



CSCE 670 - Information Storage and Retrieval

Week 11: Deep Learning for Recommendation

Yu Zhang

yuzhang@tamu.edu

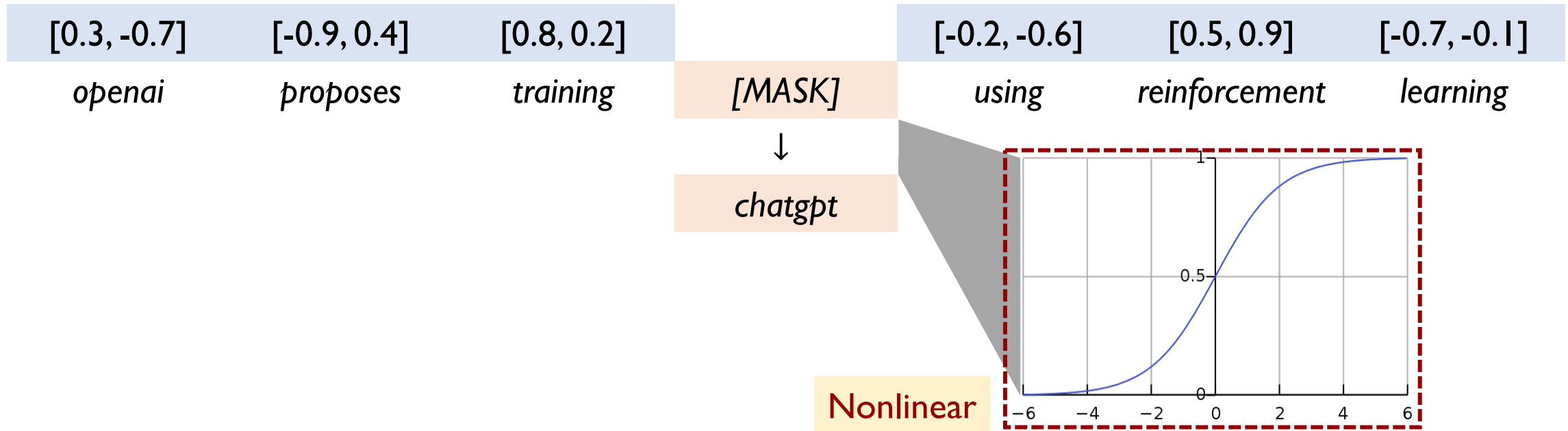
Course Website: <https://yuzhang-teaching.github.io/CSCE670-S26.html>

Quiz 3!

- All policies are the same as Quiz 1 (number of questions, time limit, grading, etc.)
- Scope:
 - Week 8 (word2vec, Neural Ranking)
 - Week 10 (BERT, BERT-Based Ranking)
 - Homework 2
 - **Neural Collaborative Filtering / Sequential Recommendation NOT included**

Neural Collaborative Filtering

Neural Word Embeddings (word2vec)



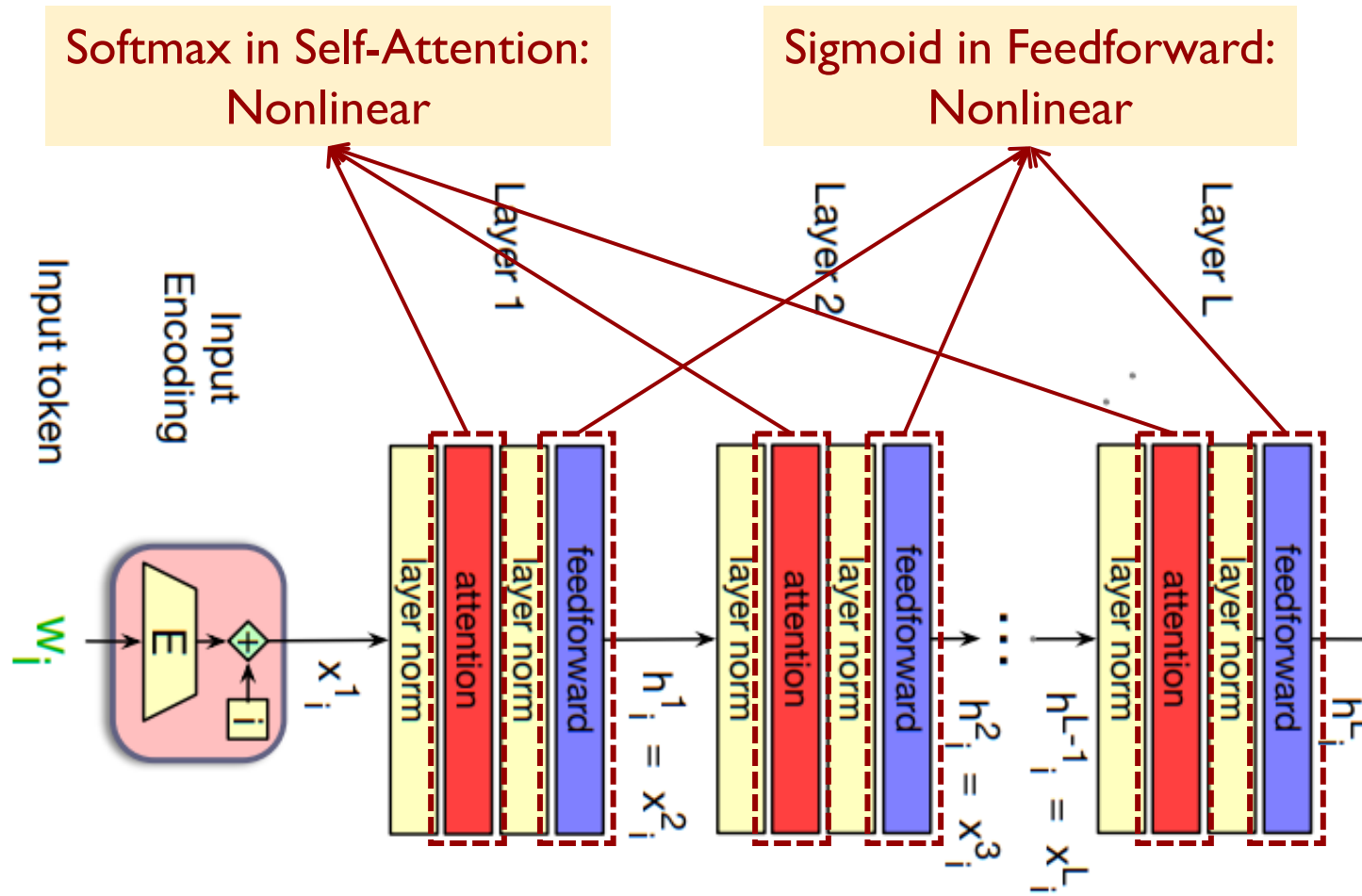
- [Levy and Golberg, NIPS 2014] word2vec is mathematically equivalent to factorizing a word-word matrix U , where

Nonlinear

$$U_{xy} = \log \frac{\#(x, y) \cdot |\mathcal{D}|}{\#(x) \cdot \#(y) \cdot b}$$

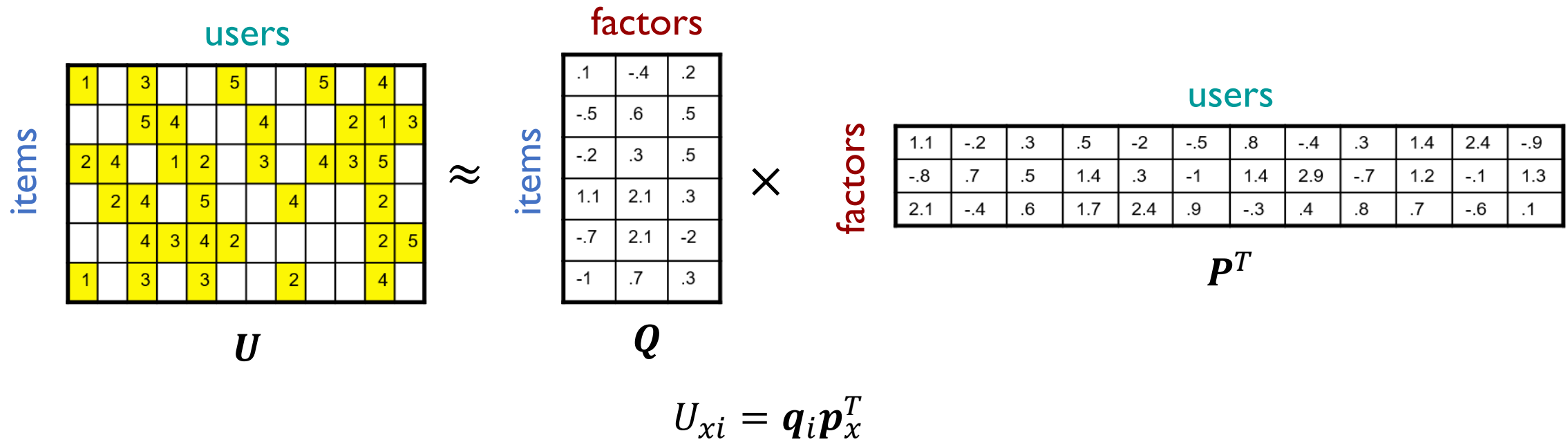
Neural Language Models (BERT)

- Transformer is a neural network



Nonlinearity in Recommender Systems?

- Do the recommender system techniques we have learned so far contain enough **nonlinearity**?
 - No, matrix factorization is essentially modeling the inner product of vectors, which is a highly linear operation.



Neural Collaborative Filtering [He et al., WWW 2017]

Neural Collaborative Filtering*

Xiangnan He
National University of
Singapore, Singapore
xiangnanhe@gmail.com

Lizi Liao
National University of
Singapore, Singapore
liaolizi.llz@gmail.com

Hanwang Zhang
Columbia University
USA
hanwangzhang@gmail.com

Liqiang Nie
Shandong University
China
nieliqiang@gmail.com

Xia Hu
Texas A&M University
USA
hu@cse.tamu.edu

Tat-Seng Chua
National University of
Singapore, Singapore
dcscts@nus.edu.sg

ABSTRACT

In recent years, deep neural networks have yielded immense success on speech recognition, computer vision and natural language processing. However, the exploration of deep neural networks on recommender systems has received relatively less scrutiny. In this work, we strive to develop techniques based on neural networks to tackle the key problem in recommendation — collaborative filtering — on the basis of implicit feedback.

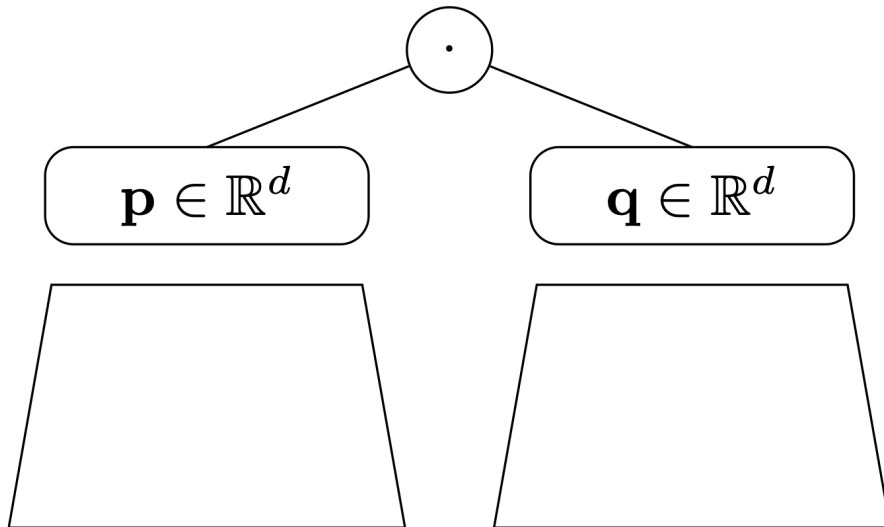
1. INTRODUCTION

In the era of information explosion, recommender systems play a pivotal role in alleviating information overload, having been widely adopted by many online services, including E-commerce, online news and social media sites. The key to a personalized recommender system is in modelling users' preference on items based on their past interactions (*e.g.*, ratings and clicks), known as collaborative filtering [31, 46]. Among the various collaborative filtering techniques, matrix

Learning the Similarity Function

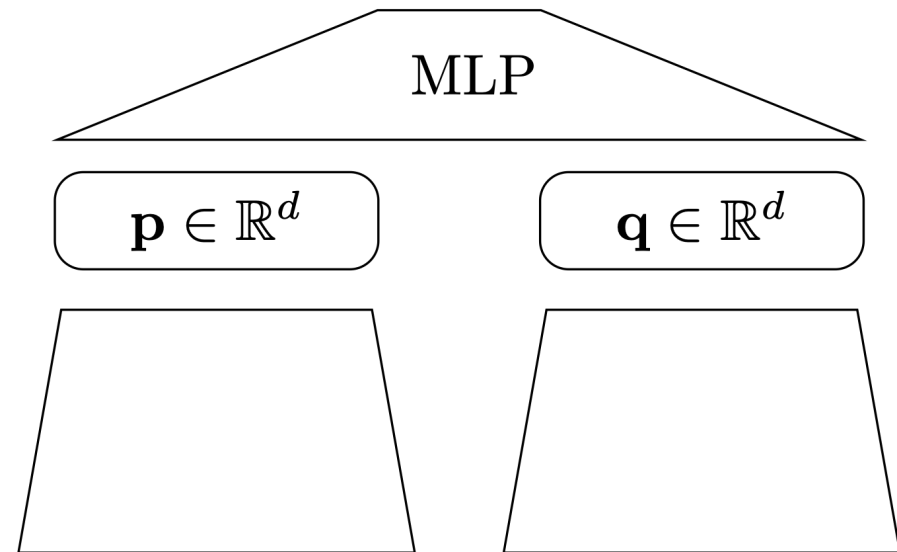
- **Predefined** similarity: inner product

$$\phi^{\text{dot}}(\mathbf{p}, \mathbf{q}) = \langle \mathbf{p}, \mathbf{q} \rangle$$



- **Learned** similarity: multi-layer perceptron (MLP)

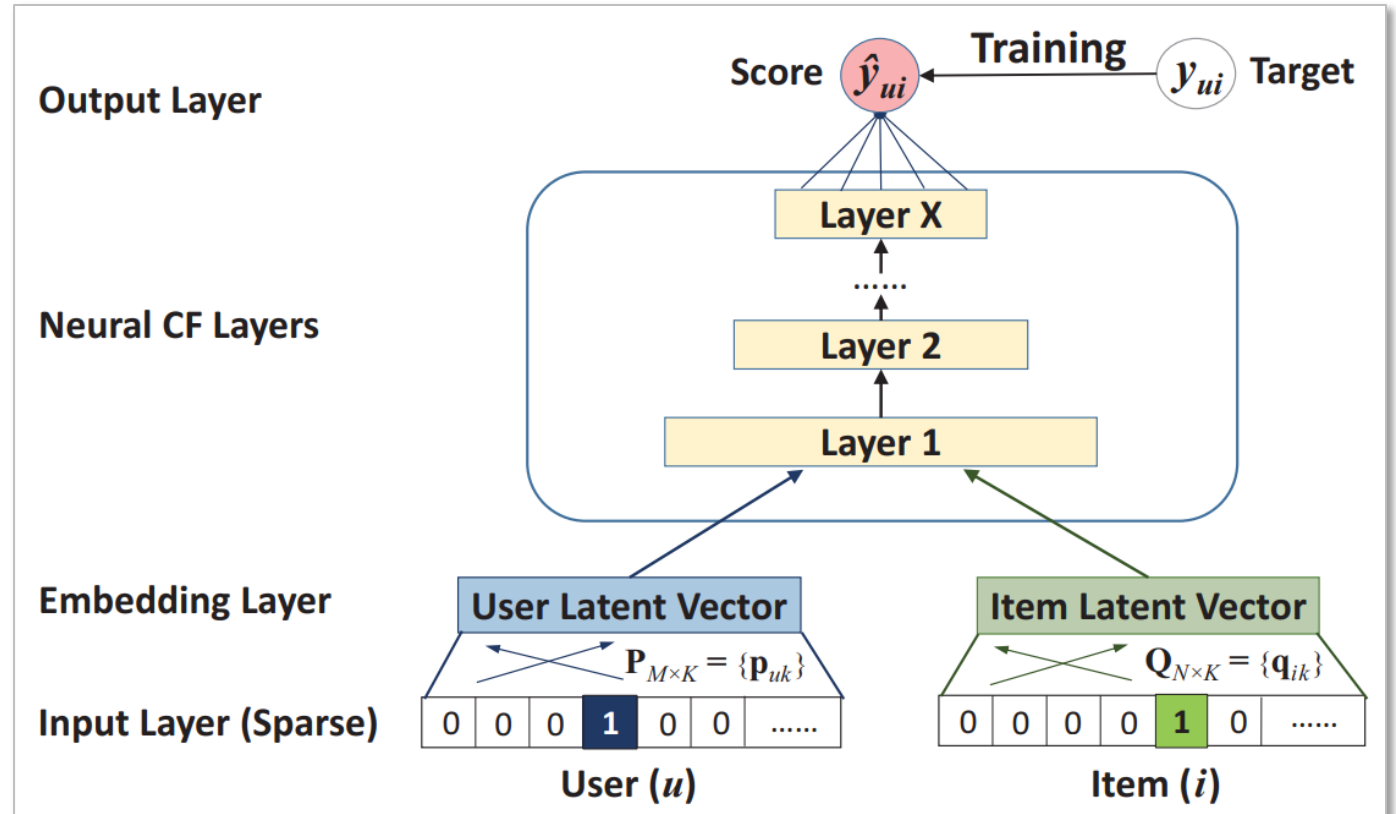
$$\phi^{\text{MLP}}(\mathbf{p}, \mathbf{q}) = \mathbf{f}_{W_l, \mathbf{b}_l}(\dots \mathbf{f}_{W_1, \mathbf{b}_1}([\mathbf{p}, \mathbf{q}]) \dots)$$



- We could use any sort of embeddings down here (e.g., from SVD or Matrix Factorization)

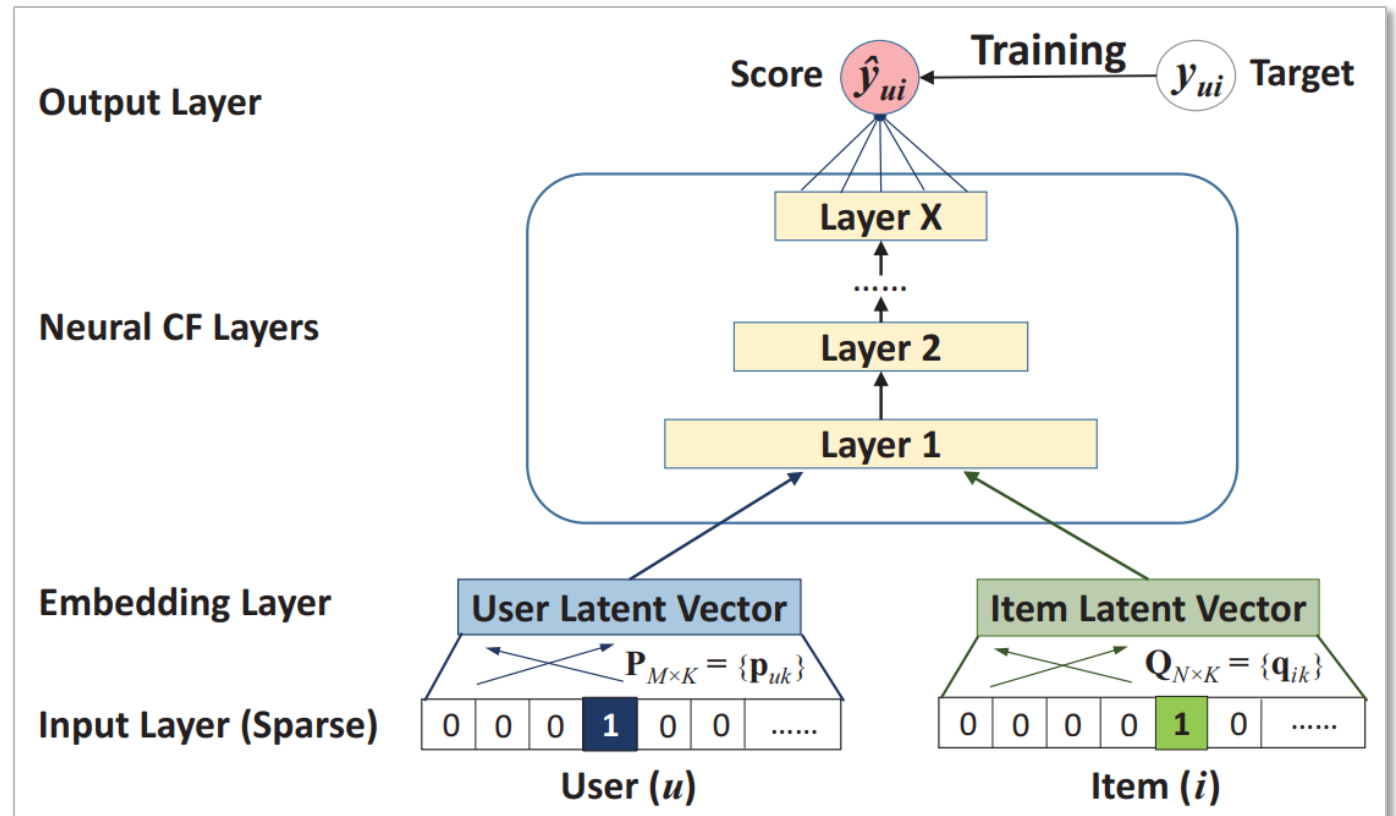
Example

- Suppose after matrix factorization,
 - User u 's vector: $\mathbf{p} = [2, 2]$
 - Item i 's vector: $\mathbf{q} = [2, 3]$
- Inner product: $\text{score}(u, i) = \mathbf{p}\mathbf{q}^T = 10$
- MLP (assuming 3 layers):
 - Input $\mathbf{x}_0 = [\mathbf{p}, \mathbf{q}] = [2, 2, 2, 3]$
 - $\mathbf{x}_1 = \text{Sigmoid}(\mathbf{W}_1\mathbf{x}_0 + \mathbf{b}_1)$
 - $\mathbf{x}_2 = \text{Sigmoid}(\mathbf{W}_2\mathbf{x}_1 + \mathbf{b}_2)$
 - $\hat{y} = \text{score}(u, i) = \text{Sigmoid}(\mathbf{W}_3\mathbf{x}_2)$



Example

- **MLP** (assuming 3 layers):
 - Input $\mathbf{x}_0 = [\mathbf{p}, \mathbf{q}] = [2, 2, 2, 3]$
 - $\mathbf{x}_1 = \text{Sigmoid}(\mathbf{W}_1 \mathbf{x}_0 + \mathbf{b}_1)$
 - $\mathbf{x}_2 = \text{Sigmoid}(\mathbf{W}_2 \mathbf{x}_1 + \mathbf{b}_2)$
 - $\hat{y} = \text{Sigmoid}(\mathbf{W}_3 \mathbf{x}_2)$
- Parameters?
 - $\mathbf{W}_1, \mathbf{b}_1, \mathbf{W}_2, \mathbf{b}_2, \mathbf{W}_3$
 - \mathbf{P}, \mathbf{Q}
- Learning objective?
 - $\text{RMSE} = (y - \hat{y})^2$
- Training data?
 - Known ratings in the user-item matrix



Combining Matrix Factorization and MLP

- Inner product: $\text{score}(u, i) = \mathbf{p}\mathbf{q}^T = 10$

- MLP (assuming 3 layers):

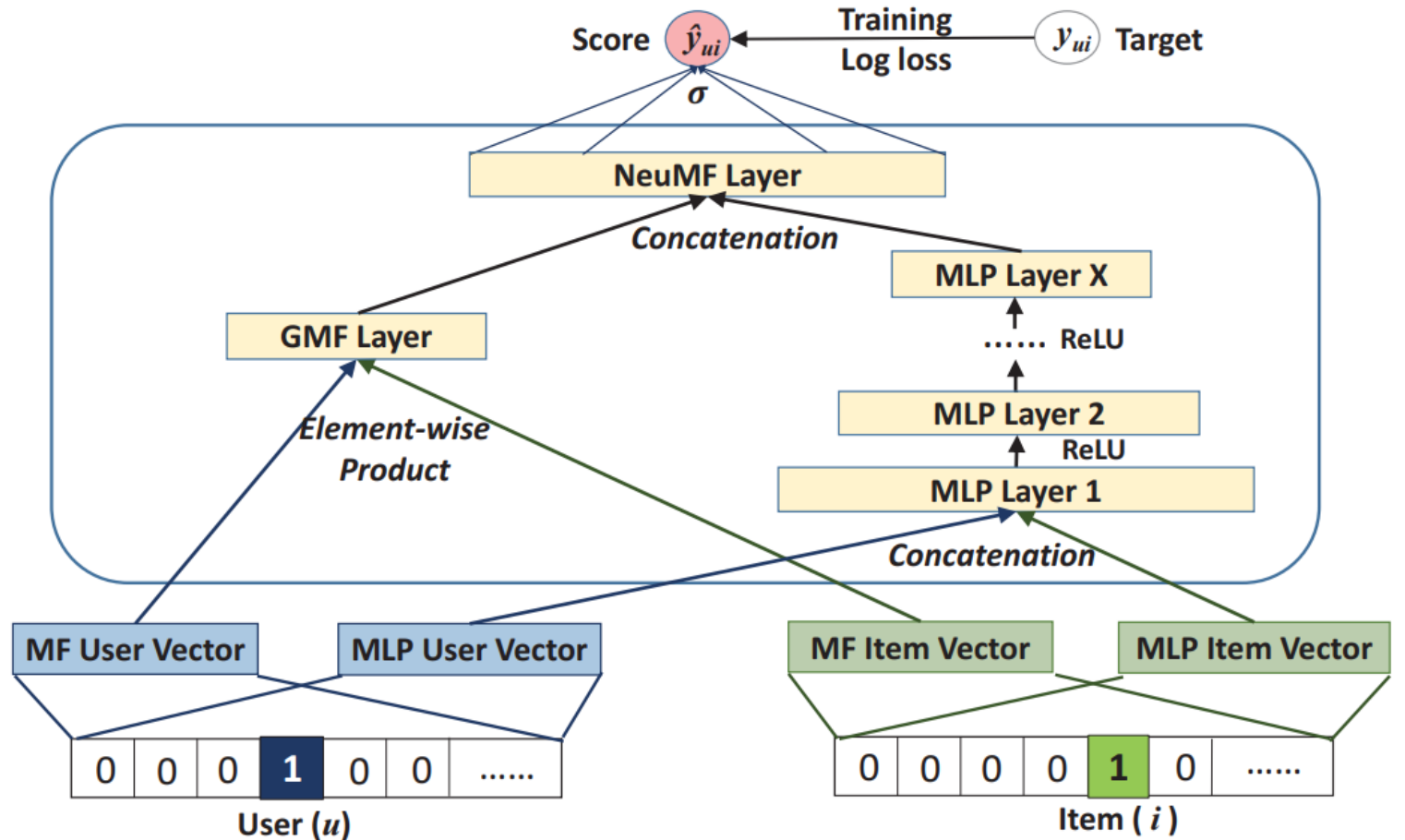
- Input $\mathbf{x}_0 = [\mathbf{p}, \mathbf{q}] = [2, 2, 2, 3]$
- $\mathbf{x}_1 = \text{Sigmoid}(\mathbf{W}_1\mathbf{x}_0 + \mathbf{b}_1)$
- $\mathbf{x}_2 = \text{Sigmoid}(\mathbf{W}_2\mathbf{x}_1 + \mathbf{b}_2)$
- $\text{score}(u, i) = \text{Sigmoid}(\mathbf{W}_3\mathbf{x}_2)$

- MLP does not explicitly model the interaction between \mathbf{p} and \mathbf{q} , but inner product does!
- Can we fuse inner product and MLP?

$$\text{score}(u, i) = \text{score}_{\text{MF}}(u, i) + \text{score}_{\text{MLP}}(u, i)$$

Combining Matrix Factorization and MLP

- User u 's vector: $\mathbf{p} = [2, 2]$
- Item i 's vector: $\mathbf{q} = [2, 3]$
- MLP: same
- Generalized Matrix Factorization (GMF):
 - **Step 1**: Element-wise product $\mathbf{p} \odot \mathbf{q} = [4, 6]$
 - **Step 2**: Nonlinear layer
 $\text{score}_{\text{MF}}(u, i) = \text{Sigmoid}(\mathbf{w} \begin{bmatrix} 4 \\ 6 \end{bmatrix})$
 - \mathbf{w} are the parameters to be learned



Generalized Matrix Factorization vs. Inner Product

- User u 's vector: $\mathbf{p} = [2, 2]$
- Item i 's vector: $\mathbf{q} = [2, 3]$
- Generalized Matrix Factorization (GMF):
 - Step 1: Element-wise product $\mathbf{p} \odot \mathbf{q} = [4, 6]$
 - Step 2: Nonlinear layer score_{MF}(u, i) = Sigmoid($\mathbf{w} \begin{bmatrix} 4 \\ 6 \end{bmatrix}$)
- Suppose $\mathbf{w} = [1, -1]$, then

$$\begin{aligned} \mathbf{w}(\mathbf{p} \odot \mathbf{q}) &= [1, -1] \begin{bmatrix} 2 \times 2 \\ 2 \times 3 \end{bmatrix} = 1 \times 2 \times 2 + (-1) \times 2 \times 3 \\ &= [2, 2] \begin{bmatrix} 1 & 0 \\ 0 & -1 \end{bmatrix} \begin{bmatrix} 2 \\ 3 \end{bmatrix} = \mathbf{p} \text{diag}\{\mathbf{w}\} \mathbf{q}^T \end{aligned}$$

$$\text{score}_{\text{MF}}(u, i) = \text{Sigmoid}(\mathbf{p} \text{diag}\{\mathbf{w}\} \mathbf{q}^T)$$

Nonlinear

Learnable Parameters

Summary

Generalized Matrix Factorization:

Input $x_0 = [p, q]$

$$\begin{aligned} \text{score}_{\text{MF}}(u, i) &= \text{Sigmoid}(w(p \odot q)) \\ &= \text{Sigmoid}(p \text{ diag}\{w\} q^T) \end{aligned}$$

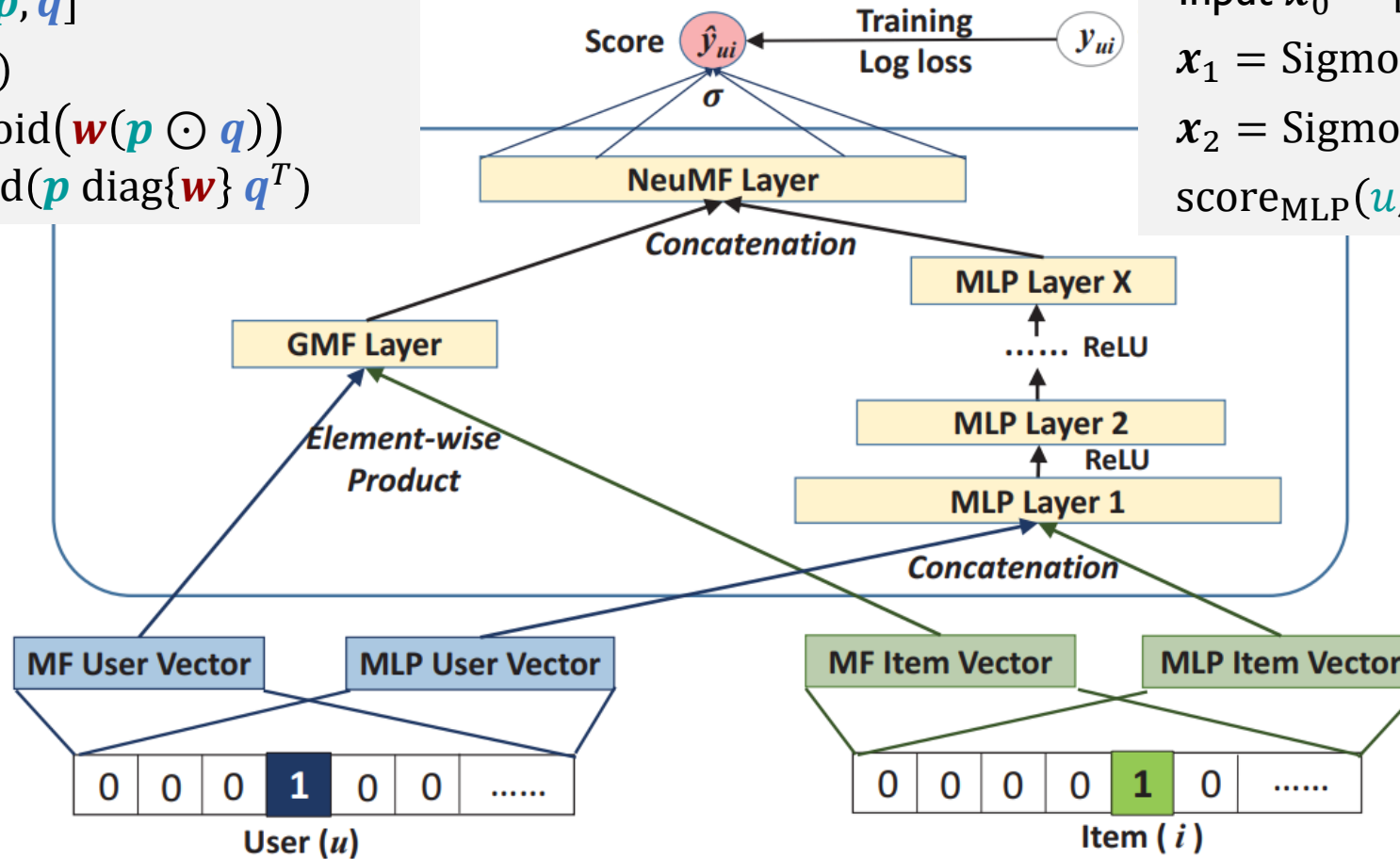
MLP (assuming 3 layers):

Input $x_0 = [p, q]$

$$x_1 = \text{Sigmoid}(W_1 x_0 + b_1)$$

$$x_2 = \text{Sigmoid}(W_2 x_1 + b_2)$$

$$\text{score}_{\text{MLP}}(u, i) = \text{Sigmoid}(W_3 x_2)$$

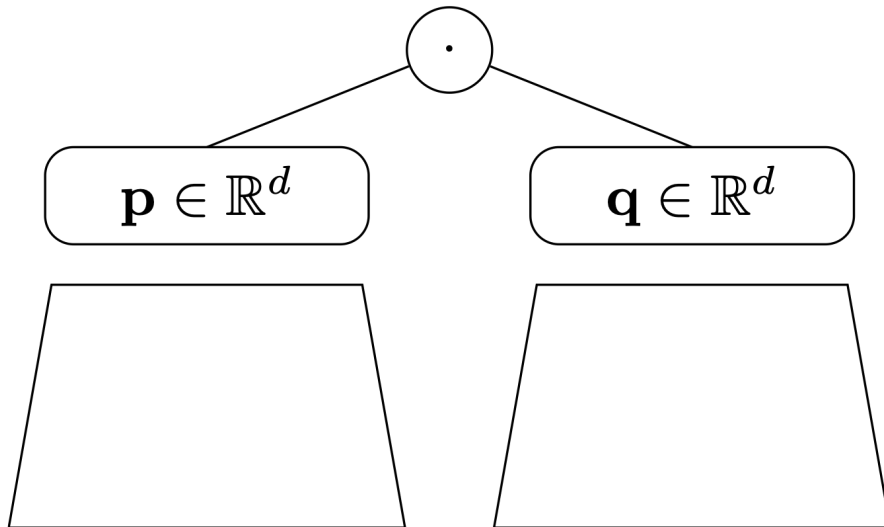


$$\text{score}(u, i) = \text{score}_{\text{MF}}(u, i) + \text{score}_{\text{MLP}}(u, i)$$

Matrix Factorization vs. Neural Collaborative Filtering

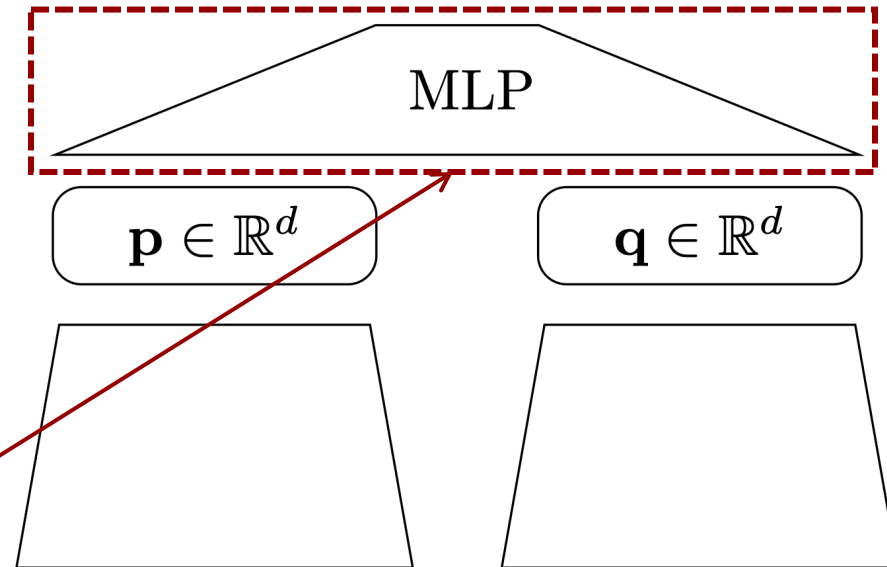
- Predefined similarity

$$\phi^{\text{dot}}(\mathbf{p}, \mathbf{q}) = \langle \mathbf{p}, \mathbf{q} \rangle$$



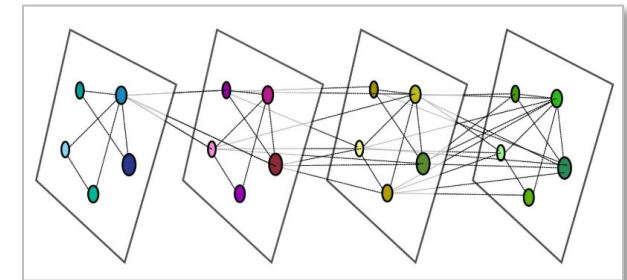
- Learned similarity

$$\phi^{\text{MLP}}(\mathbf{p}, \mathbf{q}) = \mathbf{f}_{W_l, \mathbf{b}_l}(\dots \mathbf{f}_{W_1, \mathbf{b}_1}([\mathbf{p}, \mathbf{q}]) \dots)$$



We can make this part even more complicated!

(E.g., graph neural networks)



But is inner product really that bad?

Performance Revisited [Rendle et al., RecSys 2020]

Neural Collaborative Filtering vs. Matrix Factorization Revisited

Steffen Rendle
srendle@google.com
Google Research
Mountain View, California

Li Zhang
liqzhang@google.com
Google Research
Mountain View, California

Walid Krichene
walidk@google.com
Google Research
Mountain View, California

John Anderson
janders@google.com
Google Research
Mountain View, California

ABSTRACT

Embedding based models have been the state of the art in collaborative filtering for over a decade. Traditionally, the dot product or higher order equivalents have been used to combine two or more embeddings, e.g., most notably in matrix factorization. In recent years, it was suggested to replace the dot product with a learned similarity e.g. using a multilayer perceptron (MLP). This approach is often referred to as *neural collaborative filtering* (NCF). In this work, we revisit the experiments of the NCF paper that popularized learned similarities using MLPs. First we show that with a

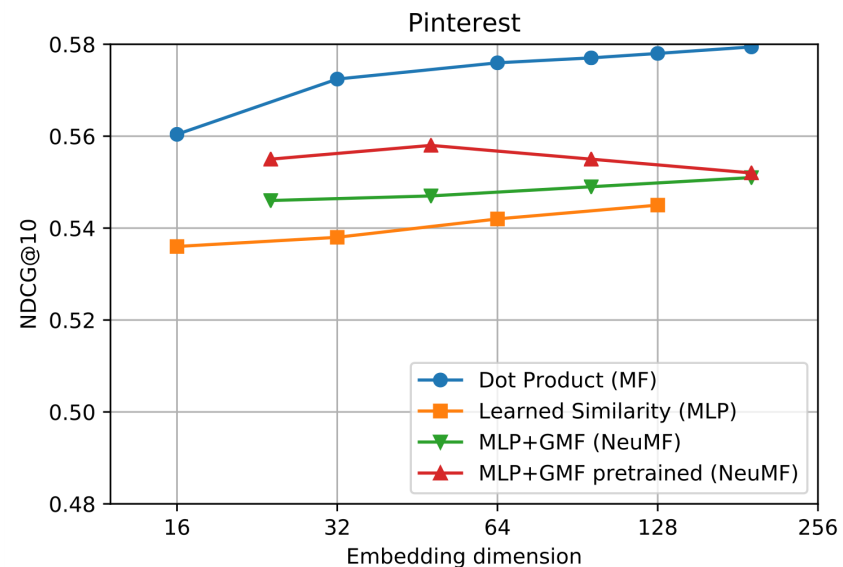
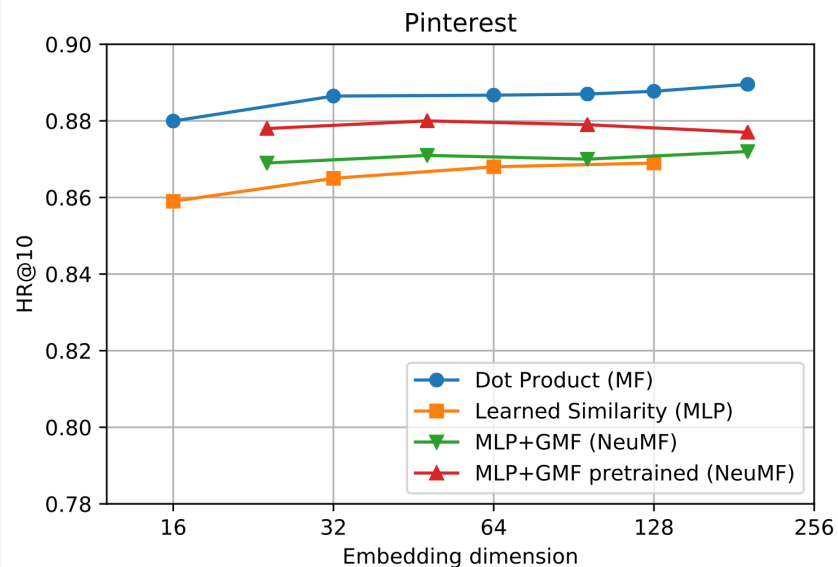
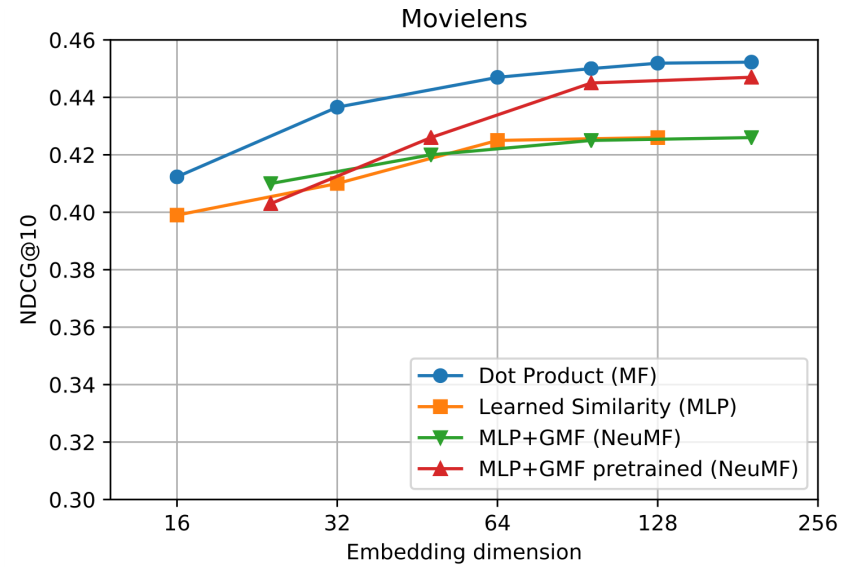
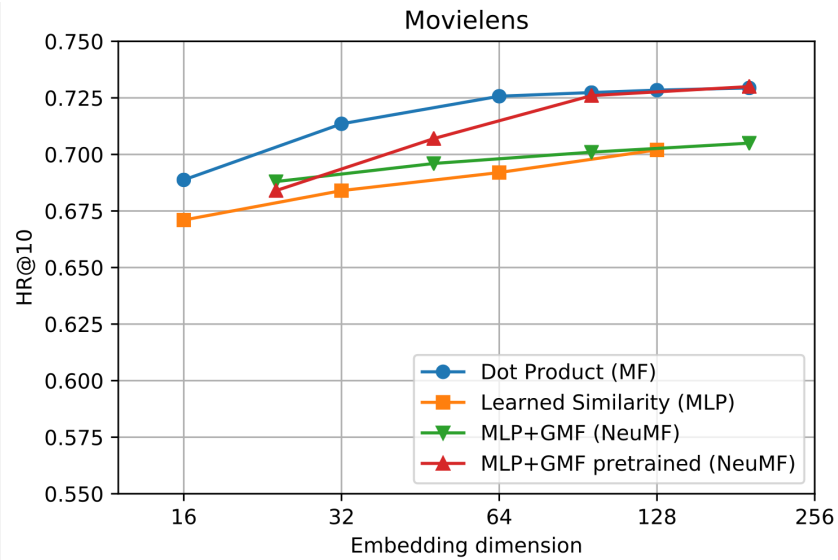
1 INTRODUCTION

Embedding based models have been the state of the art in collaborative filtering for over a decade. A core operation of most of these embedding based models is to combine two or more embeddings. For example, combining a user embedding with an item embedding to obtain a single score that indicates the preference of the user for the item. This can be viewed as a *similarity function* in the embedding space. Traditionally, a dot product or higher order products have been used for the similarity. Recently, it has become popular to learn the similarity function with a neural network. Most

Performance Revisited

- **Motivation:** Learning similarity functions (e.g., **Neural Collaborative Filtering**) has typically been regarded as the de facto standard in academic studies and is considered a strong baseline
- Experimental Setup
 - Follows the Neural Collaborative Filtering paper
 - **Datasets:** MovieLens 1M and Pinterest
 - **Compared Methods:** **Matrix Factorization** (MF) and **Neural Collaborative Filtering** variants

Results

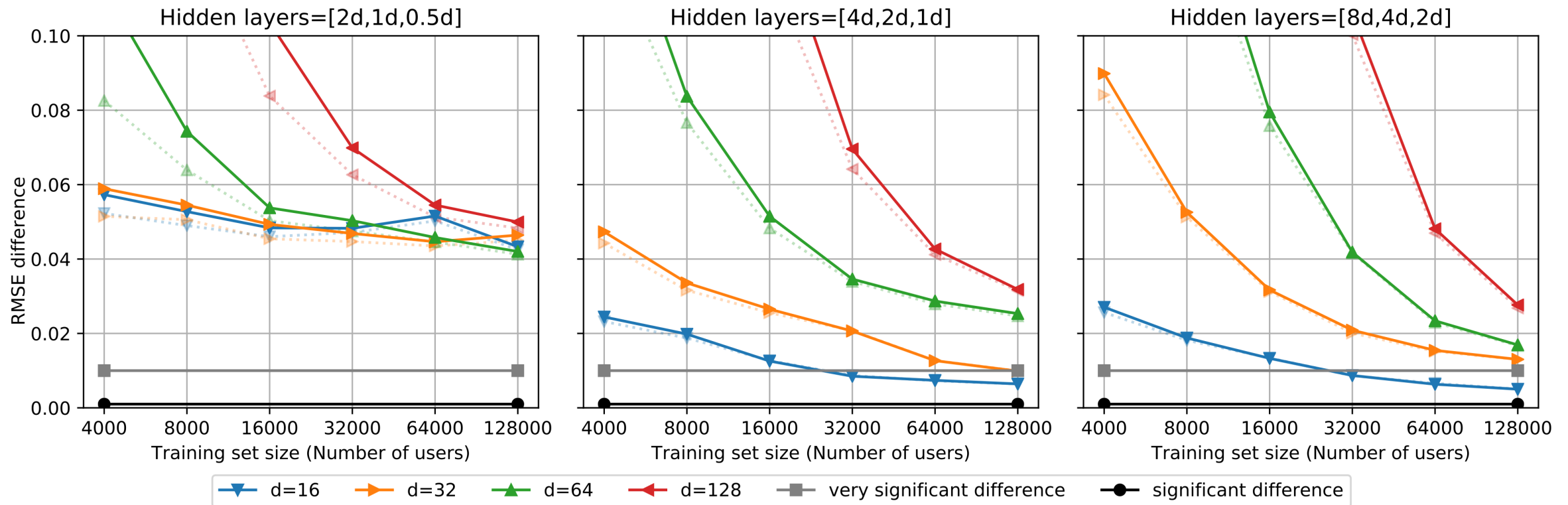


- MF substantially outperforms Neural Collaborative Filtering variants (MLP and NeuMF) on all datasets, evaluation metrics and embedding dimensions.

Why does this happen?

- Why did we believe **Neural Collaborative Filtering** could be better than **MF**?
 - Since **MLPs** can approximate any continuous function under certain conditions (known as the Universal Approximation Theorem), we expect **MLPs** to capture more complex user-item similarity patterns than a simple **inner product**
- But how much training data does **Neural Collaborative Filtering** need to capture complex user-item similarity patterns?
- Let's further simplify the learning problem: how much training data does **Neural Collaborative Filtering** need to capture even the simple **inner product** similarity?

MLPs require massive data even to learn an inner product!

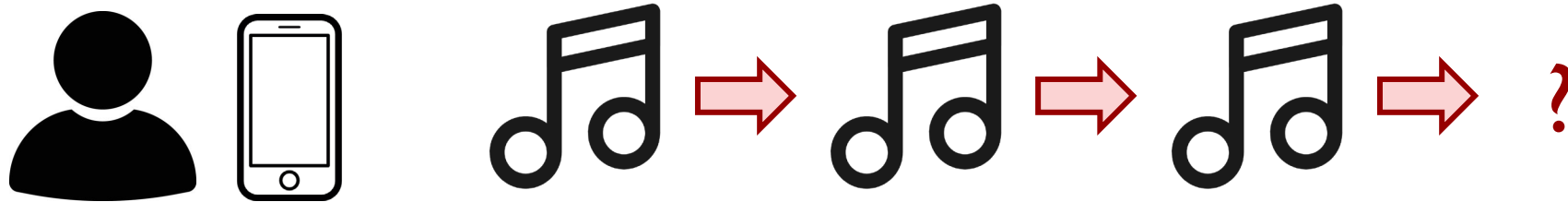


- As models become more complex, we need more training data (for example, the entire Wikipedia and BookCorpus used for BERT). However, in recommender systems, we rarely have access to such large amounts of data.

Sequential Recommendation

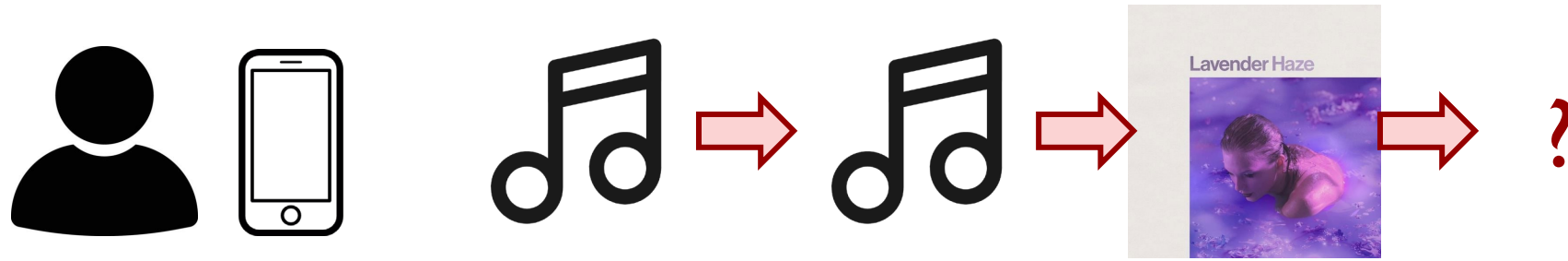
We are adding streaming to our record store!

- What is the next track to stream to a user?



We are adding streaming to our record store!

- What is the next track to stream to a user?



- Our user just listened to Taylor Swift's *Lavender Haze*
 - What should we do next?
 - Play more songs by Taylor Swift?
 - Play more color-themed songs?
 - Play more modern pop songs?
 - Play it again?

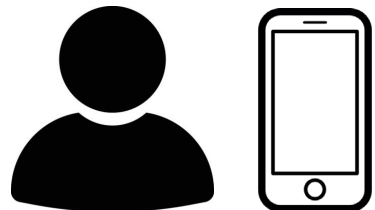
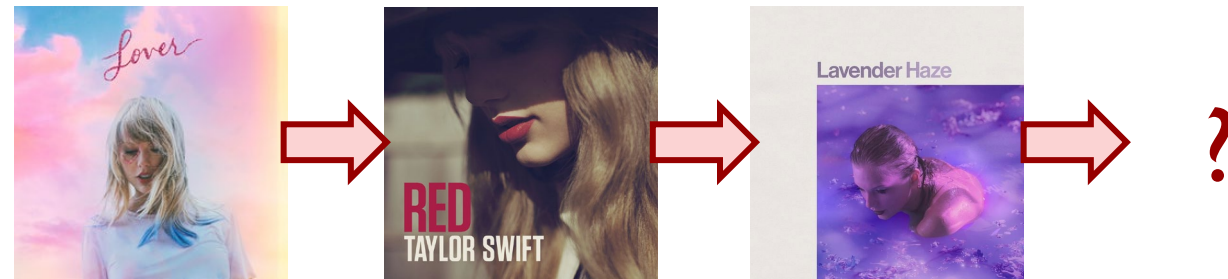
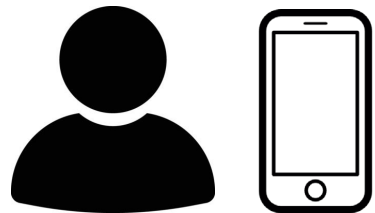
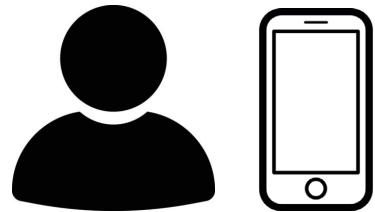
- **Question 1:** What do these options have in common?
- **Question 2:** What other factors should I consider in choosing among these options?

What RecSys techniques have we already learned?

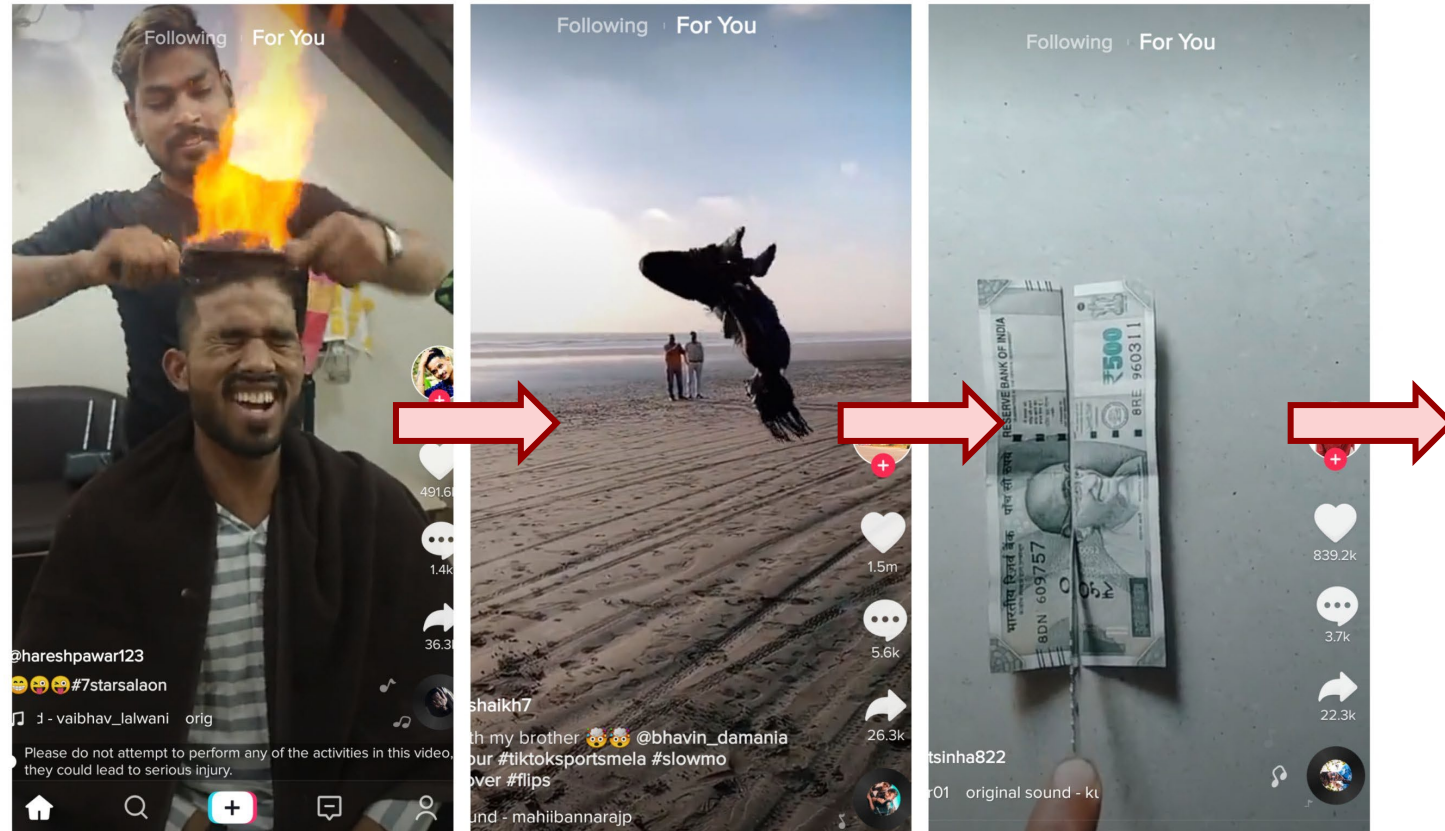
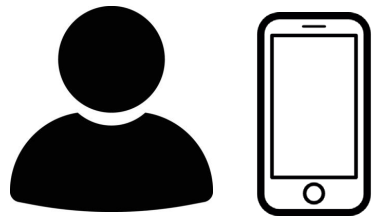
- **Content-Based Approach:** Recommend items that are **similar to** those previously viewed
 - Here, the definition of **similarity** is based on **attributes** (e.g., lyrics)
- **Item-Item Collaborative Filtering:** Given a candidate item, first identify its **similar** items, and then check whether the user has viewed any of them
 - Here, the definition of **similarity** is based on the **columns corresponding to the two items** in the user-item matrix
- **Matrix Factorization:** Find the item in the **embedding space** that are **closest to** the user and have not yet been viewed
 - A user's **embedding** should be **close to** the **embeddings** of the items they have already viewed
- **Neural Collaborative Filtering:** More sophisticated **similarity** computation methods, but still based on **embeddings**
- **Summary:** Find items **similar to** those a user has already viewed. The only difference lies in how **similarity** is defined

What are we missing?

- Understanding dynamic user intent from sequential interactions



Other Sequential RecSys Examples: TikTok

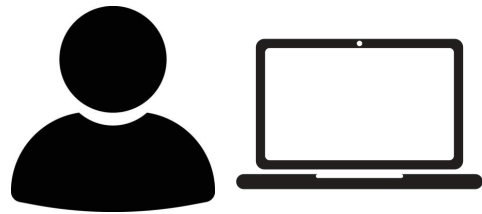


Other Sequential RecSys Examples: Netflix



- For certain products, the **order** can be crucial!
 - [movies and books] Suppose I just watched *Harry Potter 2*. Would it make sense to recommend *Harry Potter 1*?
 - [lecture videos] Suppose I just watched the video for **Lecture 3** of Stanford's CS224N. Would it make sense to recommend the video for **Lecture 2**?
 - [electronic devices] Suppose I just bought an **iPhone XR**. Would it make sense to recommend an **iPhone 8**?

Other Sequential RecSys Examples: YouTube and Amazon



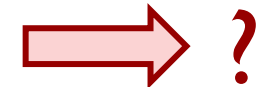
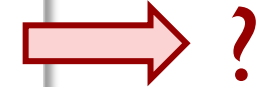
Natural Language Processing
with Deep Learning
CS224N/Ling284

Christopher Manning
Lecture 2: Word Vectors, Word Senses, and Neural Classifiers



Natural Language Processing
with Deep Learning
CS224N/Ling284

Christopher Manning
Lecture 3: Neural net learning: Gradients by hand (matrix calculus)
and algorithmically (the backpropagation algorithm)



What are we missing?

- Embedding similarity is **symmetric**.
 - If the user has just viewed item i , whether the model will recommend item j as the next item fully depends on $\text{score}(i, j) = \mathbf{e}_i^T \mathbf{e}_j$ or $\cos(\mathbf{e}_i, \mathbf{e}_j)$
 - If the user has just viewed item j , whether the model will recommend item i as the next item fully depends on $\text{score}(j, i) = \text{score}(i, j)$
- Sequential recommendation is **asymmetric**.
 - When your model recommends item j as the next item after item i , it does NOT necessarily mean that i should be recommended as the next item after j .
 - We should model something like $\text{score}(j | i) \neq \text{score}(i | j)$

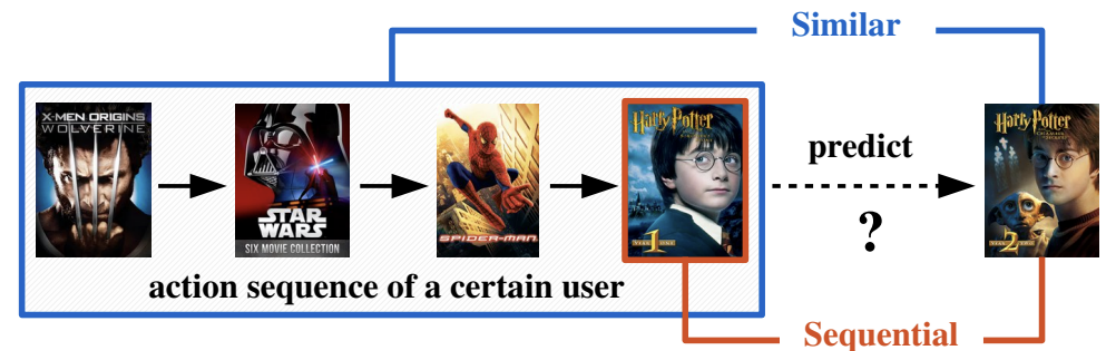
Transition matrix is asymmetric!

Fossil [He and McAuley, ICDM 2016]

Fusing Similarity Models with Markov Chains for Sparse Sequential Recommendation

Ruining He, Julian McAuley
Department of Computer Science and Engineering
University of California, San Diego
Email: {r4he, jmcauley}@cs.ucsd.edu

Abstract—Predicting personalized sequential behavior is a key task for recommender systems. In order to predict user actions such as the next product to purchase, movie to watch, or place to visit, it is essential to take into account both long-term user preferences and sequential patterns (i.e., short-term dynamics). Matrix Factorization and Markov Chain methods have emerged as two separate but powerful paradigms for modeling the two respectively. Combining these ideas has led





Rucanor Aluminium Baseball Bat

by [Rucanor](#)

★★★★☆ (29 customer reviews)

RRP: ~~£21.00~~

Price: **£18.40** & **FREE Delivery** in the UK. [Details](#)

You Save: **£2.60** (12%)

Colour: Silver

Size: 60

In stock.

Sold by [MA sports uk](#) and [Fulfilled by Amazon](#). Gift-wrap available.

- 60 CM BASEBALL BAT
- ALUMINIUM Alloy

5 new from **£14.40**

Quantity: 1 ▼

Add to Basket

or 1-Click Checkout

Buy now with 1-Click®

This is a gift

Add to Wish List ▼

More Buying Choices

ANDRA SPORTS **Add to Basket**
£14.40 + £3.99 UK delivery

sportknox **Add to Basket**
£17.45 + £3.99 UK delivery

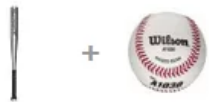
5 new from **£14.40**

Have one to sell? [Sell yours here](#)

Share

Roll over image to zoom in
[Share your own customer images](#)

Frequently Bought Together



Price For Both: **£28.39**

Add both to Basket

Show availability and delivery details

- This item:** Rucanor Aluminium Baseball Bat, Silver - 60 cm **£18.40**
- WILSON Official League Individual Baseball **£9.99**

Customers Who Bought This Item Also Bought

Page 1 of 18



REYDON Softball Ball
★★★★★ (1)
£3.99



MFH Balaclava 3 Hole Black
★★★★☆ (94)
£1.74



New Midwest Slugger
Baseball Glove Vinyl
Catching Mitt Left Hand
Junior / Senior
★★★★★ (5)



Rawlings 9" Indoor / Outdoor
T-Ball Training Baseball -
TVB
★★★★★ (3)
£1.99



- Predict a user's next step based on the next-step behavior patterns of the **overall population**
- Amazon's "*Customer Who Bought This Item Also Bought*"

Markov Chain to Model Sequential Dynamics

- Historical user-item interactions

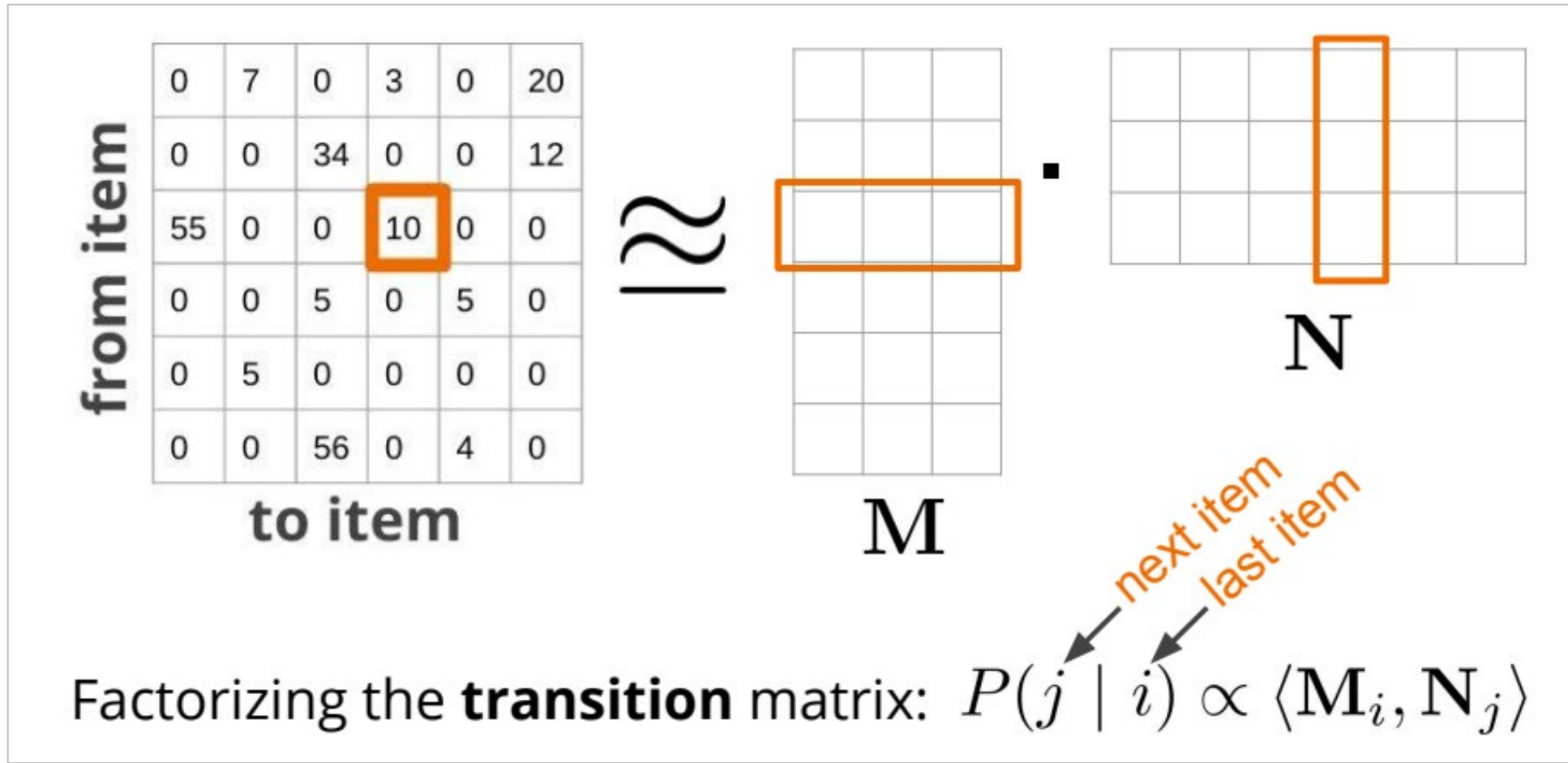
User								
<i>u1</i>	→	<i>i1</i>	→	<i>i2</i>	→	<i>i4</i>		
<i>u2</i>	→	<i>i3</i>	→	<i>i4</i>				
<i>u3</i>	→	<i>i1</i>	→	<i>i2</i>	→	<i>i4</i>	→	<i>i3</i>
<i>u4</i>	→	<i>i3</i>	→	<i>i1</i>	→	<i>i2</i>		

- **Step 1:** Construct a transition matrix

	<i>i1</i>	<i>i2</i>	<i>i3</i>	<i>i4</i>
<i>i1</i>	0	3	0	0
<i>i2</i>	0	0	0	2
<i>i3</i>	1	0	0	1
<i>i4</i>	0	0	1	0

Markov Chain to Model Sequential Dynamics

- **Step 2:** Factorize the transition matrix

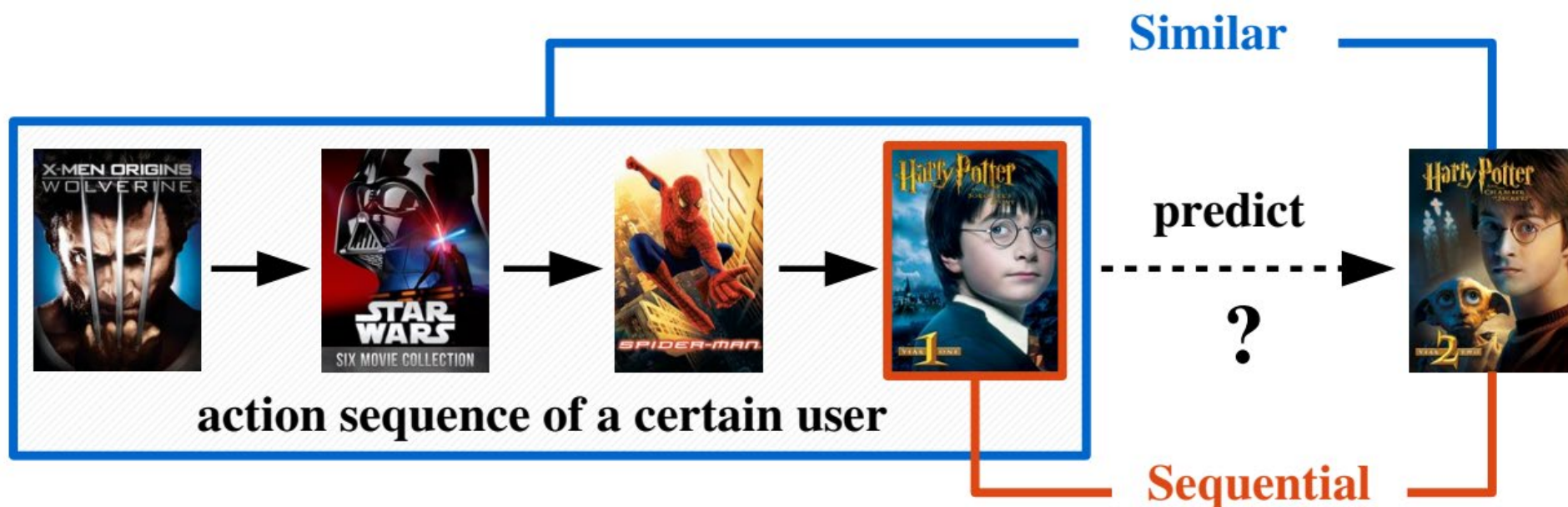


Markov Chain to Model Sequential Dynamics

- Problems with this approach?
 - NOT personalized!
 - Regardless of which user the model is considering, $\text{score}(j | i) \propto \mathbf{M}_i^T \mathbf{N}_j$
 - But \mathbf{M}_i and \mathbf{N}_j are the same for all users
- How to make the model personalized?

$$p_u(j | i) \propto \underbrace{\sum_{j' \in \mathcal{I}_u^+ \setminus \{j\}} \langle \mathbf{P}_{j'}, \mathbf{Q}_j \rangle}_{\text{user preferences}} + \underbrace{(\eta + \eta_u)}_{\text{personalized weighting factor}} \cdot \underbrace{\langle \mathbf{M}_i, \mathbf{N}_j \rangle}_{\text{sequential dynamics}}$$

$$p_u(j | i) \propto \underbrace{\sum_{j' \in \mathcal{I}_u^+ \setminus \{j\}} \langle \mathbf{P}_{j'}, \mathbf{Q}_j \rangle}_{\text{user preferences}} + \underbrace{(\eta + \eta_u)}_{\text{personalized weighting factor}} \cdot \underbrace{\langle \mathbf{M}_i, \mathbf{N}_j \rangle}_{\text{sequential dynamics}}$$



Back to Motivating Example

- Play *Lavender Haze* again? Play more songs by Taylor Swift? Play more color-themed songs?
 - Might all be likely according to the transition matrix. But if you consider user preference ...



Higher-Order Markov Chain

User								
u_1	→	i_1	→	i_2	→	i_4		
u_2	→	i_3	→	i_4				
u_3	→	i_1	→	i_2	→	i_4	→	i_3
u_4	→	i_3	→	i_1	→	i_2		

- Transition matrix:

	i_1	i_2	i_3	i_4
$i_1 \rightarrow i_2$	0	0	0	2
$i_1 \rightarrow i_3$	0	0	0	0
$i_1 \rightarrow i_4$	0	0	0	0
$i_2 \rightarrow i_1$	0	0	0	0
...

- **Sparse**
- # of rows **increases exponentially** with the order of the Markov chain
- Similar to the n-gram language model
- Can we use a **neural language model** to alleviate sparsity?

Performance

- Fossil outperforms BPR-MF (and other static baselines) because it models sequential dynamics
- Small orders seem to be enough to achieve good performance

Table IV: AUC on different datasets (higher is better). The number of latent dimensions for all comparison methods (except for POP) is set to $K = 10$. For *Fossil*, we test different orders of the Markov Chain (i.e., 1, 2, and 3). On the right we demonstrate the improvement of FPMC vs. BPR-MF, *Fossil* vs. FISM, *Fossil* vs. FPMC, and *Fossil* vs. the best baselines.

Dataset	(a)	(b)	(c)	(d)	(e)	(f-1)	(f-2)	(f-3)	% improvement			
	POP	BPR-MF	FISM	FMC	FPMC	<i>Fossil</i>	<i>Fossil</i>	<i>Fossil</i>	e vs. b	f vs. c	f vs. e	f vs. best
<i>Amazon-Office</i>	0.6427	0.6736	0.7113	0.6874	0.6891	0.7211	0.7224	0.7221	2.30%	1.56%	4.83%	1.56%
<i>Amazon-Auto</i>	0.5870	0.6379	0.6736	0.6452	0.6446	0.6910	0.6904	0.6901	1.05%	2.58%	7.20%	2.58%
<i>Amazon-Game</i>	0.7495	0.8483	0.8639	0.8401	0.8502	0.8793	0.8813	0.8817	0.22%	2.06%	3.71%	2.06%
<i>Amazon-Toy</i>	0.6240	0.7020	0.7499	0.6665	0.7061	0.7625	0.7645	0.7652	0.58%	2.04%	8.37%	2.04%
<i>Amazon-Cell</i>	0.6959	0.7212	0.7755	0.7359	0.7396	0.7982	0.8009	0.8006	2.55%	3.27%	8.29%	3.27%
<i>Amazon-Clothes</i>	0.6189	0.6513	0.7085	0.6673	0.6672	0.7255	0.7256	0.7259	2.44%	2.46%	8.80%	2.46%
<i>Amazon-Elec</i>	0.7837	0.7927	0.8210	0.7992	0.7985	0.8411	0.8438	0.8444	0.73%	2.85%	5.75%	2.85%
<i>Epinions</i>	0.4576	0.5520	0.5818	0.5532	0.5477	0.6014	0.6050	0.6048	-0.78%	3.99%	10.46%	3.99%
<i>Foursquare</i>	0.9168	0.9506	0.9230	0.9441	0.9485	0.9626	0.9621	0.9618	-0.22%	4.29%	1.49%	1.26%
Avg. ($K = 10$)	0.6751	0.7255	0.7565	0.7265	0.7324	0.7759	0.7773	0.7774	0.99%	2.79%	6.54%	2.45%
Avg. ($K = 20$)	0.6751	0.7285	0.7580	0.7293	0.7344	0.7780	0.7795	0.7788	0.88%	2.83%	6.44%	2.54%

Translation is asymmetric!

TransRec [He et al., RecSys 2017]

Translation-based Recommendation

Ruining He
UC San Diego
r4he@cs.ucsd.edu

Wang-Cheng Kang
UC San Diego
wckang@eng.ucsd.edu

Julian McAuley
UC San Diego
jmcauley@cs.ucsd.edu

ABSTRACT

Modeling the complex interactions between users and items as well as amongst items themselves is at the core of designing successful recommender systems. One classical setting is predicting users' personalized sequential behavior (or 'next-item' recommendation), where the challenges mainly lie in modeling 'third-order' interactions between a user, her previously visited item(s), and the next item to consume. Existing methods typically decompose these higher-order interactions into a combination of *pairwise* relationships, by way of which user preferences (user-item interactions) and sequential patterns (item-item interactions) are captured by separate components. In this paper, we propose a unified method, *TransRec*, to model such third-order relationships for large-scale sequen-

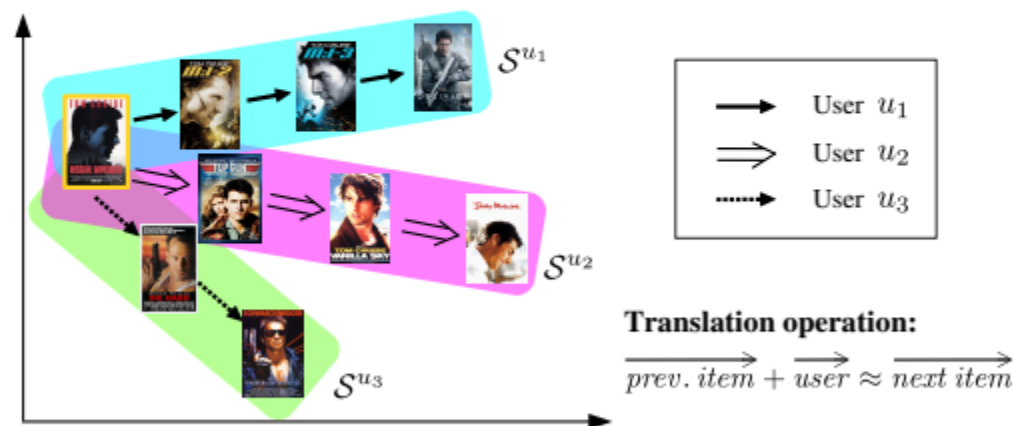
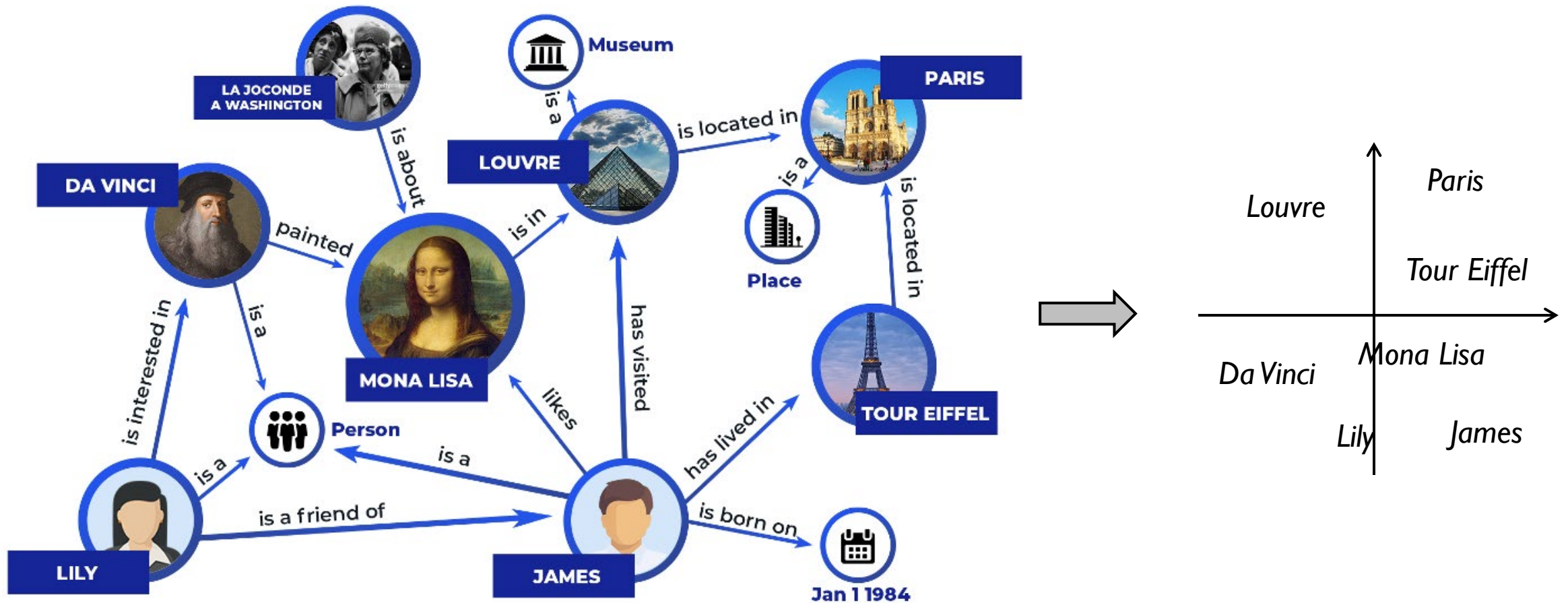


Figure 1: *TransRec* as a sequential model: Items (movies) are embedded into a 'transition space' where each user is modeled by a *translation* vector. The transition of a user from

Background: Translation-Based Embedding

- You have a knowledge graph, and you want to embed all the nodes into a vector space



Background: Translation-Based Embedding

- Each edge can be represented as a triplet (h, r, t)

head entity h	relation r	tail entity t
<i>Tour Eiffel</i>	is located in	<i>Paris</i>
<i>Da Vinci</i>	painted	<i>Mona Lisa</i>
<i>Lily</i>	is a friend of	<i>James</i>
...

- **A simple (but suboptimal) solution:** Each node should be close to its neighbors

$$\max_{e_h, e_t} \sum_{(h,r,t) \in \mathcal{G}} \frac{\exp(e_h^T e_t)}{\sum_{t'} \exp(e_h^T e_{t'})}$$

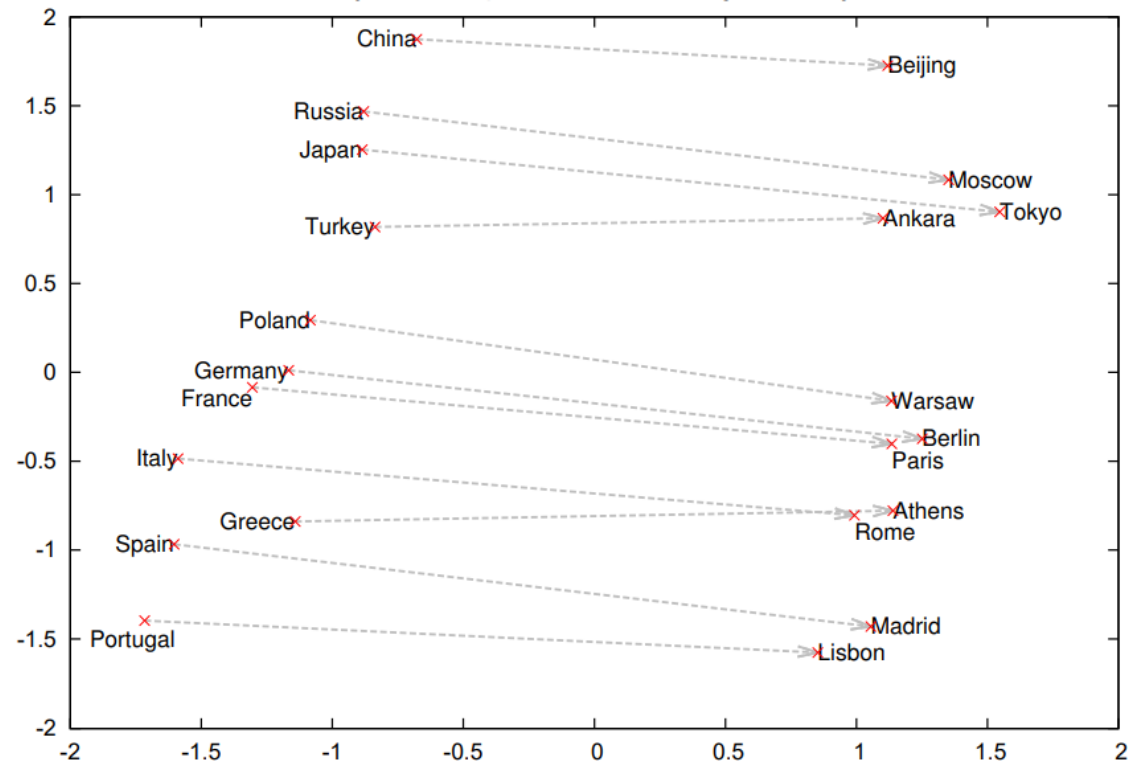
- Limitation: It ignores different relations

TransE [Bordes et al., NIPS 2013]

- Idea: Model relations as translations

head entity h	relation r	tail entity t
<i>Paris</i>	is capital of	<i>France</i>

$$e_h + e_r \approx e_t$$



TransE [Bordes et al., NIPS 2013]

- **Idea:** Model relations as translations
- **Training:**

head entity h	relation r	tail entity t
<i>Paris</i>	is capital of	<i>France</i>

head entity h	relation r	“corrupted” tail entity t'
<i>Paris</i>	is capital of	<i>UK</i>

Minimize $\|e_h + e_r - e_t\|$

Maximize $\|e_h + e_r - e_{t'}\|$
(but stop when it becomes sufficiently larger than $\|e_h + e_r - e_t\|$)

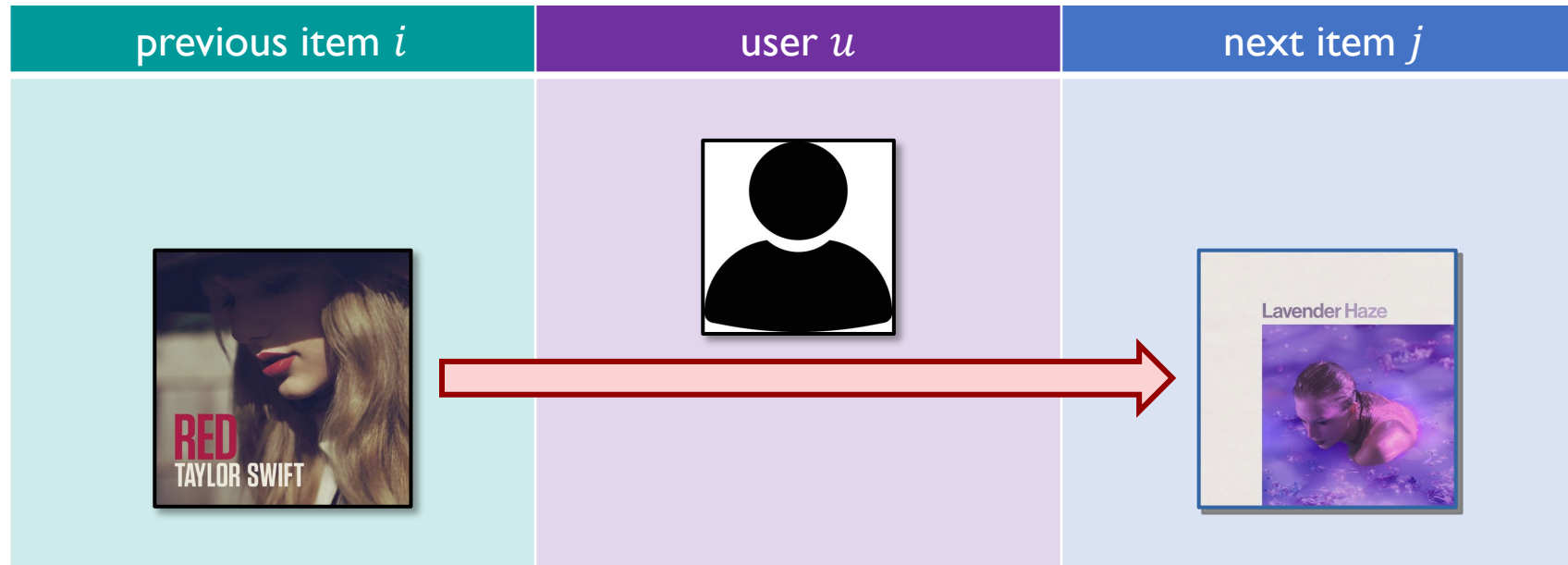
- **Testing:**

head entity h	relation r	tail entity t
<i>Ottawa</i>	is capital of	?

Find the nearest neighbor of $e_h + e_r$

Translation-Based Recommendation

- Items as points
- Users as translational vectors

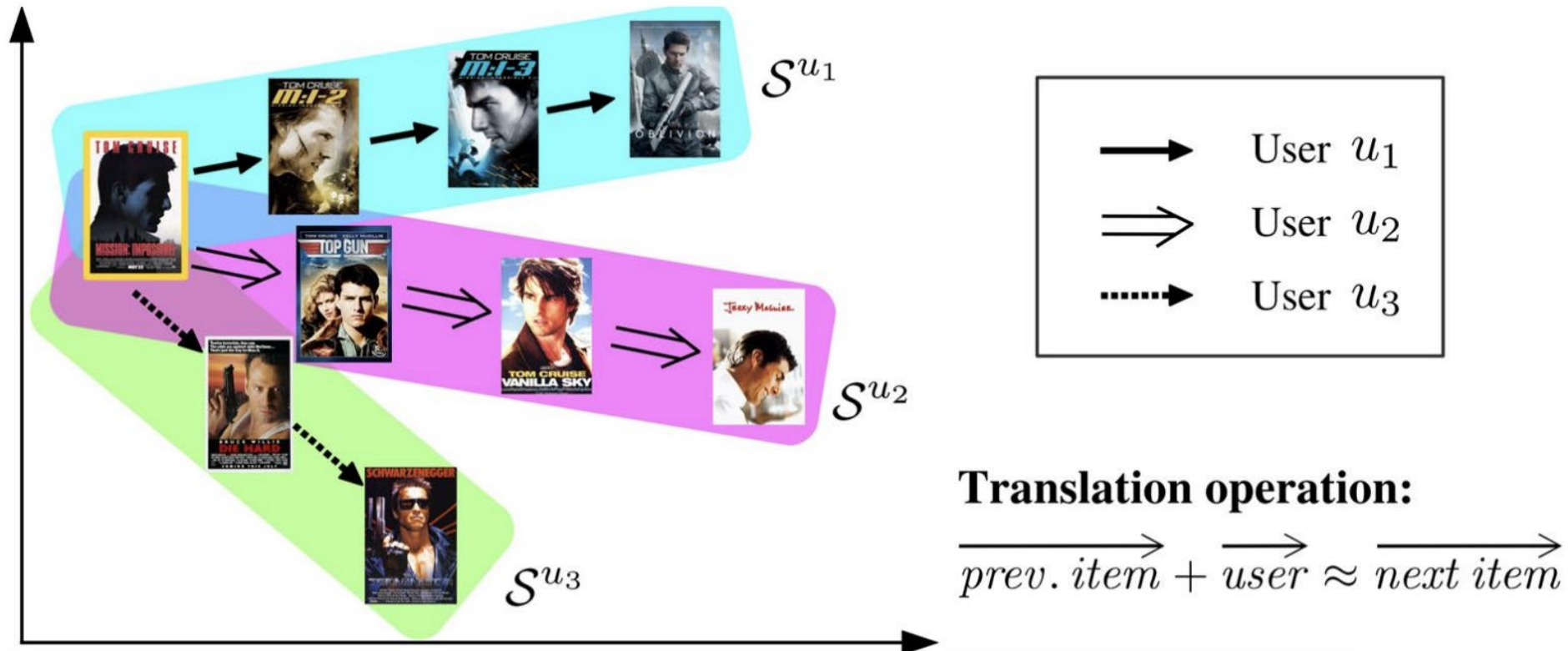


$$e_i + e_u \approx e_j$$

Translation-Based Recommendation

- Items as points
- Users as translational vectors
 - E.g., the songs progressively become more up-tempo, as the user is at a party
 - e_u : “the next song is more up-tempo than the previous song”
 - E.g., the songs gradually become more relaxing, as the user is going to sleep
 - e_u : “the next song is more relaxing than the previous song”

Translation-Based Recommendation



- **Testing time:** Find the nearest neighbor of $e_i + e_u$

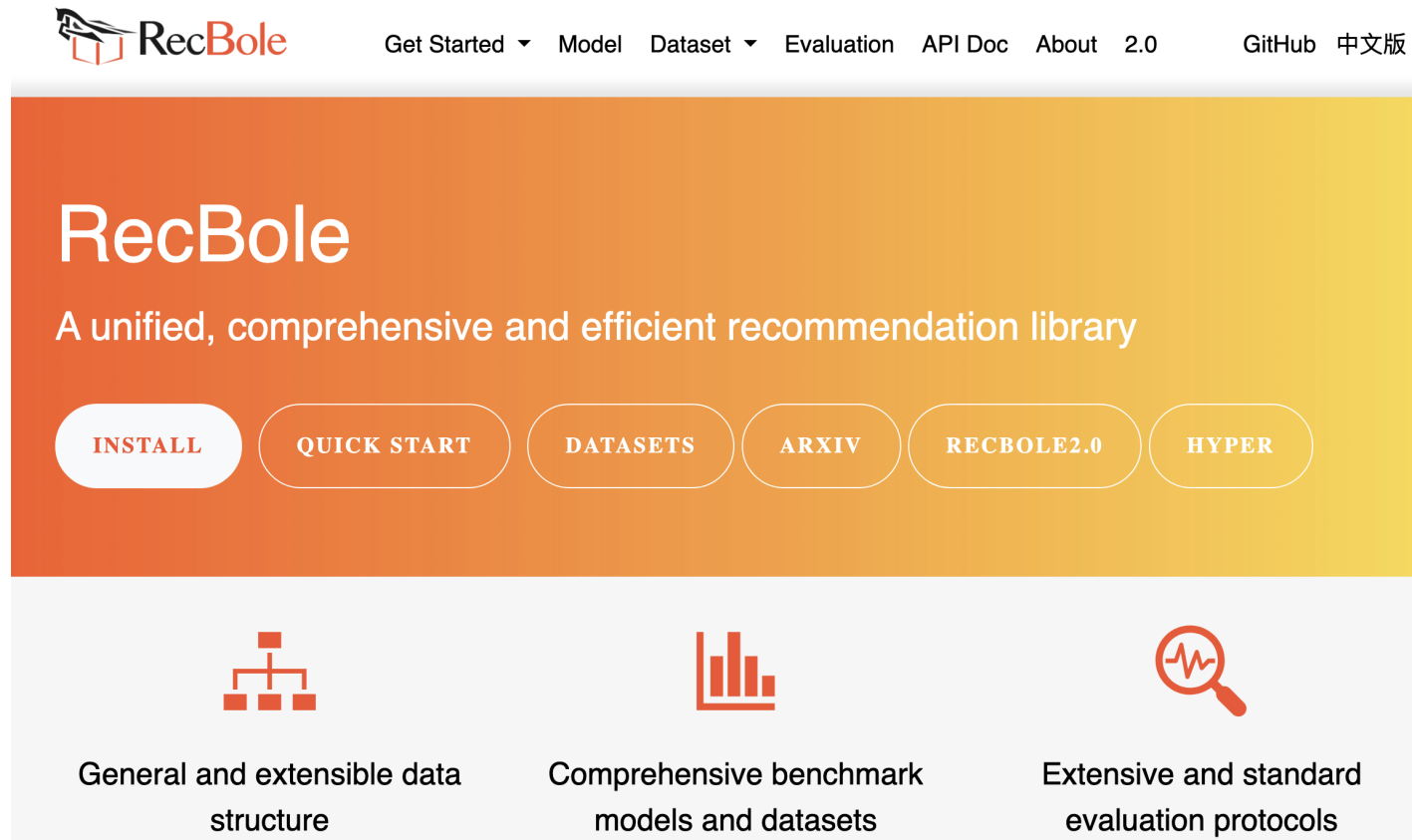
Performance

- Outperforms BPR-MF
- No comparisons with Fossil

Dataset	Metric	PopRec	BPR-MF	FMC	FPMC	HRM _{avg}	HRM _{max}	PRME	TransRec \mathcal{L}_1	TransRec \mathcal{L}_2	%Improv.
<i>Epinions</i>	<i>AUC</i>	0.4576	0.5523	0.5537	0.5517	0.6060	0.5617	0.6117	0.6063	<u>0.6133</u>	0.3%
	<i>Hit@50</i>	3.42%	3.70%	3.84%	2.93%	3.44%	2.79%	2.51%	3.18%	<u>4.63%</u>	20.6%
<i>Automotive</i>	<i>AUC</i>	0.5870	0.6342	0.6438	0.6427	0.6704	0.6556	0.6469	0.6779	<u>0.6868</u>	2.5%
	<i>Hit@50</i>	3.84%	3.80%	2.32%	3.11%	4.47%	3.71%	3.42%	5.07%	<u>5.37%</u>	20.1%
<i>Google</i>	<i>AUC</i>	0.5391	0.8188	0.7619	0.7740	0.8640	0.8102	0.8252	0.8359	<u>0.8691</u>	0.6%
	<i>Hit@50</i>	0.32%	4.27%	3.54%	3.99%	3.55%	4.59%	5.07%	6.37%	<u>6.84%</u>	32.5%
<i>Office</i>	<i>AUC</i>	0.6427	0.6979	0.6867	0.6866	0.6981	0.7005	0.7020	0.7186	<u>0.7302</u>	4.0%
	<i>Hit@50</i>	1.66%	4.09%	2.66%	2.97%	5.50%	4.17%	6.20%	<u>6.86%</u>	6.51%	10.7%
<i>Toys</i>	<i>AUC</i>	0.6240	0.7232	0.6645	0.7194	0.7579	0.7258	0.7261	0.7442	<u>0.7590</u>	0.2%
	<i>Hit@50</i>	1.69%	3.60%	1.55%	4.41%	5.25%	3.74%	4.80%	<u>5.46%</u>	5.44%	4.0%
<i>Clothing</i>	<i>AUC</i>	0.6189	0.6508	0.6640	0.6646	0.7057	0.6862	0.6886	0.7047	<u>0.7243</u>	2.6%
	<i>Hit@50</i>	1.11%	1.05%	0.57%	0.51%	1.70%	1.15%	1.00%	1.76%	<u>2.12%</u>	24.7%
<i>Cellphone</i>	<i>AUC</i>	0.6959	0.7569	0.7347	0.7375	0.7892	0.7654	0.7860	0.7988	<u>0.8104</u>	2.7%
	<i>Hit@50</i>	4.43%	5.15%	3.23%	2.81%	8.77%	6.32%	6.95%	9.46%	<u>9.54%</u>	8.8%
<i>Games</i>	<i>AUC</i>	0.7495	0.8517	0.8407	0.8523	0.8776	0.8566	0.8597	0.8711	<u>0.8815</u>	0.4%
	<i>Hit@50</i>	5.17%	10.93%	13.93%	12.29%	14.44%	12.86%	14.22%	<u>16.61%</u>	16.44%	15.0%
<i>Electronics</i>	<i>AUC</i>	0.7837	0.8096	0.8158	0.8082	0.8212	0.8148	0.8337	0.8457	<u>0.8484</u>	1.8%
	<i>Hit@50</i>	4.62%	2.98%	4.15%	2.82%	4.09%	2.59%	3.07%	4.89%	<u>5.19%</u>	12.3%
<i>Foursquare</i>	<i>AUC</i>	0.9168	0.9511	0.9463	0.9479	0.9559	0.9523	0.9565	0.9631	<u>0.9651</u>	0.9%
	<i>Hit@50</i>	55.60%	60.03%	63.00%	64.53%	60.75%	61.60%	65.32%	66.12%	<u>67.09%</u>	2.7%
<i>Flixter</i>	<i>AUC</i>	0.9459	0.9722	0.9568	0.9718	0.9695	0.9687	0.9728	0.9727	<u>0.9750</u>	0.2%
	<i>Hit@50</i>	11.92%	21.58%	22.23%	33.11%	32.34%	30.88%	<u>40.81%</u>	35.52%	35.02%	-13.0%

Homework 4 (Part 1): Compare Fossil and TransRec

- No need to implement them yourself!
- The RecBole library: <https://recbole.io/>



The screenshot shows the RecBole website homepage. At the top left is the RecBole logo, which consists of a stylized red and white icon followed by the text "RecBole". To the right of the logo is a navigation menu with the following items: "Get Started" (with a dropdown arrow), "Model", "Dataset" (with a dropdown arrow), "Evaluation", "API Doc", "About", "2.0", "GitHub", and "中文版". Below the navigation menu is a large orange and yellow gradient banner. The banner contains the text "RecBole" in large white font, followed by the tagline "A unified, comprehensive and efficient recommendation library" in smaller white font. Below the tagline are six white rounded rectangular buttons with orange text: "INSTALL", "QUICK START", "DATASETS", "ARXIV", "RECBOLE2.0", and "HYPER". Below the banner is a light gray section with three columns. Each column has an icon and a text description: 1. A tree diagram icon above the text "General and extensible data structure". 2. A bar chart icon above the text "Comprehensive benchmark models and datasets". 3. A magnifying glass over a pulse line icon above the text "Extensive and standard evaluation protocols".



Thank You!

Course Website: <https://yuzhang-teaching.github.io/CSCE670-S26.html>