



MIE524 Data Mining Recommender Systems: Neural Collaborative Filtering

Reminder: Collaborative Filtering

- **Does not require profiles (features) for items or users**
 - Different from content-based recommender systems
- **Instead, only requires user-item interaction data**
 - *Explicit* feedback
 - e.g., 0-5 stars
 - *Implicit* feedback
 - e.g., user watched movie (0/1)

Reminder: Latent Factor Models

- Approximate the rating matrix \mathbf{R} as product of matrices $\mathbf{Q} \cdot \mathbf{P}^T$
 - **Problem:** \mathbf{R} has missing entries
 - Learn \mathbf{Q} and \mathbf{P} that minimize the **reconstruction error** on known ratings and ignore the missing values

$$\hat{r}_{xi} = q_i \cdot p_x = \sum_f q_{if} \cdot p_{xf}$$

q_i = row i of \mathbf{Q}
 p_x = column x of \mathbf{P}^T

users											
items											
1	3		5		5	4					
		5	4		4		2	1	3		
2	4		1	2		3	4	3	5		
	2	4		5		4		2			
		4	3	4	2			2	5		
1	3	3		2			2		4		

\mathbf{R} \mathbf{Q}

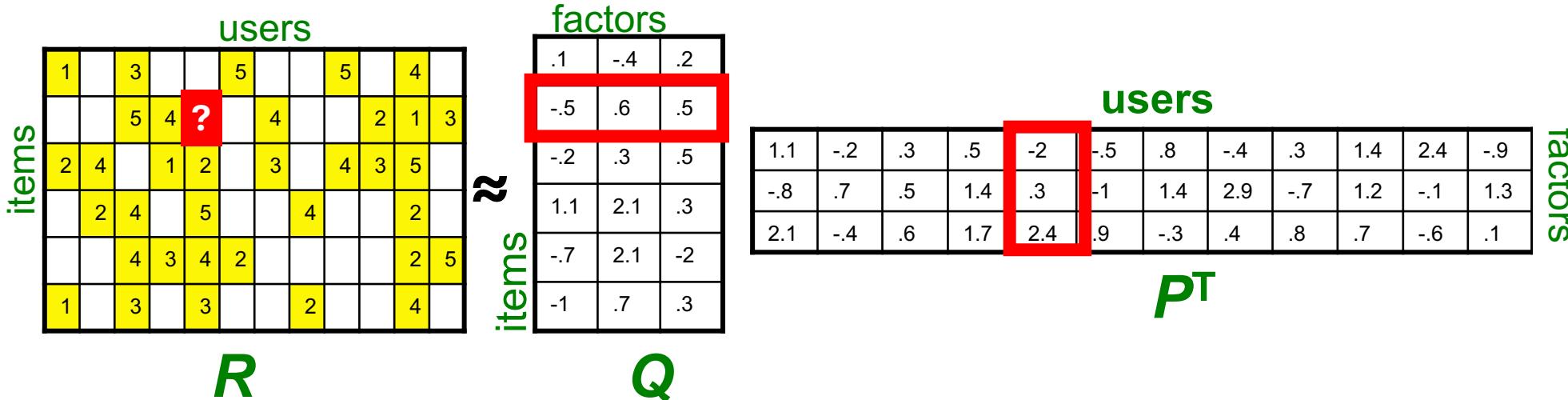
users											
items											
.1	-.4	.2									
-.5	.6	.5									
-.2	.3	.5									
1.1	2.1	.3									
-.7	2.1	-2									
-1	.7	.3									

\mathbf{P}^T factors

Reminder: Latent Factor Models

- How to estimate missing ratings?

- Estimate using latent factors $\hat{r}_{xi} = q_i \cdot p_x = \sum_f q_{if} \cdot p_{xf}$
 q_i = row i of Q
 p_x = column x of P^T



Motivation for Neural Approaches

- **Matrix Factorization is restricted by the simple choice of interaction function between users and items -- the dot (inner) product**
 - Inner product restricted to simple linear combination of multiplication of latent features
 - may not be sufficient to capture the complex structure of user interaction data
 - We've seen that even simple incorporation of user/item biases into the interaction function can improve performance
- **Can we use (deep) neural networks to learn the interaction function from data?**

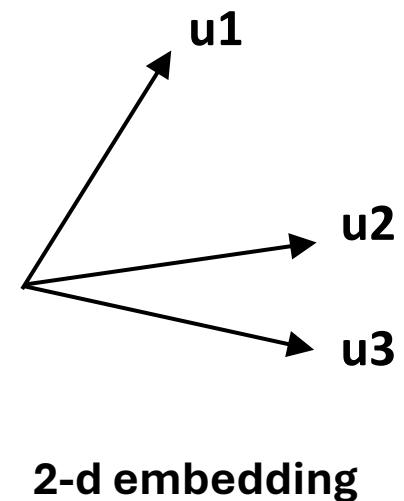
Motivating Example

$$\hat{y}_{ui} = \mathbf{U}_i^T \mathbf{V}_j \simeq \cos(\mathbf{U}_i, \mathbf{V}_j) \quad (\text{E.g., assuming a unit length})$$

	\mathbf{i}_1	\mathbf{i}_2	\mathbf{i}_3	\mathbf{i}_4	\mathbf{i}_5
\mathbf{u}_1	1	1	1	0	1
\mathbf{u}_2	0	1	1	0	0
\mathbf{u}_3	0	1	1	1	0

$$\begin{aligned}\text{sim}(\mathbf{u}_1, \mathbf{u}_2) &= 0.5 \\ \text{sim}(\mathbf{u}_3, \mathbf{u}_1) &= 0.4 \\ \text{sim}(\mathbf{u}_3, \mathbf{u}_2) &= 0.66\end{aligned}$$

$$\text{sim}(\mathbf{u}_i, \mathbf{u}_j) = \frac{|\mathbf{u}_i \cap \mathbf{u}_j|}{|\mathbf{u}_i \cup \mathbf{u}_j|}$$



Motivating Example

$$\hat{y}_{ui} = \mathbf{U}_i^T \mathbf{V}_j \simeq \cos(\mathbf{U}_i, \mathbf{V}_j) \quad (\text{E.g., assuming a unit length})$$

	\mathbf{i}_1	\mathbf{i}_2	\mathbf{i}_3	\mathbf{i}_4	\mathbf{i}_5
\mathbf{u}_1	1	1	1	0	1
\mathbf{u}_2	0	1	1	0	0
\mathbf{u}_3	0	1	1	1	0
\mathbf{u}_4	1	0	1	1	1

$$\text{sim}(\mathbf{u}_1, \mathbf{u}_2) = 0.5$$

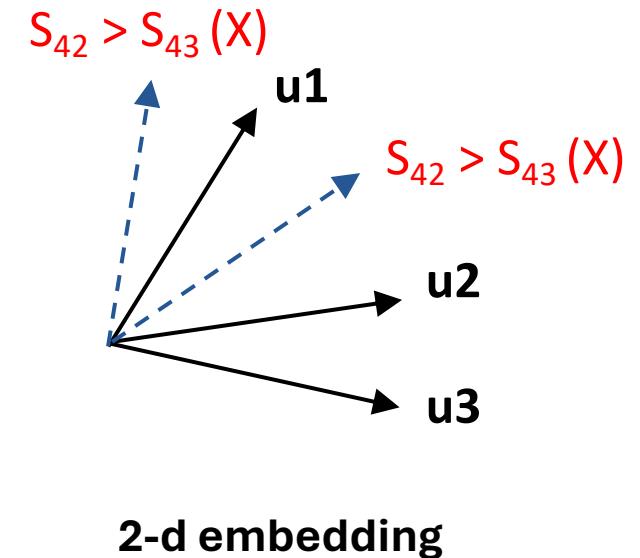
$$\text{sim}(\mathbf{u}_3, \mathbf{u}_1) = 0.4$$

$$\text{sim}(\mathbf{u}_3, \mathbf{u}_2) = 0.66$$

$$\text{sim}(\mathbf{u}_4, \mathbf{u}_1) = 0.6$$

$$\text{sim}(\mathbf{u}_4, \mathbf{u}_2) = 0.2$$

$$\text{sim}(\mathbf{u}_4, \mathbf{u}_3) = 0.4$$



$$\text{sim}(\mathbf{u}_i, \mathbf{u}_j) = \frac{|\mathbf{u}_i \cap \mathbf{u}_j|}{|\mathbf{u}_i \cup \mathbf{u}_j|}$$

Motivating Example

$$\hat{y}_{ui} = \mathbf{U}_i^T \mathbf{V}_j \simeq \cos(\mathbf{U}_i, \mathbf{V}_j) \quad (\text{E.g., assuming a unit length})$$

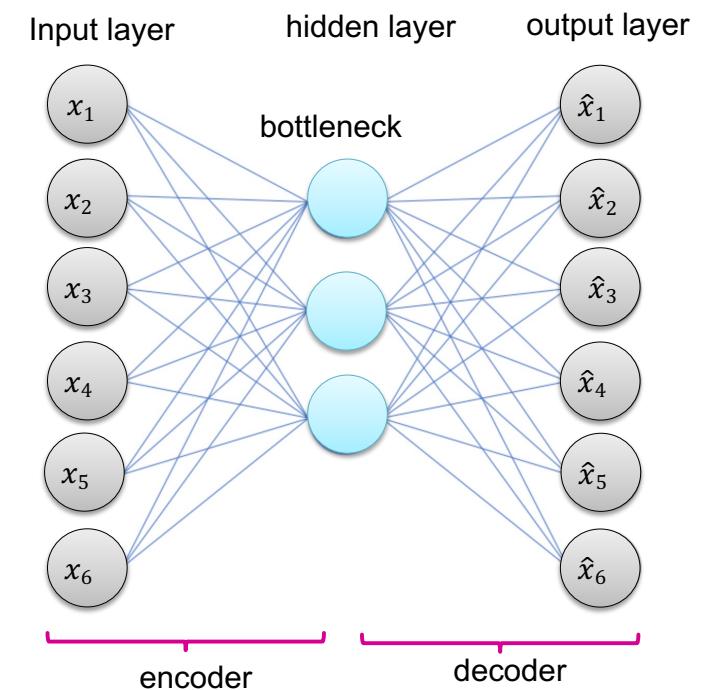
S S (x)

How to resolve?

- Increase the number of latent dimensions
 - May lead to overfitting in sparse environment
- Replacing the simple, fixed inner product with interaction function learned from the data
 - AutoRec: Autoencoders Meet Collaborative Filtering
 - NCF: Neural Collaborative Filtering

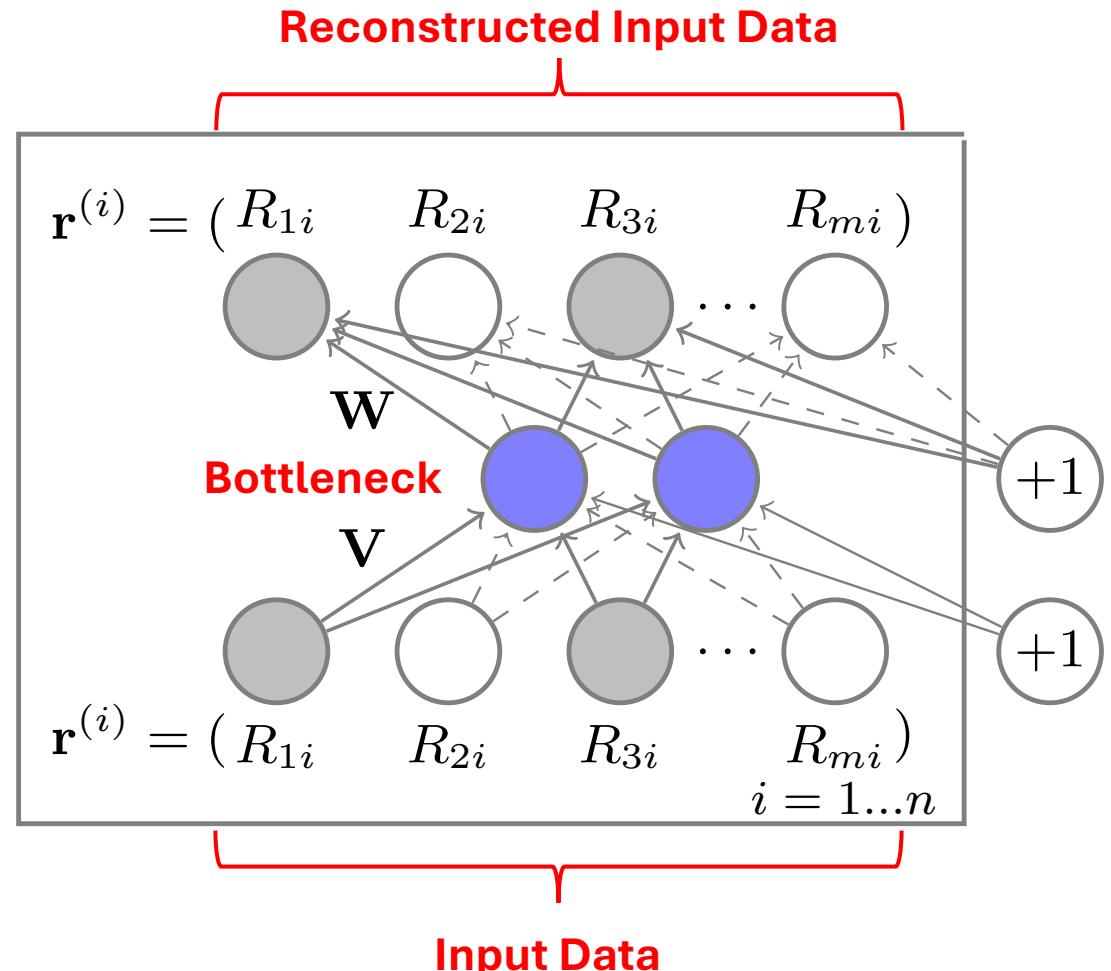
Autoencoders + CF = AutoRec

- **Autoencoders (recall lecture 3)**
 - Learning concise representation by training to reconstruct input
 - Representations can be used for downstream tasks
 - Can be used to denoise input data
 - Reconstruct data from noised input
 - input: data + noise, output: data
 - **In recommender systems:**
 - **Can we learn to ‘reconstruct’ missing rating data?**



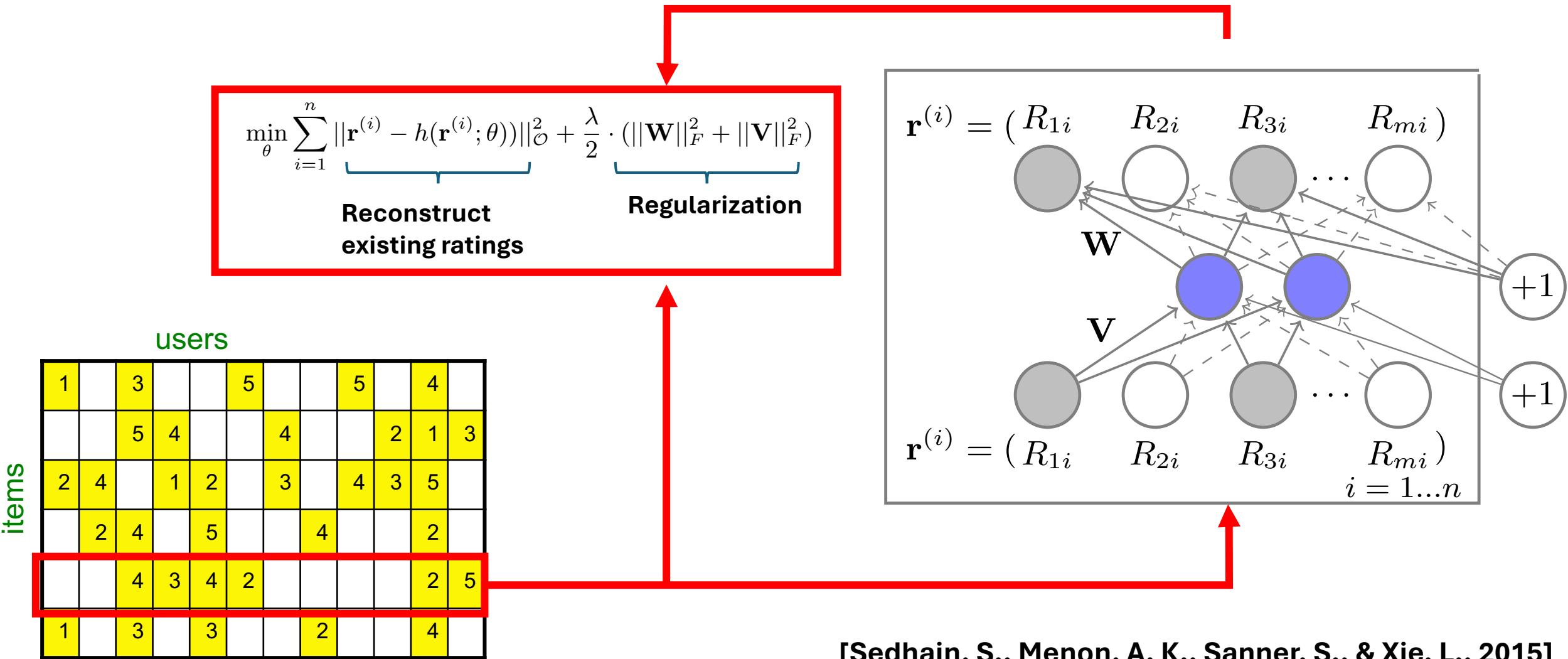
Autoencoders + CF = AutoRec

users									
items	1	2	3	4	5	6	7	8	9
1			3			5		5	4
2				5	4			4	
3						4		2	1
4							2	1	3
5	2	4		1	2		3		
6						4	3	5	
7							4		
8								2	
9									5
10	1								
11		3			3				
12						2			
13							4		
14								2	5
15	1								

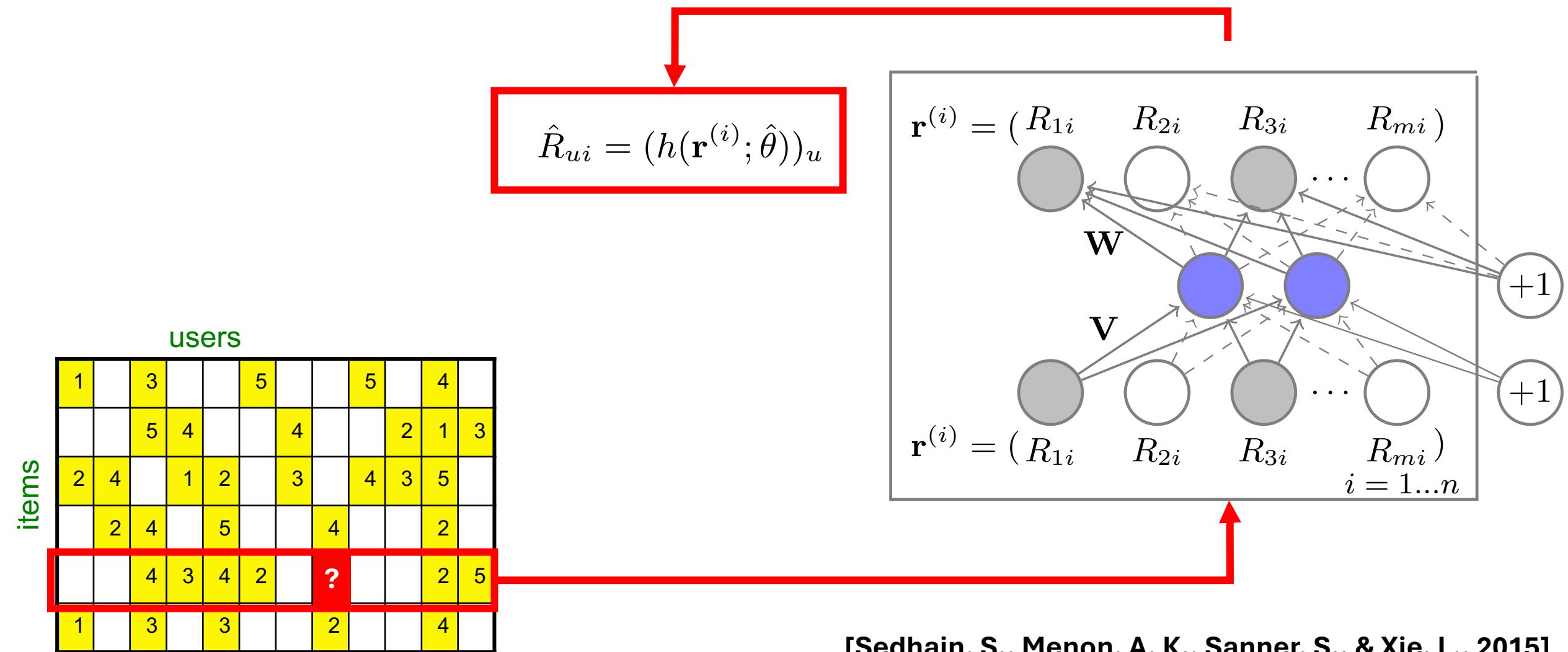


[Sedhain, S., Menon, A. K., Sanner, S., & Xie, L., 2015]

AutoRec: Training on existing ratings



AutoRec: Predicting missing ratings

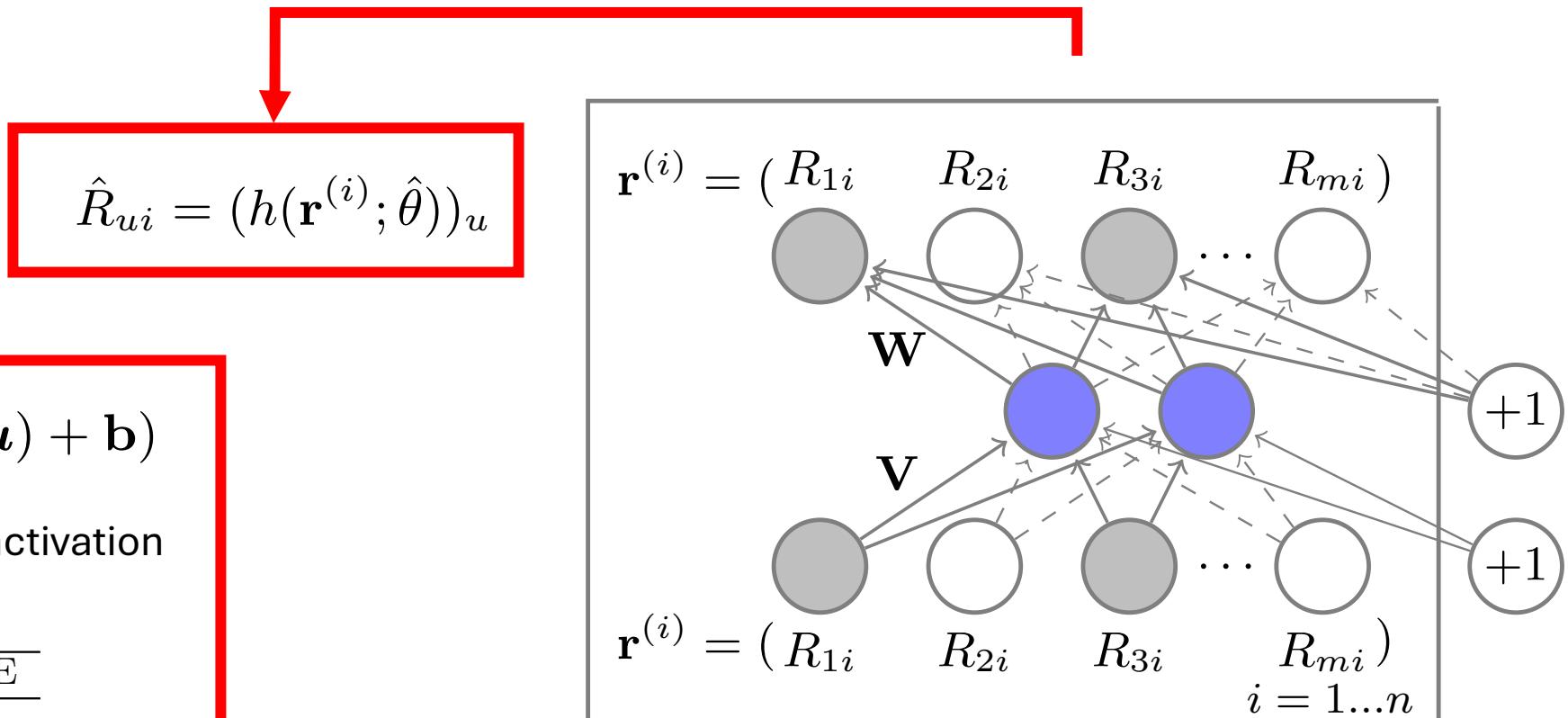


AutoRec: Predicting missing ratings

$h(\mathbf{r}; \theta) = f(\mathbf{W} \cdot g(\mathbf{V}\mathbf{r} + \boldsymbol{\mu}) + \mathbf{b})$

Choice of linear and non-linear activation functions:

$f(\cdot)$	$g(\cdot)$	RMSE
Identity	Identity	0.872
Sigmoid	Identity	0.852
Identity	Sigmoid	0.831
Sigmoid	Sigmoid	0.836



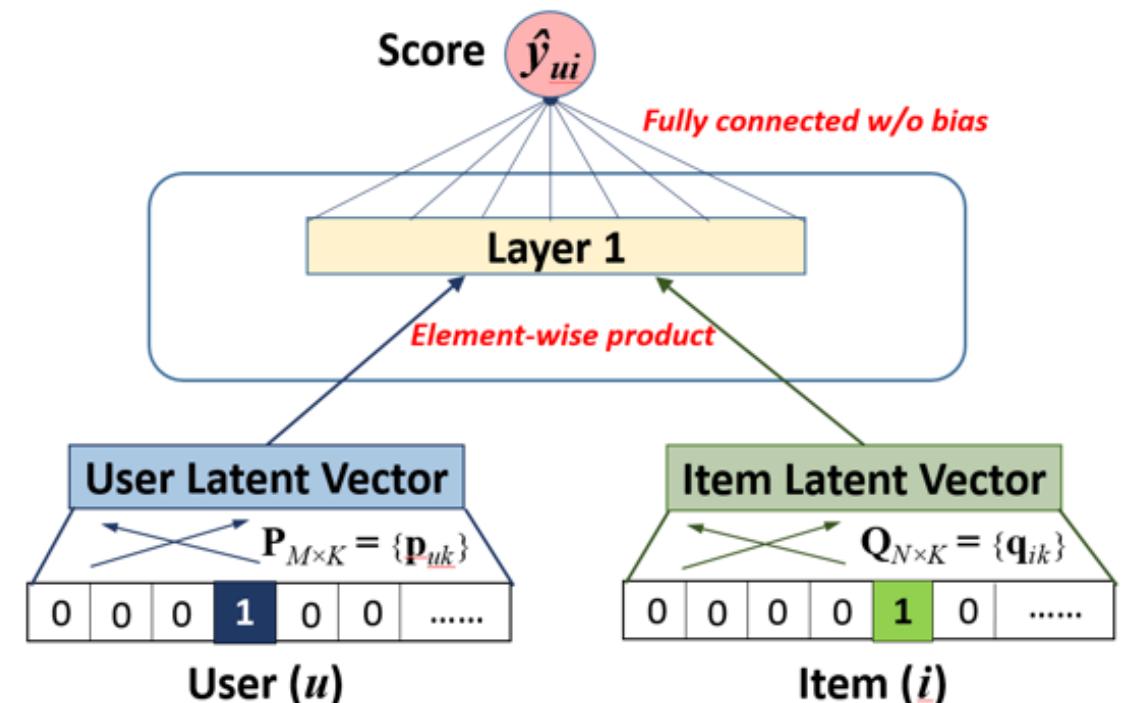
Neural Collaborative Filtering

- Three approaches:
- **Generalized Matrix Factorization (GMF)**
 - a simple neural network that generalize and extend matrix factorization
- **Multi-layer perceptron (MLP)**
 - a deep neural network that incorporates nonlinearities in an MLP architecture
- **Neural Collaborative Filtering (NCF)**
 - a framework that combines both GMF and MLP

Generalized Matrix Factorization (GMF)

- K-dimensional “Shallow” (lookup) embedding for each user and each item (similar to Funk’s SVD).
- Layer 1: element-wise product
- Output layer: fully-connected w/o bias

$$\hat{y}_{ui} = a_{out}(\mathbf{h}^T (\mathbf{p}_u \odot \mathbf{q}_i))$$

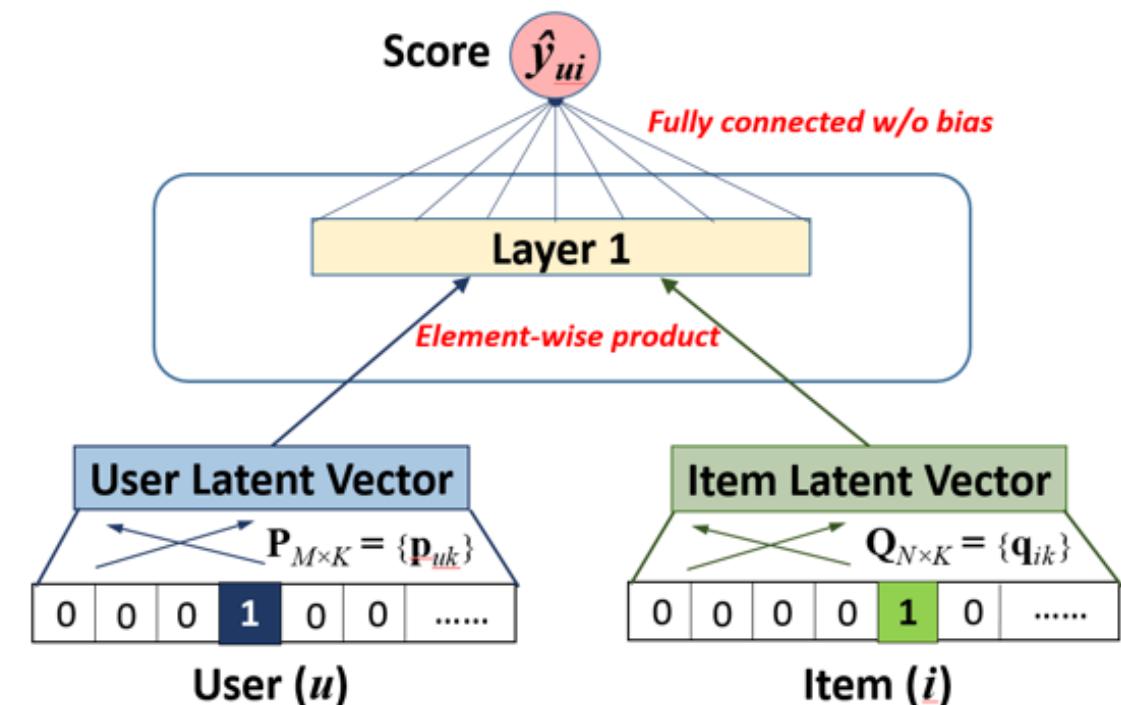


Generalized Matrix Factorization (GMF)

- K-dimensional “Shallow” (lookup) embedding for each user and each item (similar to Funk’s SVD).
- Layer 1: element-wise product
- Output layer: fully-connected w/o bias

$$\hat{y}_{ui} = a_{out}(\mathbf{h}^T (\mathbf{p}_u \odot \mathbf{q}_i))$$

Matrix Factorization
(e.g., in Funk’s SVD)
if $\mathbf{h}=[1, \dots, 1]$



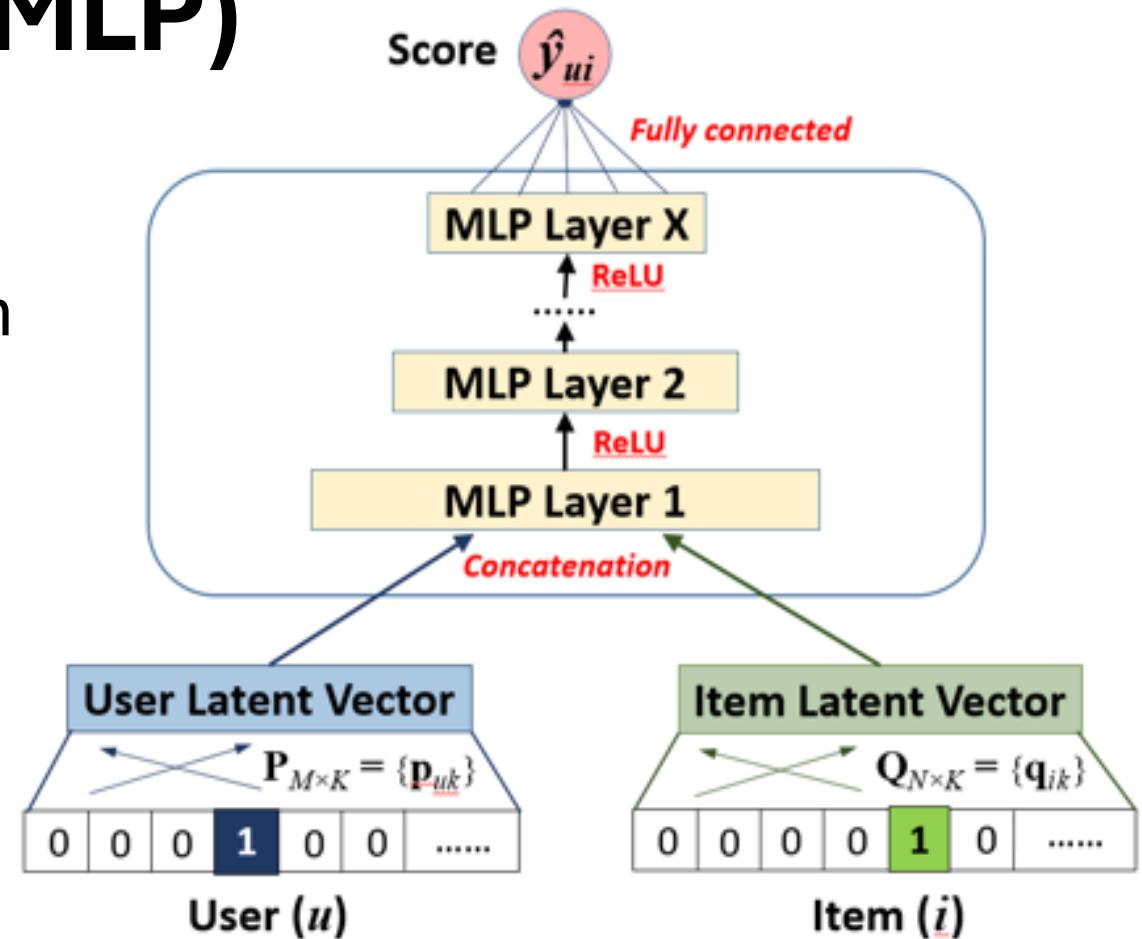
Multi-Layer Perceptron (MLP)

- K-dimensional “Shallow” (lookup) embedding for each user and each item (similar to Funk’s SVD).

- Layer 1: concatenate user-item embeddings

$$\mathbf{z}_1 = \phi_1(\mathbf{p}_u, \mathbf{q}_i) = \begin{bmatrix} \mathbf{p}_u \\ \mathbf{q}_i \end{bmatrix}$$

- Remaining layers: fully-connected w/ ReLU



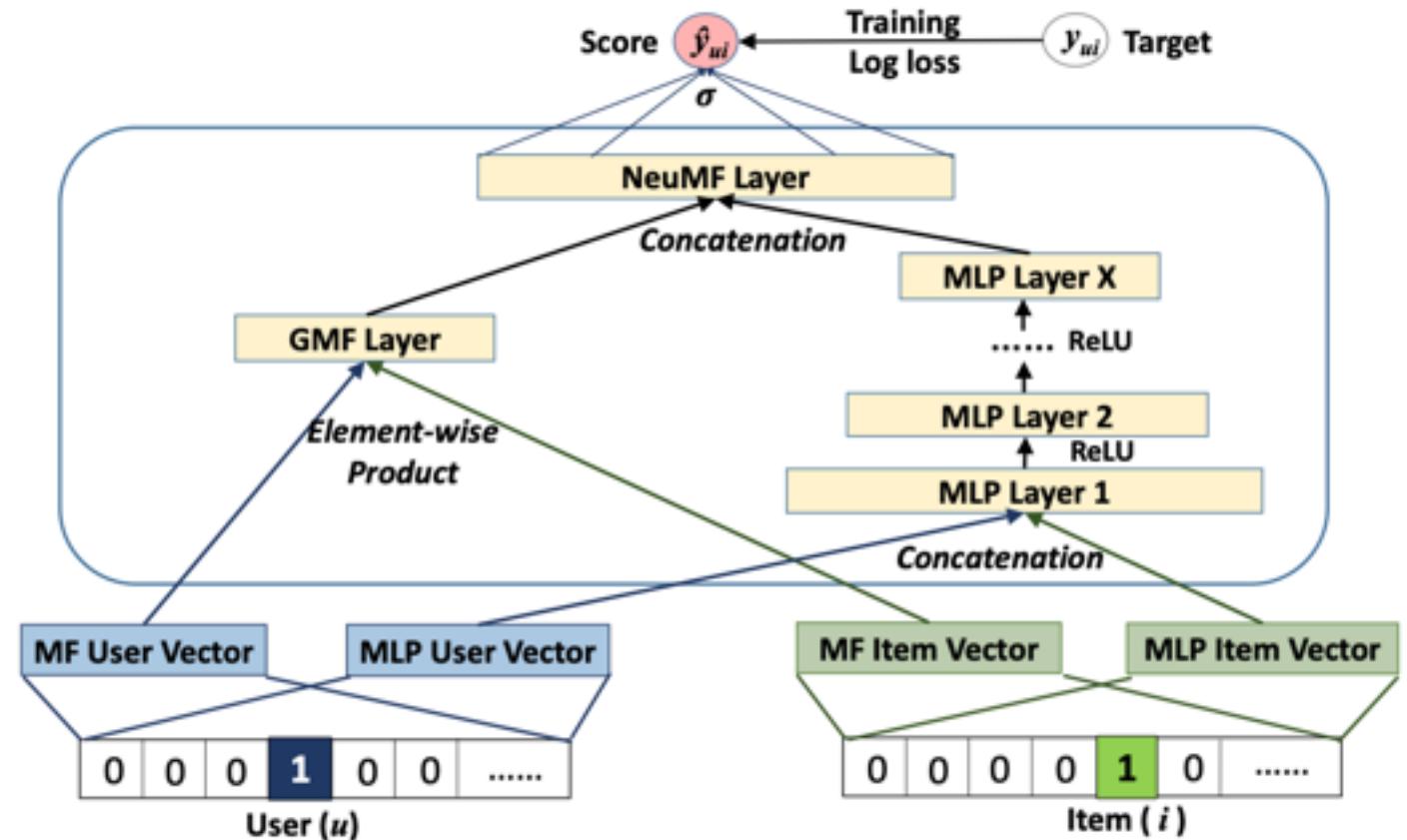
(G)MF vs. MLP

- MF uses dot (inner) product as the interaction function
 - Latent factors are independent of each other
 - Empirically, has good generalization ability for CF
- MLP uses nonlinear functions to learn interaction function
 - Latent factors are not independent of each other
 - The interaction is learned from data, which should have better representation ability
 - Risk of overfitting
- Can we combine the two models together?

Neural Collaborative Filtering (NCF)

- NCF combines GMF and MLP by allowing them to learn different sets of embeddings

- **Output:**
 - Concatenate GMF and last MLP layer
 - Fully-connected layer to output score



Training NCF Models

- For explicit feedback (e.g., ratings on 1-5 scale)

- Use regression loss:

$$L_r = \sum_{(u,i) \in \mathcal{Y}} \underbrace{(y_{ui} - \hat{y}_{ui})^2}_{\text{Accuracy}} + \lambda_\Theta \underbrace{\|\Theta\|^2}_{\text{Regularization}}$$

- For implicit feedback (e.g., watched video 0/1)
 - Use binary classification loss

- Train using SGD-based algorithm (e.g., Adam)

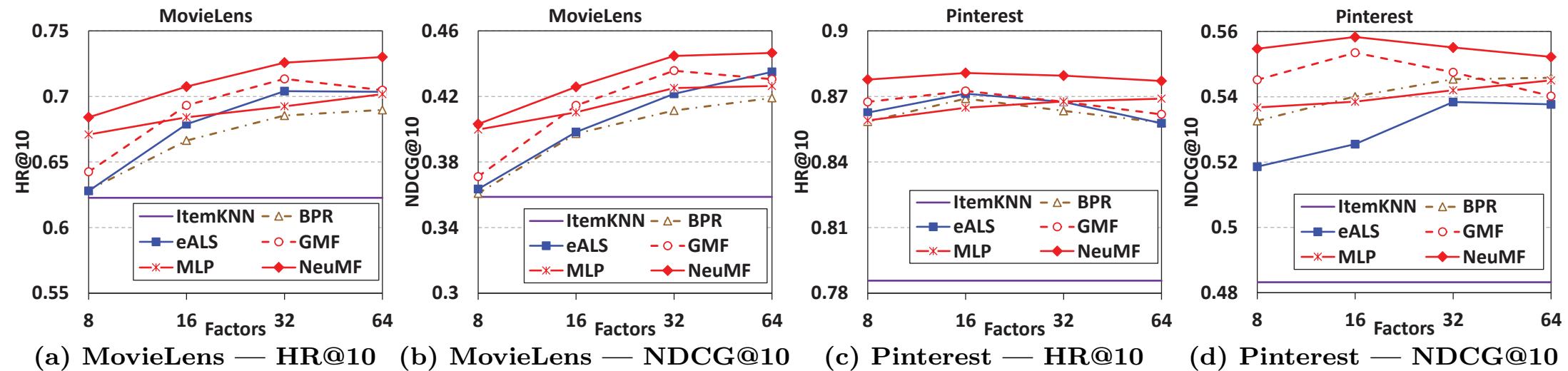
Training NCF Models

- “Tricks” and Enhancements:
 - **Pre-training:**
 - Train GMF model and MLP model on their own
 - Then initialize the weights of the combined NCF model using the weights learned independently for GMF and MLP
 - **Negative sampling for implicit feedback:**
 - Implicit feedback is trained as a binary classification
 - Typically, the large majority of interactions are (implied) negatives
 - Instead, uniformly sample negatives while controlling ratio of negatives to positives

NCF Performance

- **Two datasets from MovieLens and Pinterest:**
 - Focus on implicit feedback 0/1 (e.g., will the user watch a movie)
- **Evaluation:**
 - Leave-one-out: hold out latest rating for each user as test set
 - **Hit Ratio @10**
 - Does the test item appears in the top 10 recommendations for user?
 - **NDCG @10**
 - account for the position of “hit”: higher score for hits at top ranks

NCF Performance

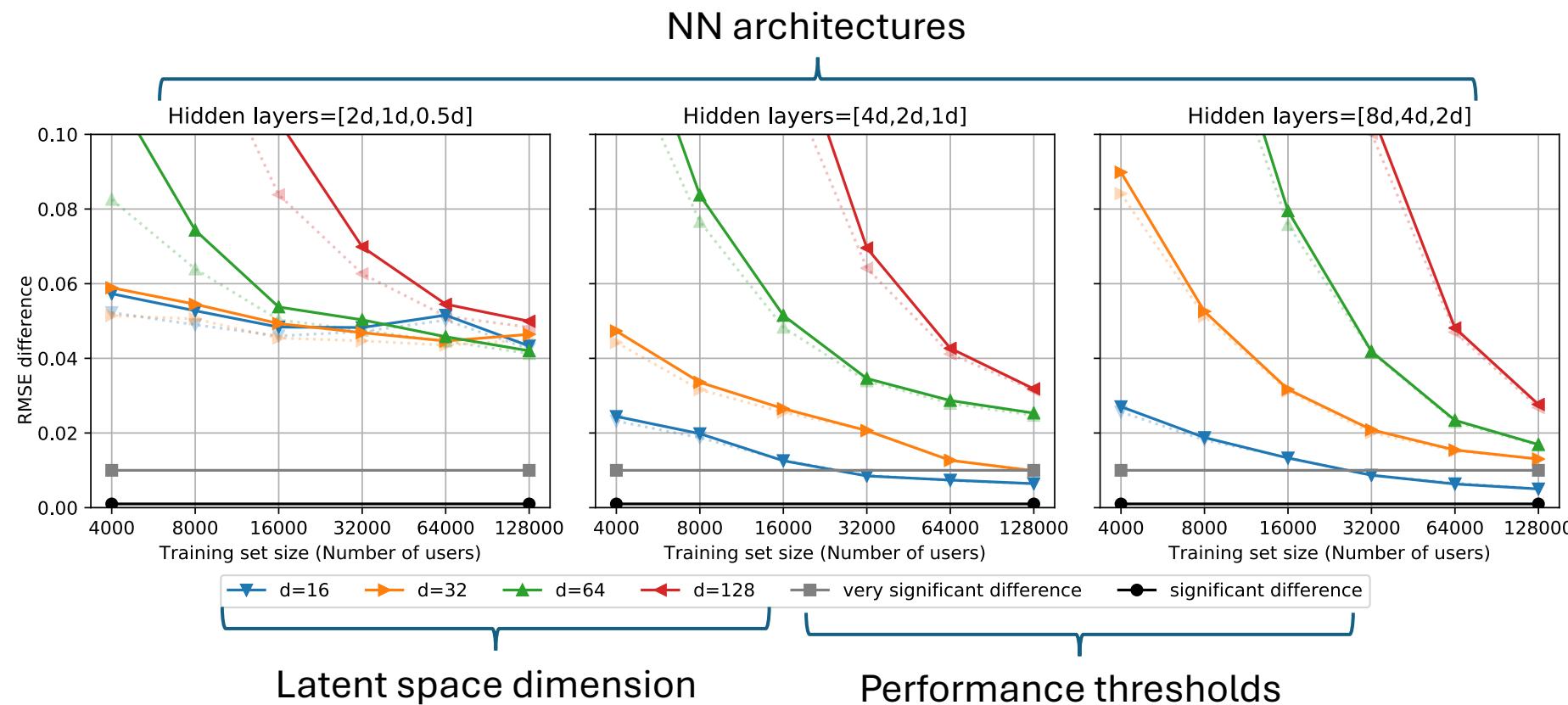


- NCF outperforms the baselines
- NCF > GMF > MLP (why GMF outperforms MLP?)

Why GMF alone outperforms MLP alone?

- #1: Overfitting
 - MLP learns more expressive interaction functions
 - Susceptible to fit noise
 - Generalize more poorly than MF
- #2: Learning dot product is hard
 - Dot product is empirically a good interaction function
 - In theory, MLP can be approximate any function. In practice, however, it is non-trivial to learn a dot product.

Learning dot product using MLP is non-trivial

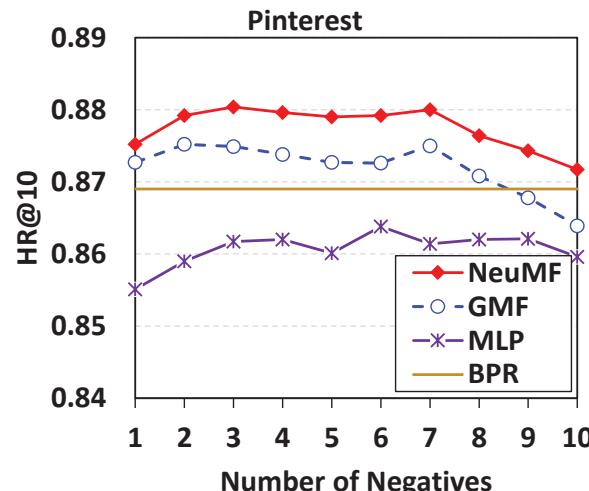


Impact of Enhancements

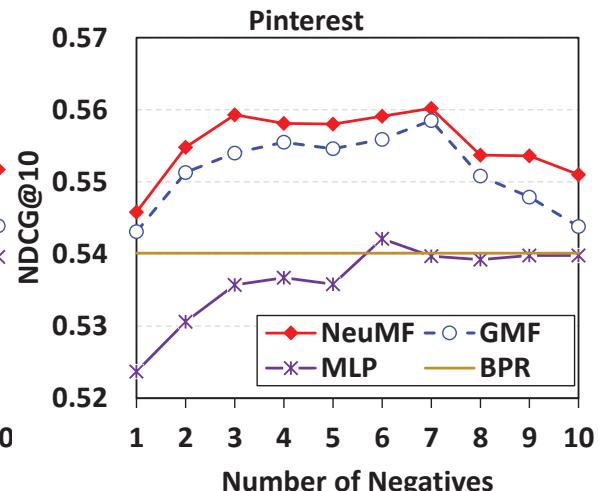
Impact of Pre-training

Factors	With Pre-training		Without Pre-training	
	HR@10	NDCG@10	HR@10	NDCG@10
MovieLens				
8	0.684	0.403	0.688	0.410
16	0.707	0.426	0.696	0.420
32	0.726	0.445	0.701	0.425
64	0.730	0.447	0.705	0.426
Pinterest				
8	0.878	0.555	0.869	0.546
16	0.880	0.558	0.871	0.547
32	0.879	0.555	0.870	0.549
64	0.877	0.552	0.872	0.551

Impact of Negative Sampling



(c) Pinterest — HR@10



(d) Pinterest — NDCG@10

Note on Reproducibility

- Neural collaborative filtering is very active research area
 - Many complex approaches (compared to the simple MF)
- A study tried to reproduce the reported results for 26 papers
 - Only 12/16 papers had reproducible setup (code + dataset)
 - “None of the computationally complex neural methods was actually consistently better than already existing learning-based techniques, e.g., using matrix factorization or linear models”
- You can test this in the assignment!

Summary

- **Neural architectures for collaborative filtering:**
 - AutoRec: autoencoder-based
 - NCF: MLP + (G)MF
- **Deeper more expressive models can improve performance**
 - Combining with MF can help generalization
- **In Assignment 5:**
 - Implement and evaluate NCF on MovieLens (explicit ratings)
 - Compare with Funk's SVD