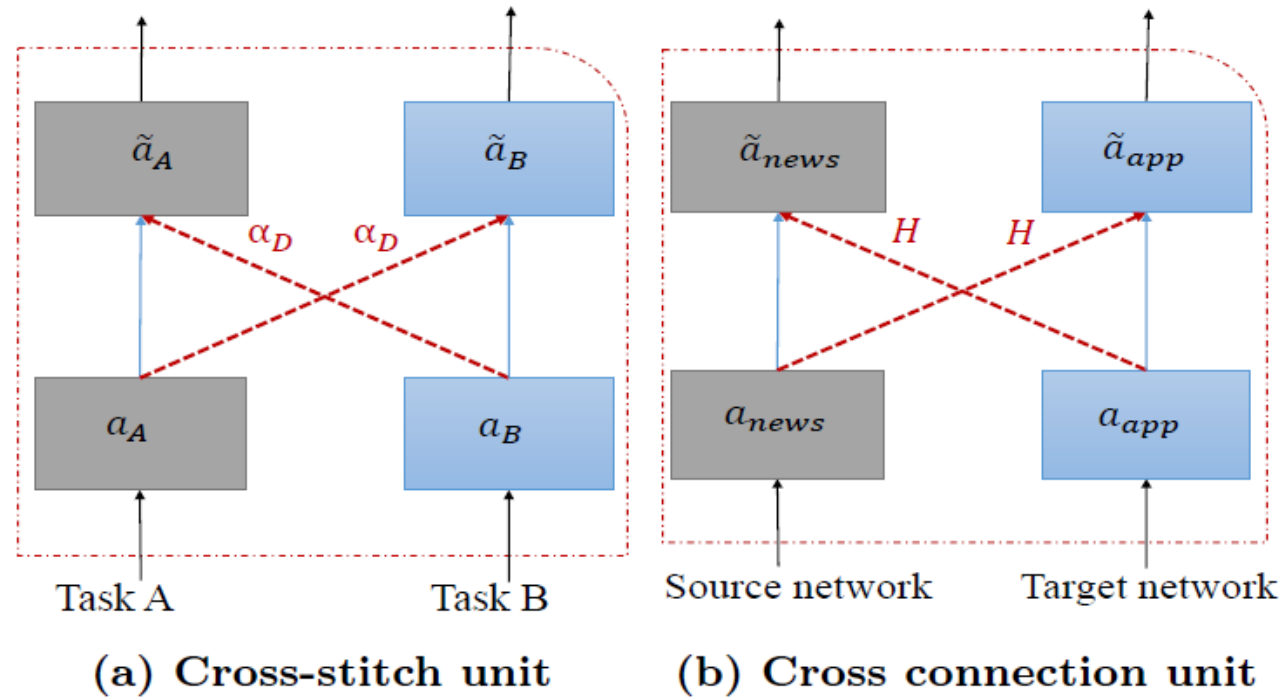


Collaborative Cross Networks (CoNet)

- A novel deep transfer learning method
- Alleviate the data sparsity issue faced by deep collaborative filtering
 - By transferring knowledge from a related source domain
- Relax strong assumptions faced by 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

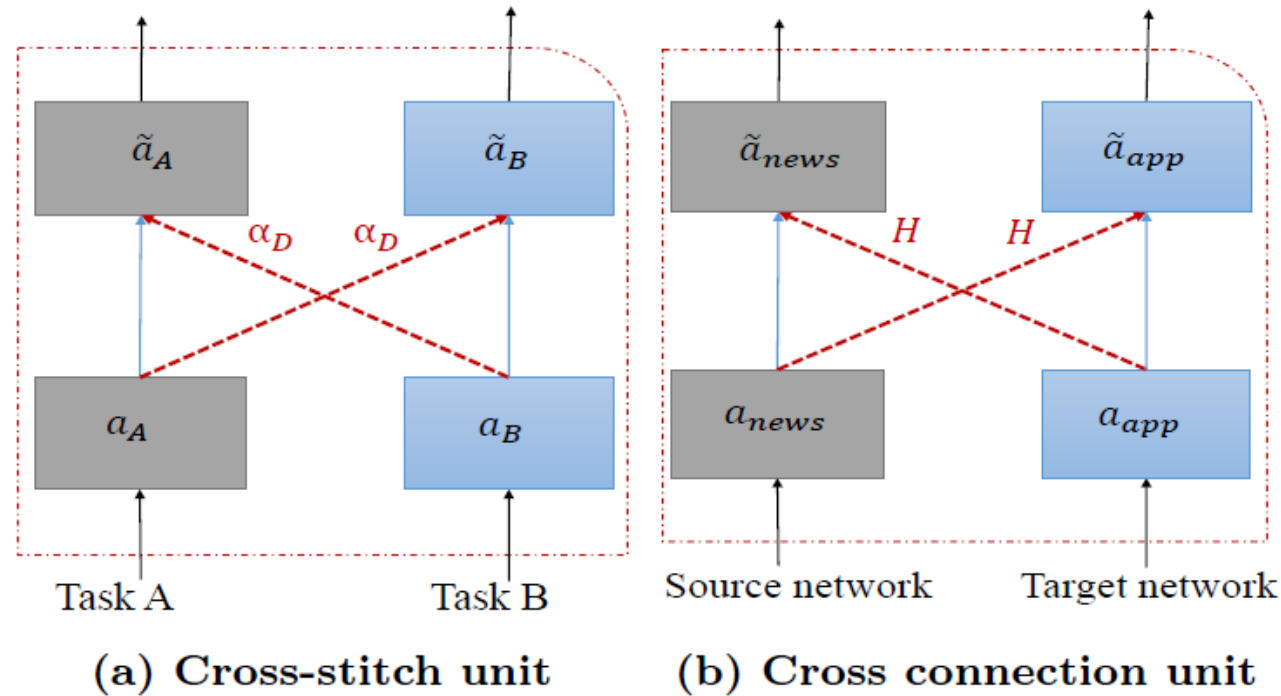
- We can relax the SA 1 assumption (equally important)



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

Idea 2: Selecting representations via sparsity-induced regularization

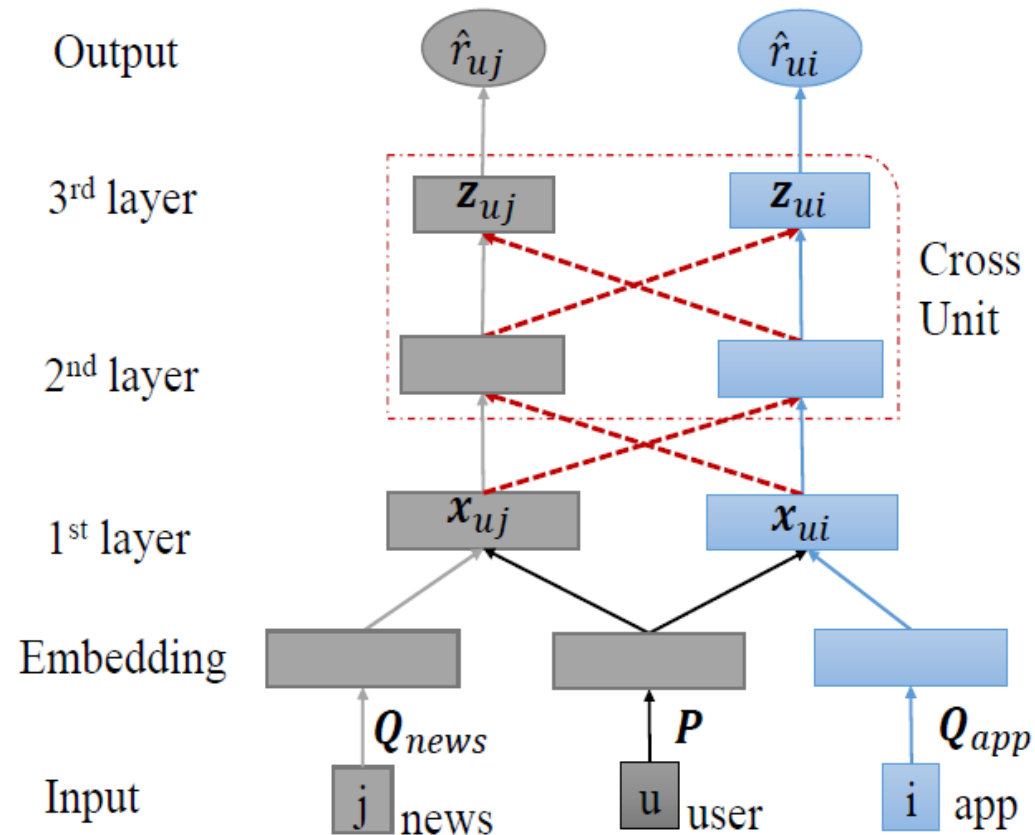
- We can relax the SA 2 assumption (all useful)



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

Architecture of CoNet

- 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 Tensor Flow (<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%

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

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

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

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

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]

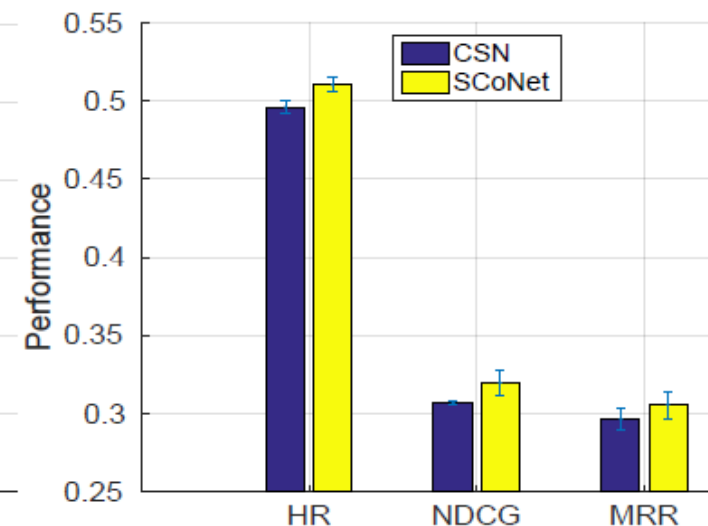
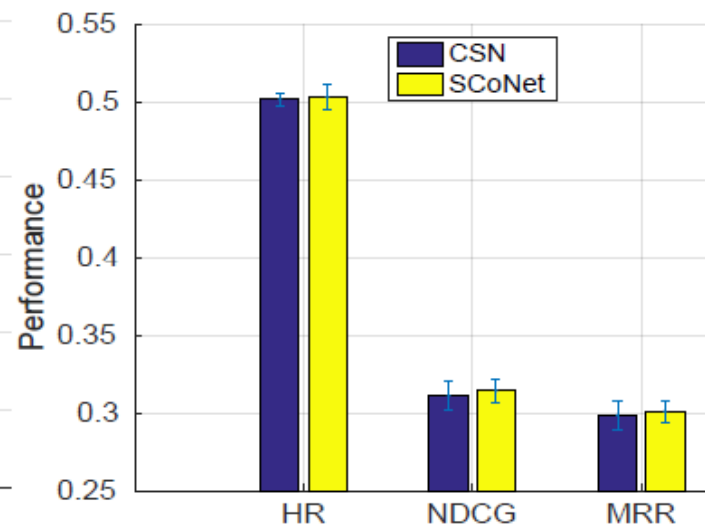
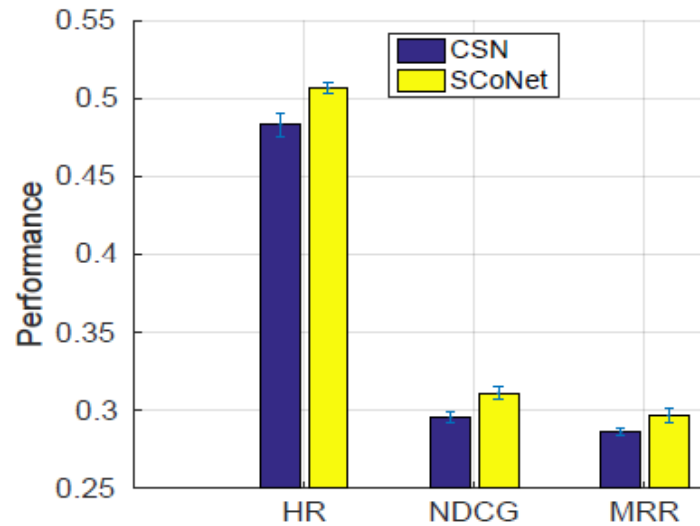
Comparing Different Approaches

- 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



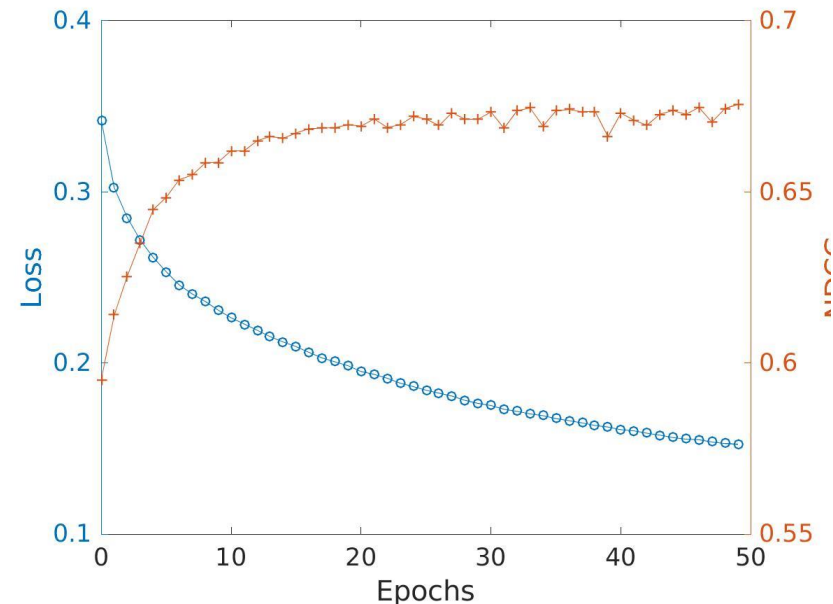
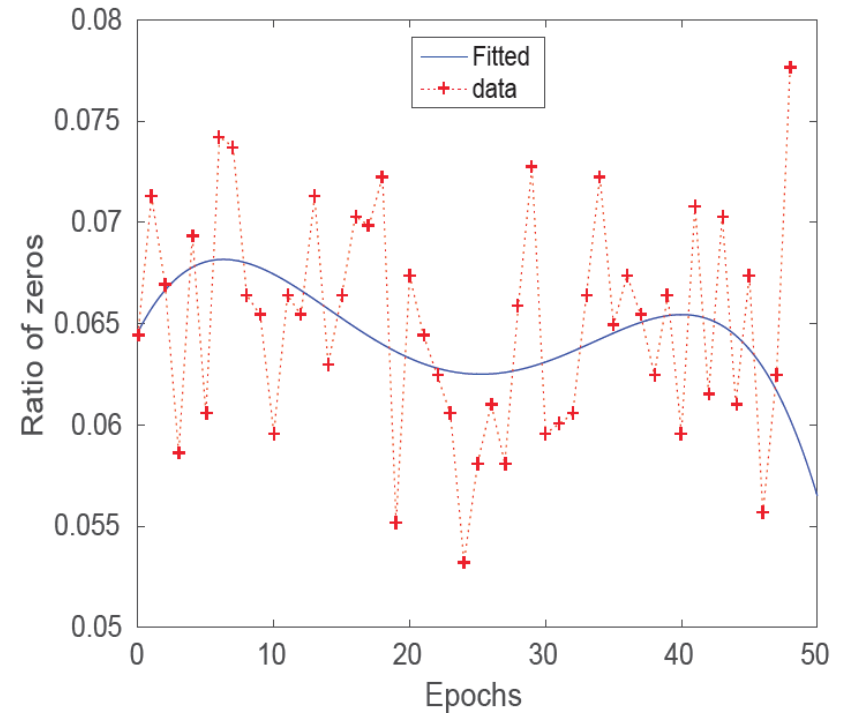
Benefit of Transferring Knowledge

- 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



Summary

- 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,

Hybrid filtering

Another Solution: Hybrid Filtering (Collaborative + Content)

- Item reviews justify ratings
- Item content reveals topic semantics

OliviuNea... ★★★★★ Rating

iPhone 6 16GB - A jump into the best Smartphone available place. 17.11.2014

I am a tech freak, I have owned every iPhone this, but I also owned almost every flagship. I rarely keep smartphones more than 6 months or sell them and put a little extra so I can buy I bought about 3 weeks ago. I used to have the everything about it, it was small and beautiful, had the opportunity to exchange it for an iPhone photos and videos I disliked the design of the bigger phone and hated how I had problems walking, always in need for 2 hands was one

[Add to my Circle of Trust](#)
[Subscribe to reviews](#)

About me: Exams coming up next, sorry for my absence.

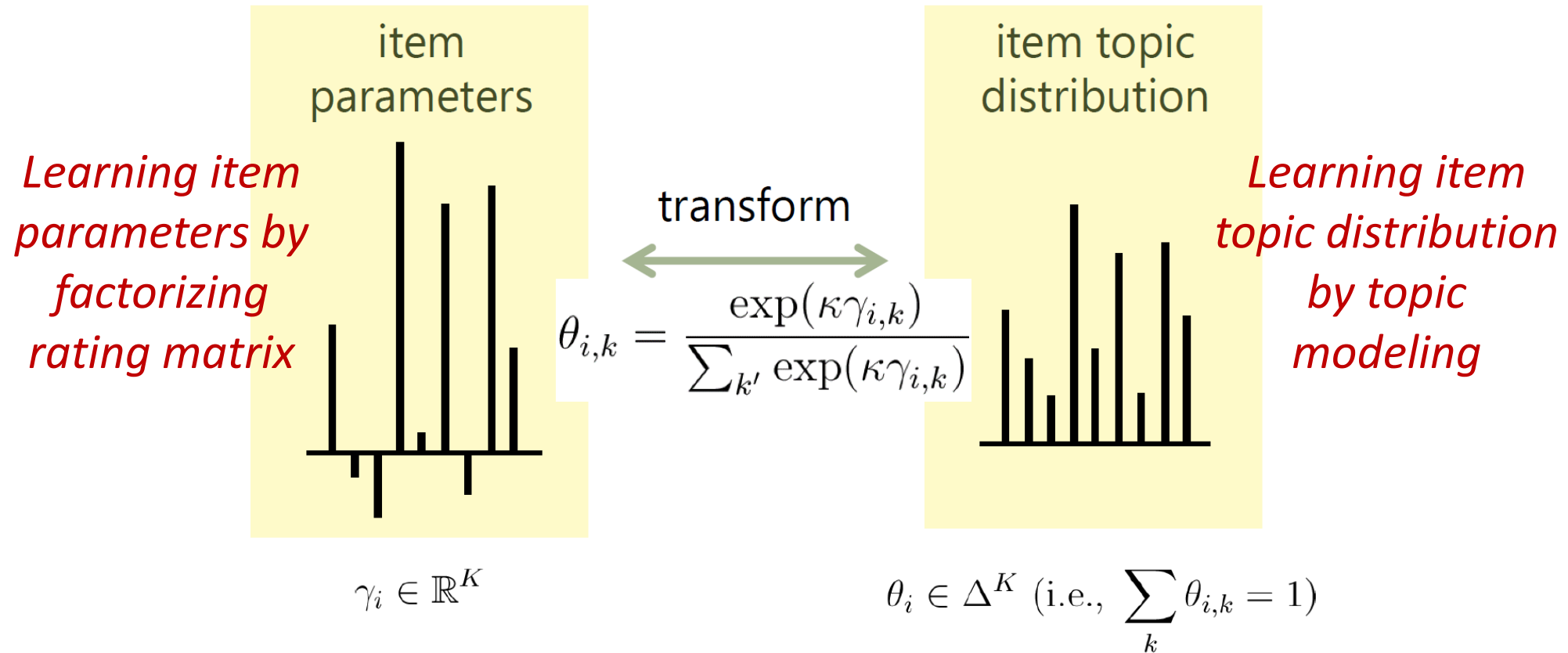
Member since: 12.10.2014
Reviews: 30

already read



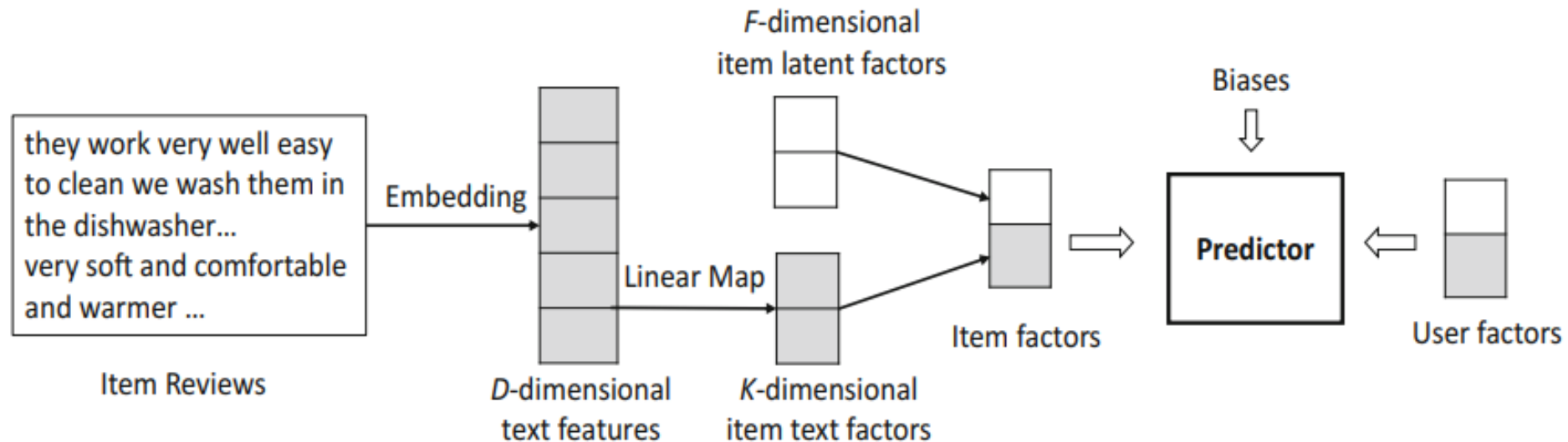
Topic Modelling: Hidden Factors & Topics

- using a transform that aligns latent rating and review terms, so that both are determined by a single parameter



Pre-extracted Word-embedding Features

- Basic MF factorizes ratings into user/item *latent* factors
- Another MF factorizes reviews into user/item *text* factors



$$f_i \equiv \frac{1}{|d_i|} \sum_{w \in d_i} e_w \quad P_u^T Q_i + \theta_u^T (H f_i)$$

Personalized Neural Embeddings (PNE)

- The way of pre-extracted embeddings *separates* the extraction of text features from the learning of user-item interaction
- These two processes cannot benefit from each other and *errors* in the previous step maybe propagate to the successive steps
- PNE learns embeddings of users, items, and words jointly, and predict user preferences on items based on these learned representations
- PNE estimates the probability that a user will like an item by two terms — *behavior* factors and *semantic* factors

Architecture of PNE

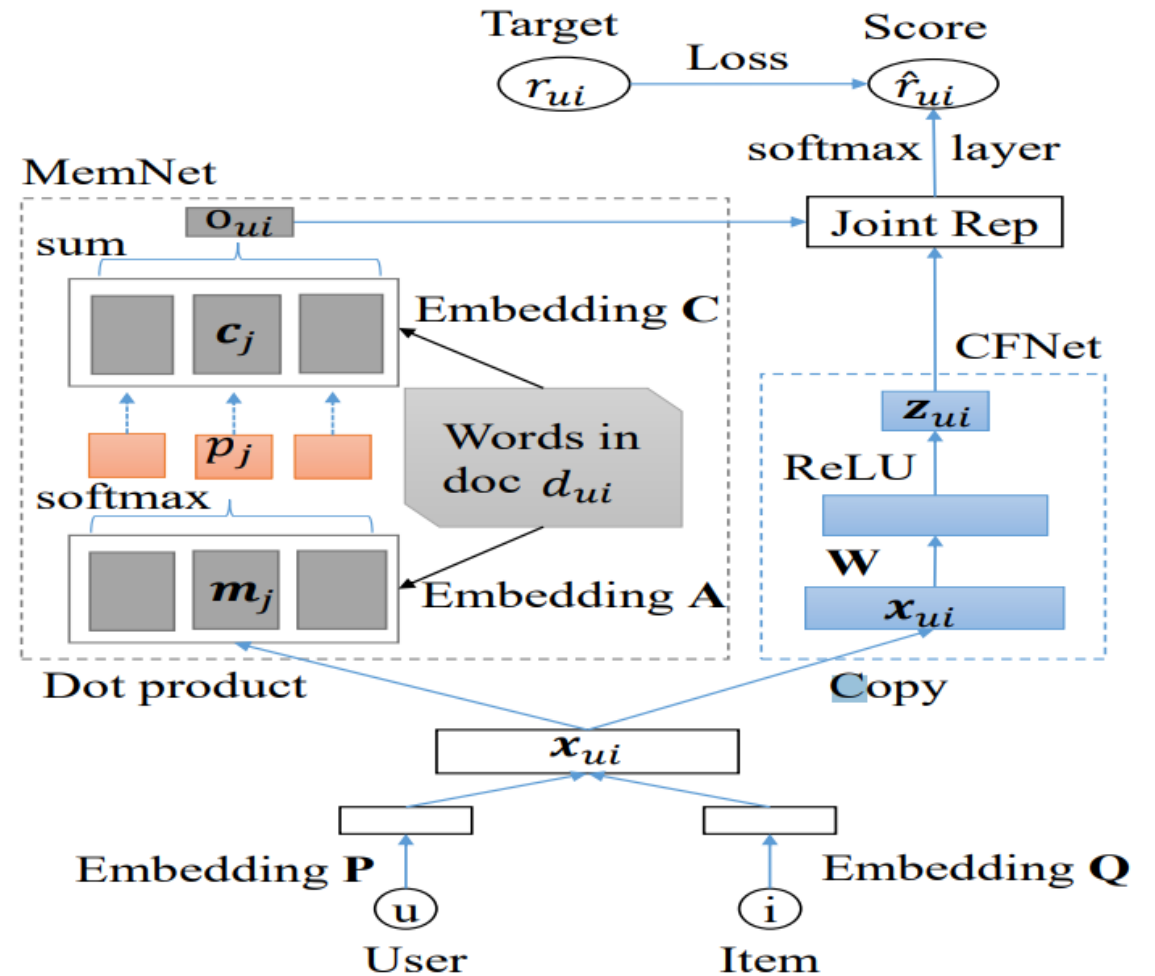
- Behavior factors: same with neural CF

$$z_{ui}^{\text{behavior}} = \text{ReLU}(W x_{ui} + b)$$

- Semantic factors: relevance of a user to a word is learned by attention mechanism

$$z_{ui}^{\text{semantic}} = \sum_{j: w_j \in d_{ui}} \text{Softmax}(a_j^{u,i}) c_j$$

$$a_j^{u,i} = x_{ui}^T m_j^{u,i}$$



Dataset and Baselines

- Datasets

- Amazon reviews
- Cheetah news

Dataset	#user	#item	#rating	#word	#density	avg. words
Amazon	8,514	28,262	56,050	1,845,387	0.023%	65.3
Cheetah	15,890	84,802	477,685	612,839	0.035%	7.2

- Baselines

Baselines	Shallow method	Deep method
CF	BPR	MLP
CF w/ text	HFT, TBPR	LCMR, PNE (ours)

Comparing Different Approaches

- PNE vs MLP: Since CFNet of PNE is a neural CF (with one hidden layer), results show the benefit of exploiting unstructured text to alleviate the data sparsity issue faced by CF methods

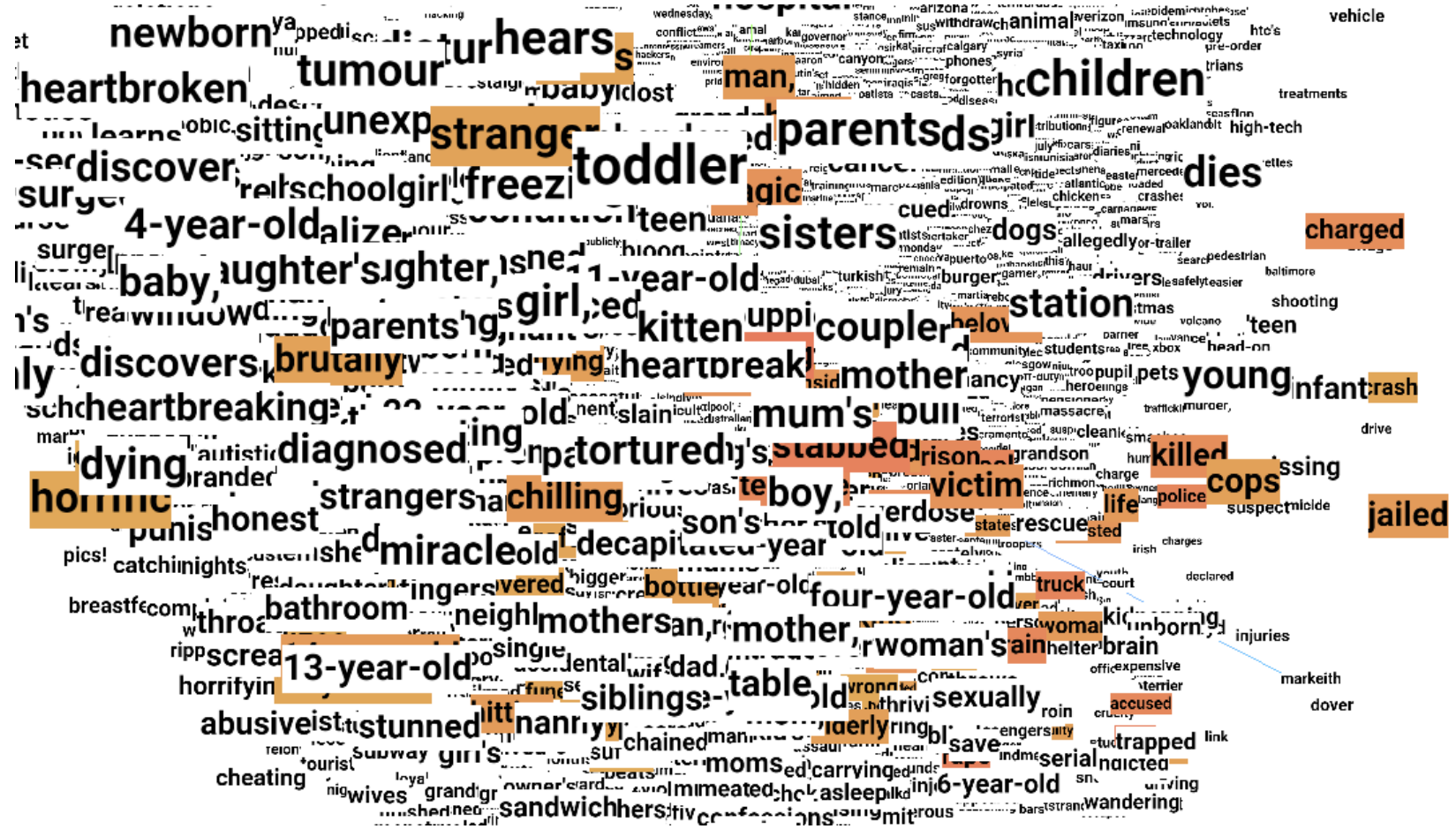
- PNE vs HFT/TBPR: Results show the benefit of integrating content text through MemNet (and exploiting interactions through neural CF)

TopK	Metric	Method					
		BPR	HFT	TBPR	MLP	LCMR	PNE
5	HR	8.10	10.77	15.17	21.00*	20.24	23.52
	NDCG	5.83	8.15	12.08	14.86*	14.51	16.46
	MRR	5.09	7.29	11.04	12.83*	12.63	14.13
10	HR	12.04	13.60	17.77	28.36*	28.36*	31.86
	NDCG	7.10	9.07	12.91	16.97*	16.78	19.15
	MRR	5.61	7.67	11.38	13.71*	13.56	15.24
20	HR	18.21	27.82	22.68	38.20	39.51*	42.21
	NDCG	8.64	12.52	14.14	18.99	19.18*	21.75
	MRR	6.02	8.54	11.71	14.26*	14.20	15.95

- PNE vs LCMR: Since MemNet of PNE is the same with Local MemNet of LCMR (with one-hop), results show the design of CFNet of PNE is more reasonable than that of Centralized MemNet of LCMR

PNE Learns Meaningful Word Embeddings

- Nearest neighbors of drug: *shot, shoots, gang, murder, killing, rape, stabbed, truck, school, police, teenage*
- Google word2vec: *drugs, heroin, addiction, abuse, fda, alcoholism, cocaine, lsd, alcohol, schedule, substances*



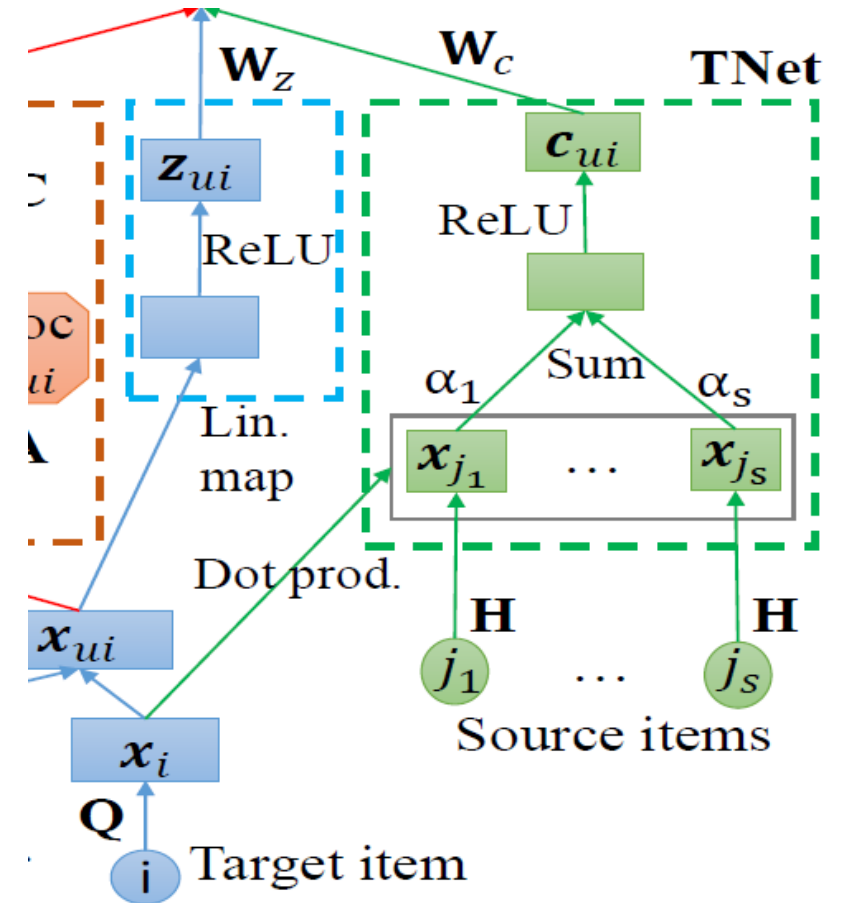
Transfer meets hybrid

Transfer Meets Hybrid: A Synthetic Approach for Cross-Domain Collaborative Filtering with Text

- Hybrid filtering methods integrate content information, e.g. product reviews and news titles
- Cross-domain methods leverage knowledge from a related domain, e.g. from Apps to News
- TMH attentively extracts useful content from unstructured text via a memory network and ...
- ... selectively transfers knowledge from a source domain via a transfer network

Architecture of TMH

- A MemNet: Matching Word Semantics with User Preferences
 - Same with MemNet of PNE and Local MemNet of LCMR
- A TransNet: Selecting Source Items to Transfer by a way of coarse-to-fine
 - *Coarse*: transfer source items such that this user has interacted in source domain
 - *Fine*: similarities between target item and coarse source items by content-based addressing
 - Finally: transfer vector is a weighted sum of the corresponding source item embeddings



$$\mathbf{c}_{ui} = \text{ReLU}\left(\sum_j \alpha_j^{(i)} \mathbf{x}_j\right)$$

Datasets

Dataset	Domain	Statistics	Amount
Mobile News	Shared	#Users	15,890
	Target	#News	84,802
		#Reads	477,685
		Density	0.035%
		#Words	612,839
		Avg. Words Per News	7.2
	Source	#Apps	14,340
#Installations		817,120	
Density		0.359%	
Amazon Product	Shared	#Users	8,514
	Target	#Clothes (Men)	28,262
		#Ratings/#Reviews	56,050
		Density	0.023%
		#Words	1,845,387
		Avg. Words Per Review	32.9
	Source	#Products (Sports)	41,317
#Ratings/#Reviews		81,924	
Density		0.023%	

Baselines

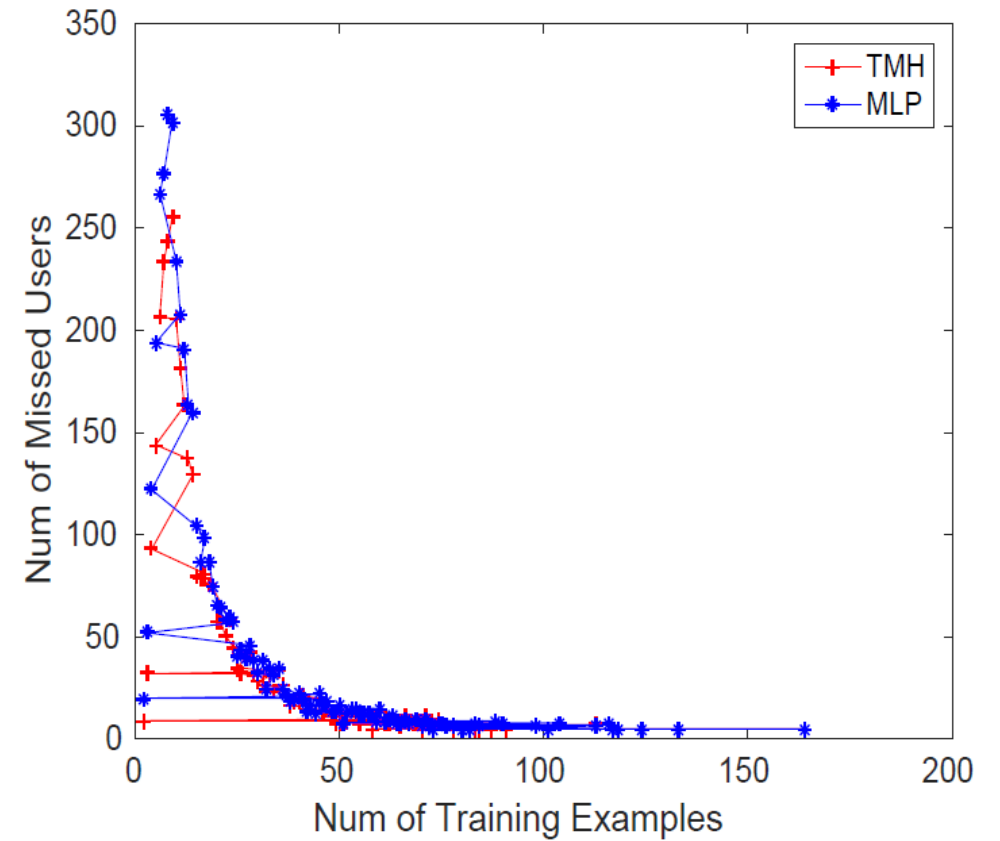
Baselines	Shallow method	Deep method
Single-domain	BPRMF [41]	MLP [17]
Cross-domain	CDCF [31], CMF [42]	MLP++, CSN [34]
Hybrid	HFT [33], TextBPR [16, 20]	LCMR [19]
Cross + Hybrid	CDCF++	TMH (ours)

Results on Amazon Dataset

Method	<i>topK = 5</i>			<i>topK = 10</i>			<i>topK = 20</i>		
	HR	NDCG	MRR	HR	NDCG	MRR	HR	NDCG	MRR
BPRMF	.0810	.0583	.0509	.1204	.0710	.0561	.1821	.0864	.0602
CDCF	.1295	.0920	.0797	.2070	.1167	.0897	.3841	.1609	.1015
CMF	.1498	.0950	.0771	.2224	.1182	.0863	.3573	.1521	.0957
HFT	.1077	.0815	.0729	.1360	.0907	.0767	.2782	.1252	.0854
TextBPR	.1517	.1208	.1104	.1777	.1291	.1138	.2268	.1414	.1171
CDCF++	.1314	.0926	.0800	.2102	.1177	.0901	.3822	.1605	.1016
MLP	.2100	.1486	.1283	.2836	.1697	.1371	.3820	.1899	.1426
MLP++	.2263	.1626	.1417	.2992	.1862	.1514	.3810	.2069	.1570
CSN	.2340*	.1680*	.1462*	.3018*	.1898*	.1552*	.3944*	.2091*	.1605*
LCMR	.2024	.1451	.1263	.2836	.1678	.1356	.3951	.1918	.1420
TMH	.2575	.1796	.1550	.3490	.2077	.1666	.4443	.2311	.1727
Our improve	10.04%	6.90%	6.01%	15.63%	9.43%	7.34%	12.65%	10.52%	7.60%

Improvement on Cold Users (and Items)

- Missed Hit Users (MHU) distribution on Cheetah Mobile
- We expect that cold users in MHUs can be reduced by using TMH
- The more amount we can reduce, the more effective that TMH can alleviate the cold-user start issues
- MHUs are most of cold users who have few training examples.
- #cold-users in MHUs of MLP is higher than that of TMH.
- TMH reduces #cold-users from 1,385 to 1,145 on Mobile, achieving relative 20.9% reduction



Future works

- Data privacy
 - Source domain cannot share the raw data, but model parameters

Thanks!

Q & A