

# DSC250: Advanced Data Mining

## Recommender Systems

**Zhiting Hu**

Lecture 16, Feb 27, 2025

**UC San Diego**

**HALICIOĞLU DATA SCIENCE INSTITUTE**

# Outline

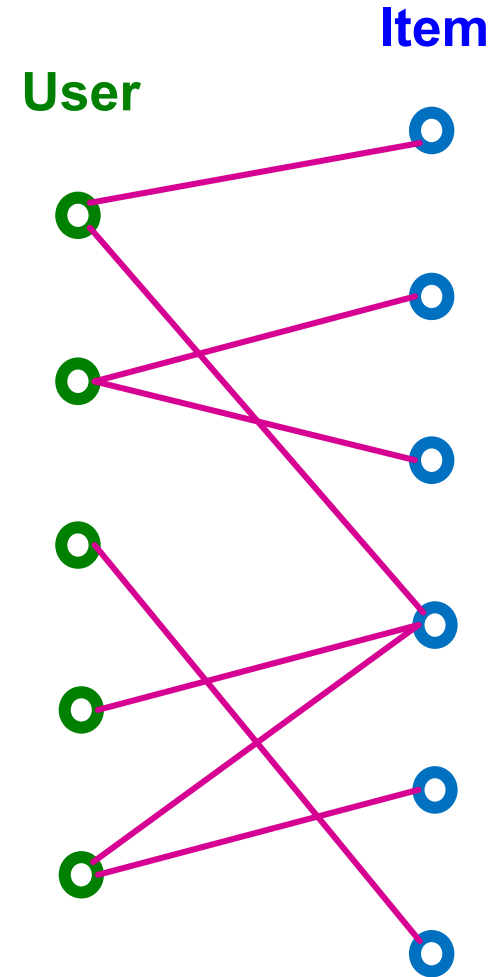
- Recommender Systems
- Presentation
  - Aryan Bansal, Yash Vishe: "Scaling Smart: Accelerating Large Language Model Pre-training with Small Model Initialization"
  - Thuy Nguyen, Taylor Martinez: "LLM Lies: Hallucinations are not Bugs, but Features as Adversarial Examples"
  - Samyak Karnavat, Derrick Yao: "Utility Engineering: Analyzing and Controlling Emergent Value Systems in AIs"
  - Har Simrat Singh, Mohammad Kianezhad: "Chain of Agents: Large Language Models Collaborating on Long-Context Tasks"

## Recap: Formulating the RecSys problem (I): Matrix Completion

Users	Movie1	Movie2	Movie3	Movie4	Movie5	Movie6	...
User1	?	?	4	?	1	?	...
User2	2	5	2	?	?	2	...
User3	?	?	5	3	2	4	...
User4	1	?	?	4	?	?	...
User5	2	3	?	?	?	?	...
...	...	...	...	...	...	...	...

# Formulating the RecSys problem (II): Link Prediction

- i Recommender system can be naturally modeled as a **bipartite graph**
  - § A graph with two node types: **users** and **items**.
  - § **Edges** connect users and items
    - § Indicates user-item interaction (e.g., click, purchase, review etc.)
    - § Often associated with timestamp (timing of the interaction).



# Formulating the RecSys problem (II): Link Prediction

## i Given

§ Past user-item interactions

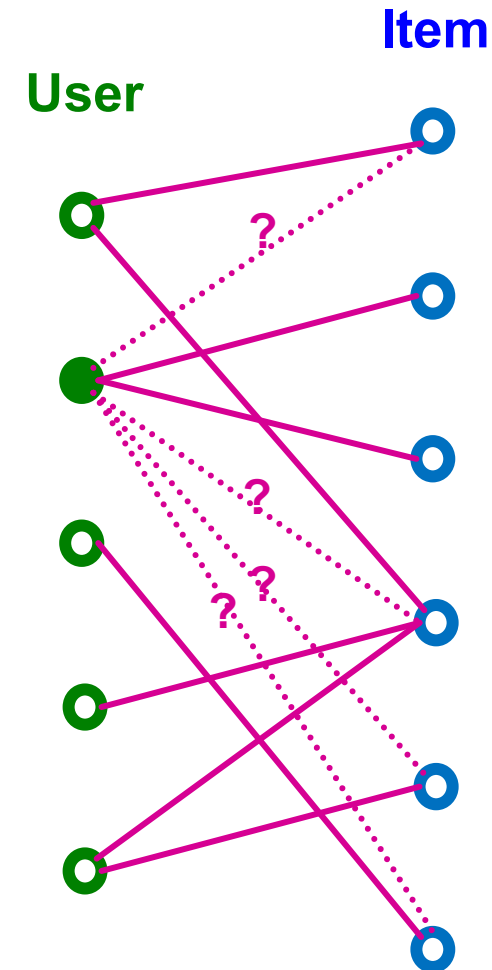
## i Task

§ Predict new items each user will interact in the future.

§ Can be cast as **link prediction** problem.

§ Predict new user-item interaction edges given the past edges.

§ For  $! \in \#$ ,  $\% \in \&$ , we need to get a real-valued **score**  $' (!, \%)$ .



# Methods

- Collaborative filtering
- Content-based recommendation
- Hybrid methods

# Methods

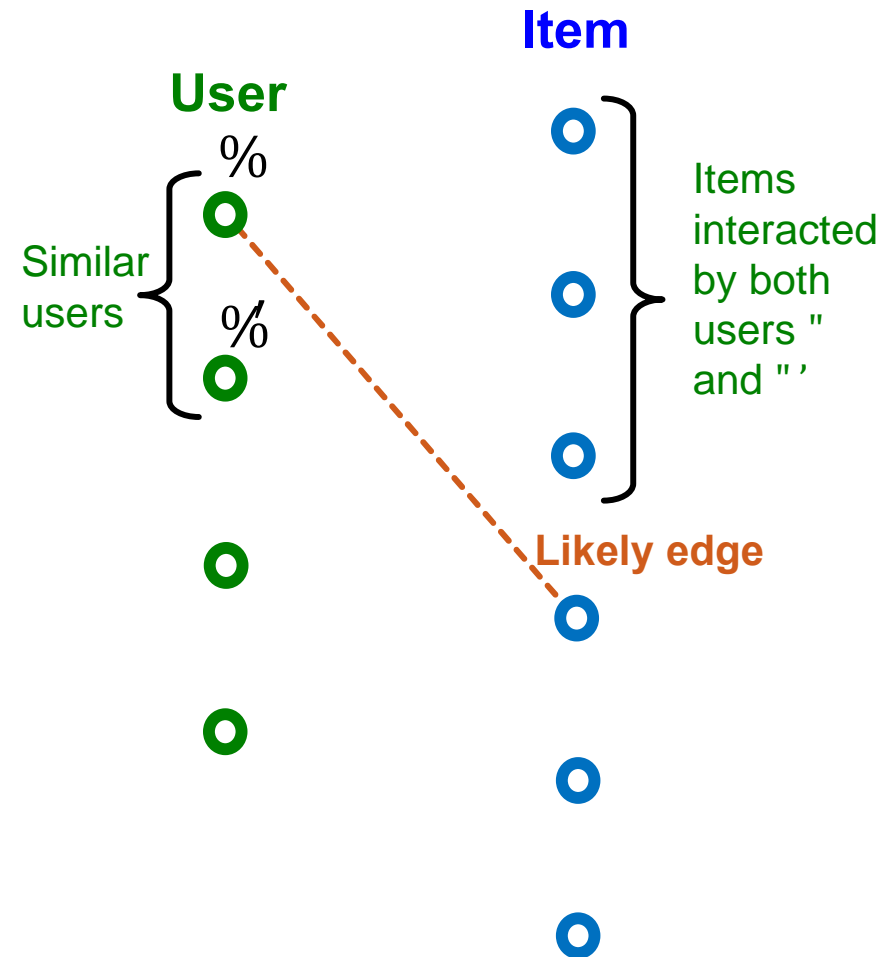
- Collaborative filtering
- Content-based recommendation
- Hybrid methods

# Collaborative Filtering (CF)

## i Underlying idea: Collaborative filtering

§ Recommend items for a user by **collecting preferences of many other similar users.**

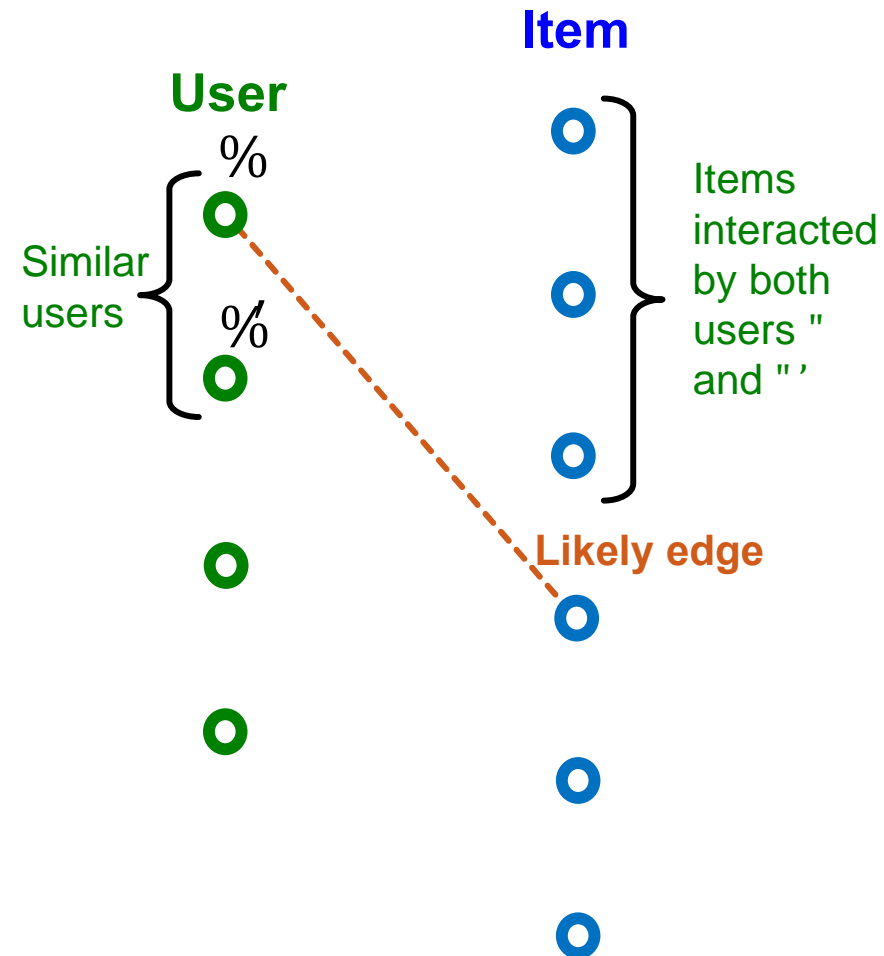
§ **Similar users tend to prefer similar items.**





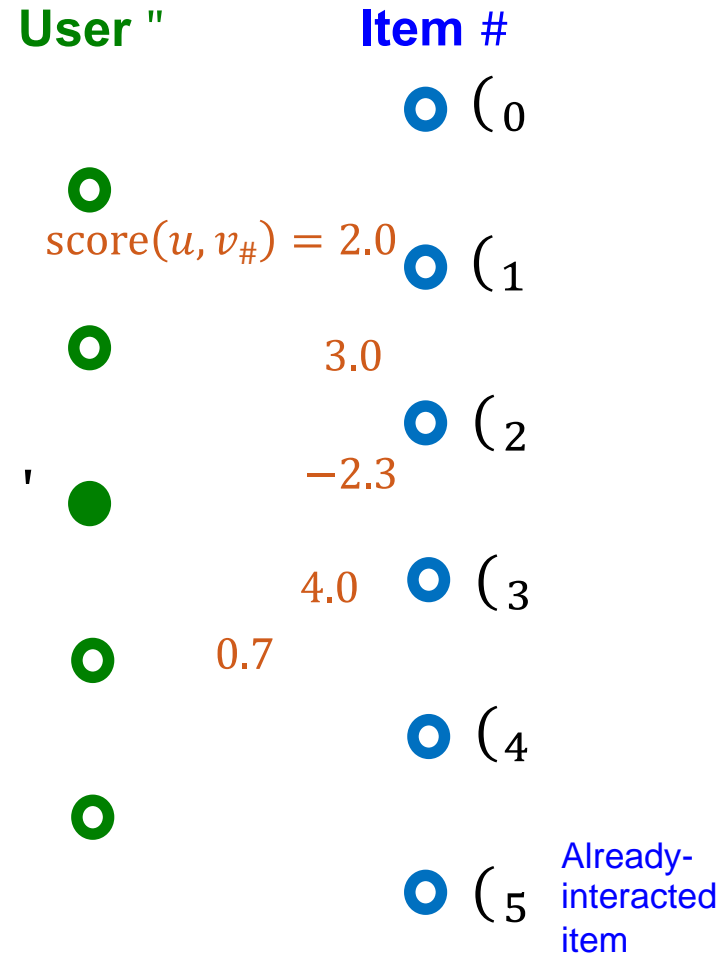
# Collaborative Filtering (CF)

- i **Underlying idea:**
  - § **Collaborative filtering**
  - § Recommend items for a user by **collecting preferences of many other similar users.**
  - § **Similar users tend to prefer similar items.**
- i **Key question: How to capture similarity between users/items?**



# Embedding-Based Models

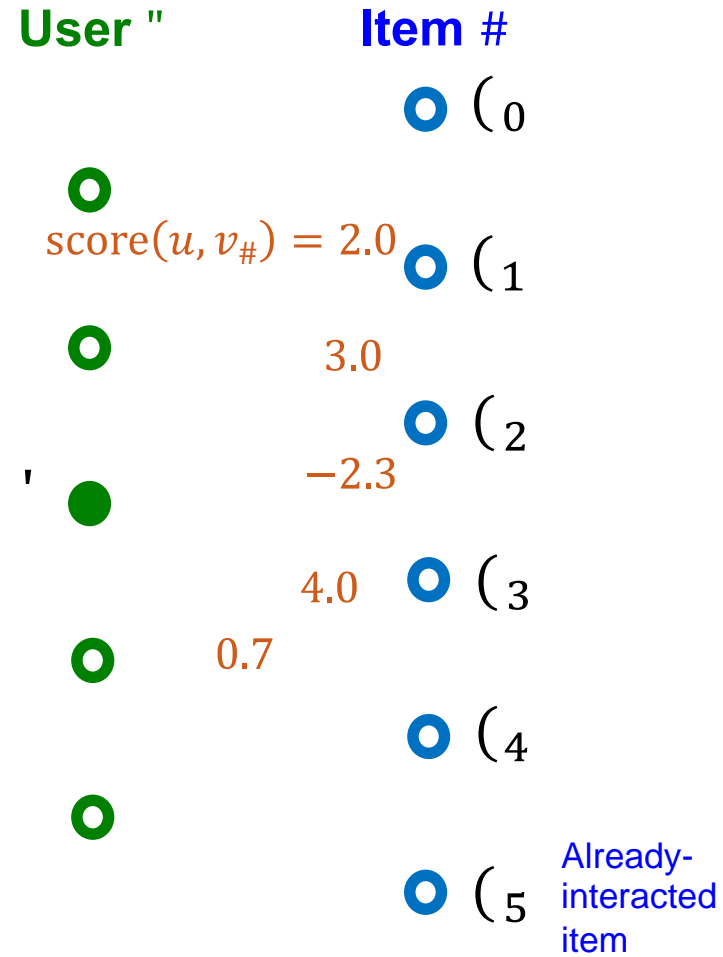
- i To get the top- $\%$  items, we need a score function for user-item interaction:
  - § For  $! \in \#$ ,  $\% \in \&$ , we need to get a real-valued scalar **score(!, %)**.



For  $K = 2$ , recommended items for user  $u$  would be  $\{v_3, v_4\}$ .

# Embedding-Based Models

- i To get the top- $\%$  items, we need a score function for user-item interaction:
  - § For  $! \in \#$ ,  $\% \in \&$ , we need to get a real-valued scalar **score**( $!$ ,  $\%$ ).
  - § \* **items with the largest scores for a given user !** (excluding **already-interacted items**) are then recommended.



For  $K = 2$ , recommended items for user  $u$  would be  $\{v_{\$}, v_{\%}\}$ .

# Embedding-Based Models

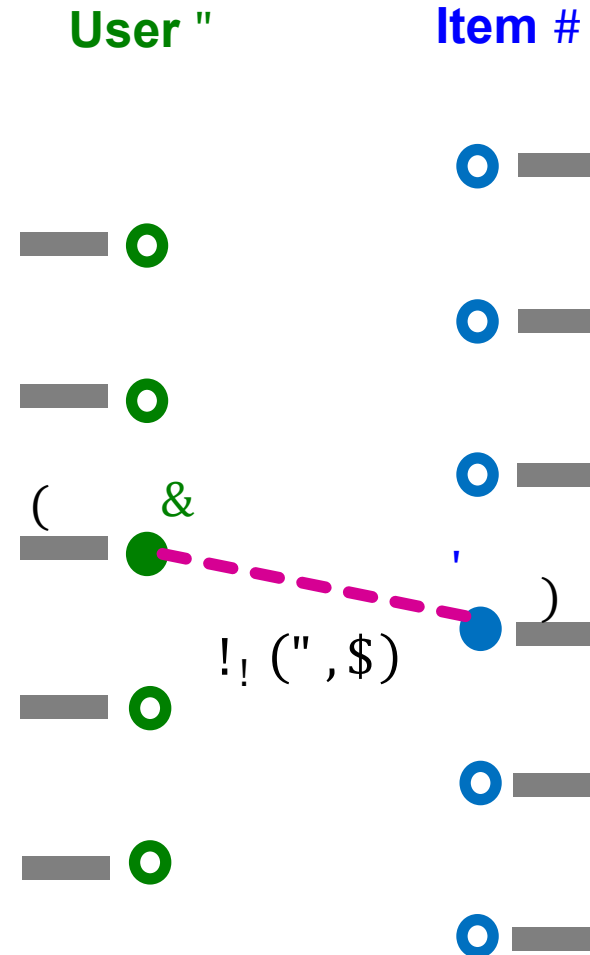
i We consider **embedding-based models** for scoring user-item interactions.

§ For each user  $u \in \mathcal{U}$ , let  $u \in \mathbb{R}^d$  be its  $d$ -dimensional embedding.

§ For each item  $i \in \mathcal{I}$ , let  $i \in \mathbb{R}^d$  be its  $d$ -dimensional embedding.

§ Let  $f_{\theta}(\cdot, \cdot): \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$  be a parametrized function.

§ Then,  $\text{score}(u, i) \equiv f_{\theta}(u, i)$



# Embedding-Based Models: Training Objective

- i Embedding-based models have three kinds of parameters:
  - § An encoder to generate user embeddings  $\{u_i\}_{i \in \mathcal{U}}$
  - § An encoder to generate item embeddings  $\{v_j\}_{j \in \mathcal{I}}$
  - § Score function  $f(u_i, v_j)$

# Embedding-Based Models: Training Objective

- i Embedding-based models have three kinds of parameters:
  - § An encoder to generate user embeddings  $\{; \}$ ,  $\epsilon_+$
  - § An encoder to generate item embeddings  $\{>\}$ ,  $\epsilon_-$
  - § Score function  $'\%(\cdot, \cdot)$

# Embedding-Based Models: Training Objective

- i Embedding-based models have three kinds of parameters:
  - § An encoder to generate user embeddings  $\{u_i\}_{i \in \mathcal{U}}$
  - § An encoder to generate item embeddings  $\{v_j\}_{j \in \mathcal{I}}$
  - § Score function  $f(u_i, v_j)$
- i **Training objective:** Optimize the model parameters to achieve high “**accuracy**” on *seen (i.e., training)* user-item interactions
  - We hope this objective would lead to high accuracy on *unseen (i.e., test)* interactions.

# Why Embedding Models Work?

- i Embedding-based models can capture similarity of users/items!
  - § **Low-dimensional embeddings *cannot* simply memorize all user-item interaction data.**
  - § Embeddings are forced to **capture similarity between users/items to fit the data.**
  - § This allows the models to make effective prediction on *unseen* user-item interactions.



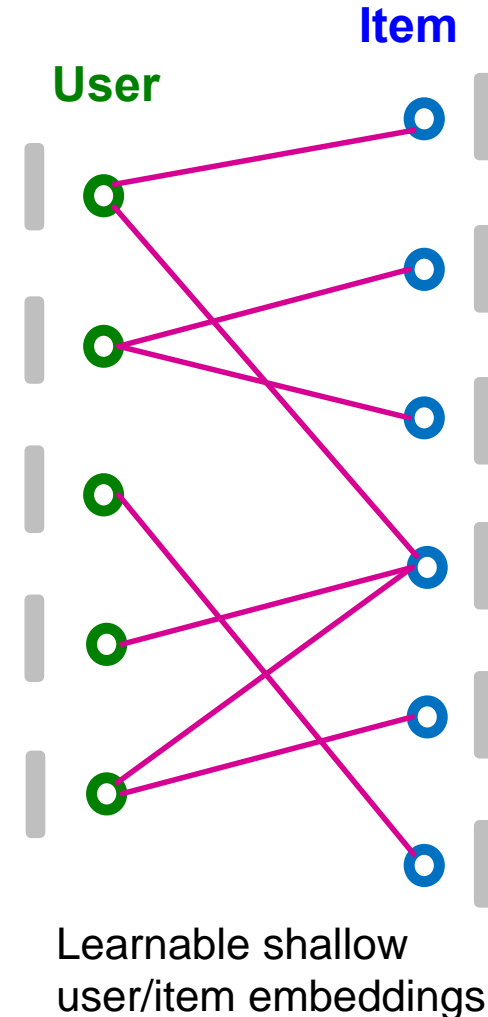
# Conventional Embedding-based CF

- Conventional collaborative filtering model is based on **shallow encoders**:

- § Use shallow encoders for users and items:

- § For every  $u \in \mathcal{U}$  and  $i \in \mathcal{I}$ , we prepare shallow learnable embeddings  $u, i \in \mathbb{R}^Q$ .

- § Score function for user  $u$  and item  $i$  is  $s(u, i) \equiv K_u^T K_i$ .



## Limitations of Shallow Encoders

- i The model itself does *not explicitly* capture graph structure
- § The graph structure is *only implicitly* captured in the training objective.

## Limitations of Shallow Encoders

- i The model itself does *not explicitly* capture graph structure
  - § The graph structure is *only implicitly* captured in the training objective.
- i Only the **first-order graph structure** (i.e., edges) is captured in the training objective.
  - § **High-order graph structure** (e.g., +-hop paths between two nodes) is *not explicitly captured*.

## We want a model that ...

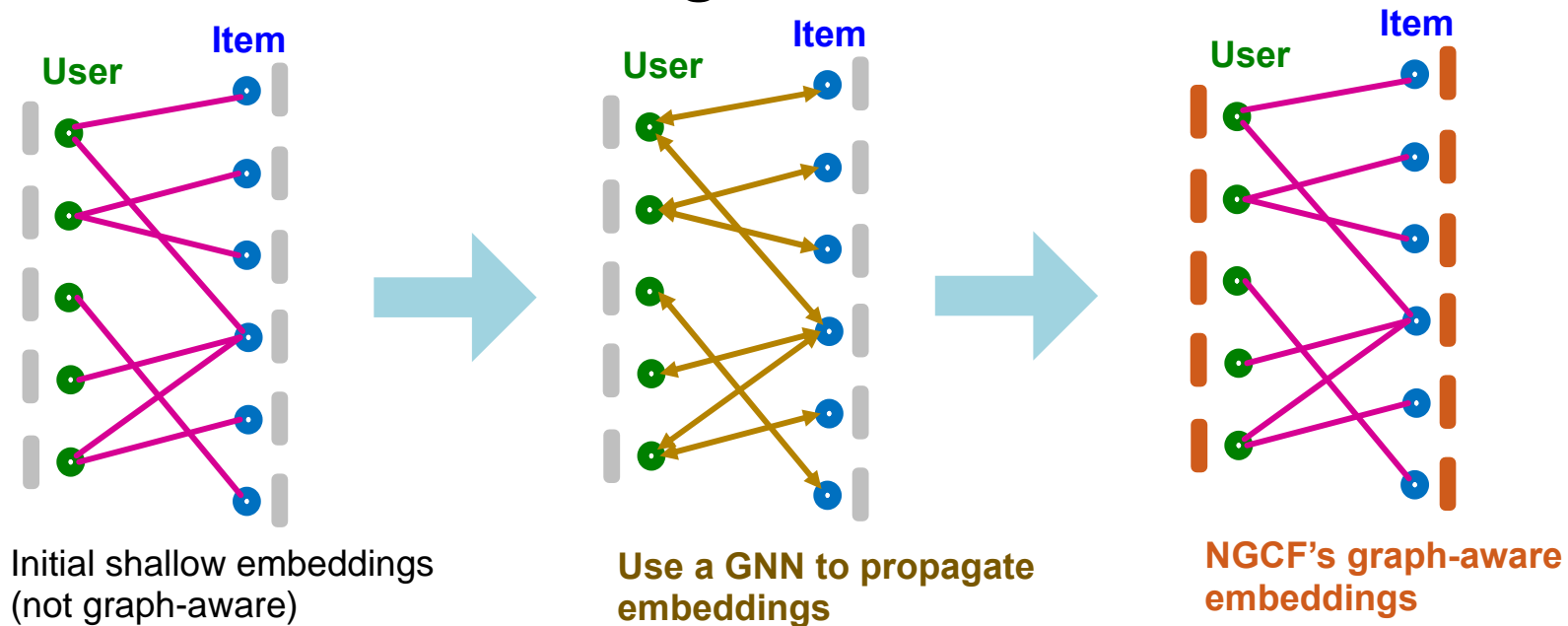
- i We want a model that...
  - § **explicitly captures graph structure** (beyond implicitly through the training objective)
  - § captures **high-order graph structure** (beyond the first-order edge connectivity structure)

## We want a model that ...

- i We want a model that...
  - § **explicitly captures graph structure** (beyond implicitly through the training objective)
  - § captures **high-order graph structure** (beyond the first-order edge connectivity structure)
- i **GNNs are a natural approach to achieve both!**
  - § **Neural Graph Collaborative Filtering (NGCF)** [Wang et al. 2019]
  - § **LightGCN** [He et al. 2020]
    - § A simplified and improved version of NGCF

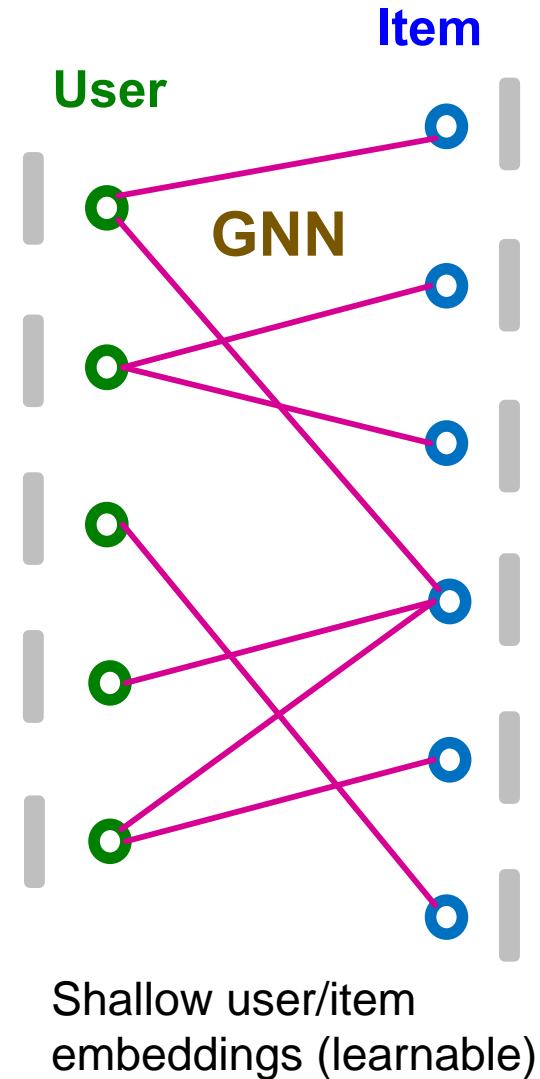
# NGCF: Overview

- Neural Graph Collaborative Filtering (NGCF) *explicitly* incorporates high-order graph structure when generating user/item embeddings.
- Key idea:** Use a GNN to generate graph-aware user/item embeddings.



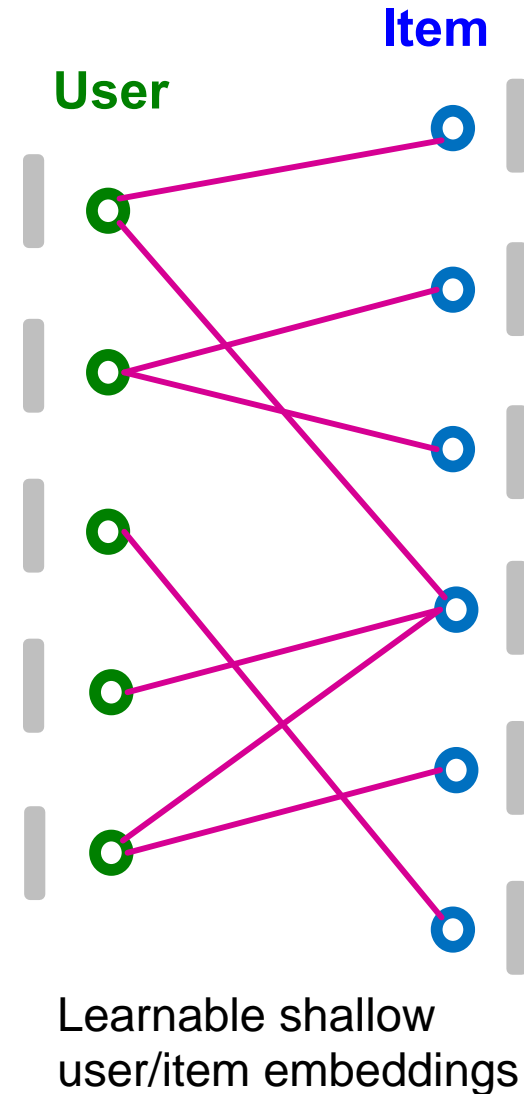
# NGCF

- Given: User-item bipartite graph.
- NGCF framework:
  - Prepare shallow learnable embedding for each node.
  - Use multi-layer GNNs to propagate embeddings along the bipartite graph.
    - High-order graph structure is captured.
  - Final embeddings are *explicitly* graph-aware!
- Two kinds of learnable params are jointly learned:
  - Shallow user/item embeddings
  - GNN's parameters



# NGCF: Initial Node Embeddings

- i Set the shallow learnable embeddings as the initial node features:
  - § For every user  $! \in \#$ , set  $\$!^{(\#)}$  as the user's shallow embedding.
  - § For every item  $\% \in \&$ , set  $\$\%^{(\#)}$  as the item's shallow embedding.





# NGCF: Neighbor Aggregation

- i Iteratively update node embeddings using neighboring embeddings.

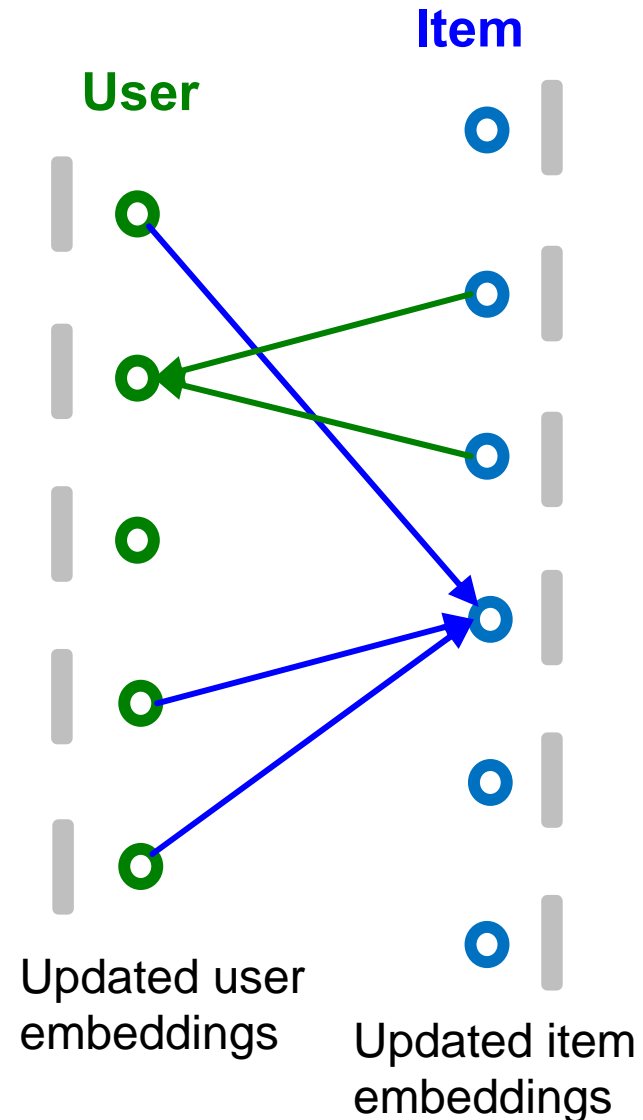
$$h_u^{(45)} = \text{COMBINE} \left( h_u^{(4)}, \text{AGGR} \left( \{h_i^{(4)}\}_{i \in \mathcal{N}(u)} \right) \right)$$

$$h_i^{(45)} = \text{COMBINE} \left( h_i^{(4)}, \text{AGGR} \left( \{h_u^{(4)}\}_{u \in \mathcal{N}(i)} \right) \right)$$

High-order graph structure is captured through iterative neighbor aggregation.

Different architecture choices are possible for AGGR and COMBINE.

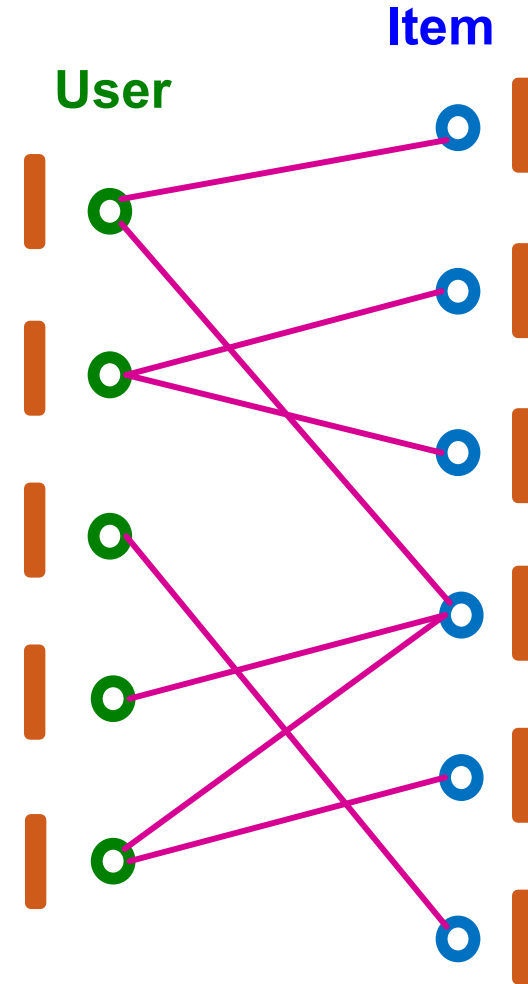
- AGGR( $\cdot$ ) can be MEAN( $\cdot$ )
- COMBINE(P, Q) can be  $\text{ReLU}(\text{Linear}(\text{Concat}(P, Q)))$



# NGCF: Final Embeddings and Score Function

- After  $\rho$  rounds of neighbor aggregation, we get the **final user/item embeddings**  $u_i^{(\rho)}$  and  $v_j^{(\rho)}$ .
- For all  $u \in U, v \in V$ , we set  $\hat{u} \leftarrow u_i^{(\rho)}, \hat{v} \leftarrow v_j^{(\rho)}$ .
- Score function is the inner product

$$\text{score}(u, v) = \hat{u}^T \hat{v}$$



Final user/item embeddings (graph-aware)

## NGCF: Summary

- i Conventional collaborative filtering uses shallow user/item embeddings.
  - § The embeddings do **not explicitly model graph structure**.
  - § The training objective **does not model high-order graph structure**.
- i **NGCF uses a GNN to propagate the shallow embeddings.**
  - § The embeddings are **explicitly aware of high-order graph structure**.

# Issues of Collaborative Filtering

- **Cold Start**: There needs to be enough other users already in the system to find a match.
- **Sparsity**: If there are many items to be recommended, even if there are many users, the user/ratings matrix is sparse, and it is hard to find users that have rated the same items.
- **First Rater**: Cannot recommend an item that has not been previously rated.
  - New items
  - Esoteric items
- **Popularity Bias**: Cannot recommend items to someone with unique tastes.
  - Tends to recommend popular items.

# Methods

- Collaborative filtering
- **Content-based recommendation**
- Hybrid methods

# Content-based Recommendation

- Collaborative filtering does **NOT** require any information about content,
  - However, it might be reasonable to exploit such information
  - E.g. recommend fantasy novels to people who liked fantasy novels in the past
- What do we need:
  - Information about the available items such as the genre ("content")
  - *user profile* describing what the user likes (the preferences)
- The task:
  - Learn user preferences
  - Locate/recommend items that are "similar" to the user preferences

# Content-based Recommendation

User profile

Title	Genre	Author	Type	Price	Keywords
The Night of the Gun	Memoir	David Carr	Paperback	29.90	Press and journalism, drug addiction, personal memoirs, New York
The Lace Reader	Fiction, Mystery	Brunonia Barry	Hardcover	49.90	American contemporary fiction, detective, historical
Into the Fire	Romance, Suspense	Suzanne Brockmann	Hardcover	45.90	American fiction, Murder, Neo-nazism
...					

Item

Title	Genre	Author	Type	Price	Keywords
...	Fiction, Suspense	Brunonia Barry, Ken Follet, ..	Paperback	25.65	detective, murder, New York

- Simple approach
  - Compute the similarity of an unseen item with the user profile based on the keyword overlap (e.g. using the Dice coefficient)
  - $\text{sim}(b_i, b_j) = \frac{2 * |\text{keywords}(b_i) \cap \text{keywords}(b_j)|}{|\text{keywords}(b_i)| + |\text{keywords}(b_j)|}$
- Other advanced similarity measure

# Methods

- Collaborative filtering
- Content-based recommendation
- **Hybrid methods**
  - Combining both user-item interaction and other external sources of information
  - E.g., Factorization Machines



**Questions?**