# Stanford CS224W: GNNs for Recommender Systems

CS224W: Machine Learning with Graphs
Jure Leskovec, Stanford University
http://cs224w.stanford.edu

# Stanford CS224W: Recommender Systems: Task and Evaluation

CS224W: Machine Learning with Graphs
Jure Leskovec, Stanford University
http://cs224w.stanford.edu
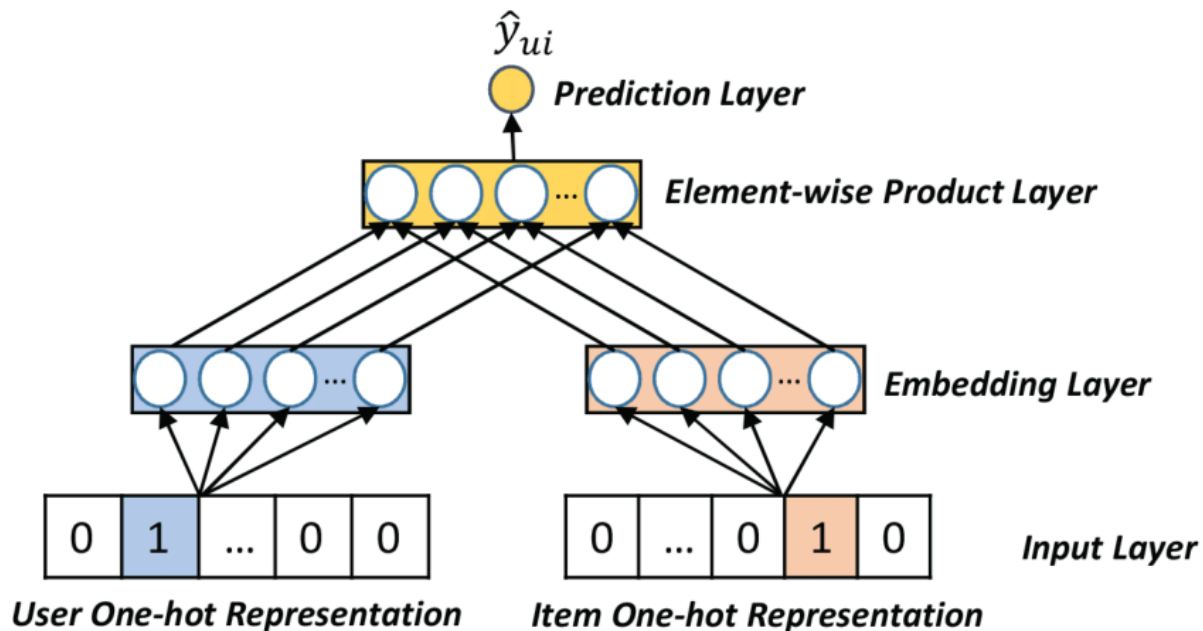
# Preliminary of Recommendation

- **Information Explosion in the era of Internet**
  - 10K+ movies in Netflix
  - 12M products in Amazon
  - 70M+ music tracks in Spotify
  - 10B+ videos on YouTube
  - 200B+ pins (images) in Pinterest
- **Personalized recommendation (i.e., suggesting a small number of interesting items for each user)** is critical for users to effectively explore the content of their interest.
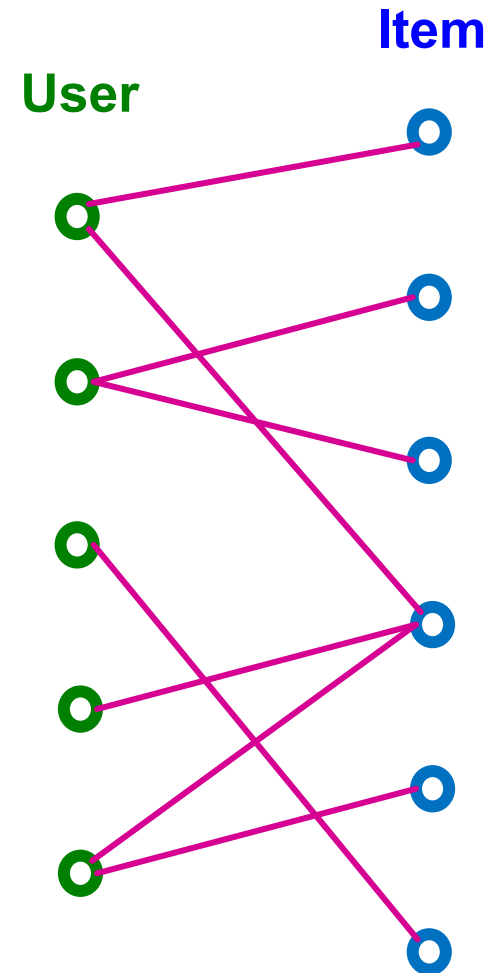
# Matrix Factorization



The embeddings are learned such that their dot product is a good approximation of the user-item matrix

# Matrix-factorization as a shallow neural network model

# Recommender System as a Graph

- 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.)
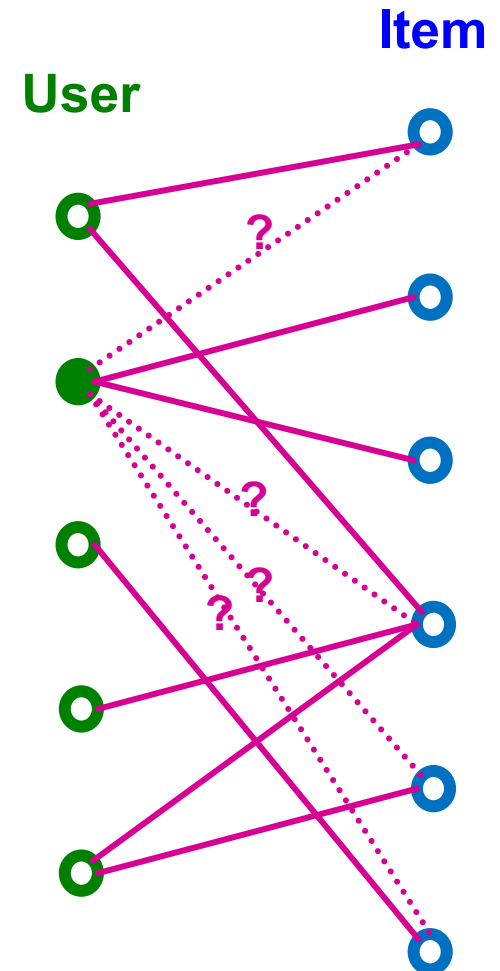    - Often associated with timestamp (timing of the interaction).



**User**   **Item**

# Recommendation Task

- **Given**
  - Past user-item interactions
- **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.

# Stanford CS224W: Recommender Systems: Embedding-Based Models

CS224W: Machine Learning with Graphs
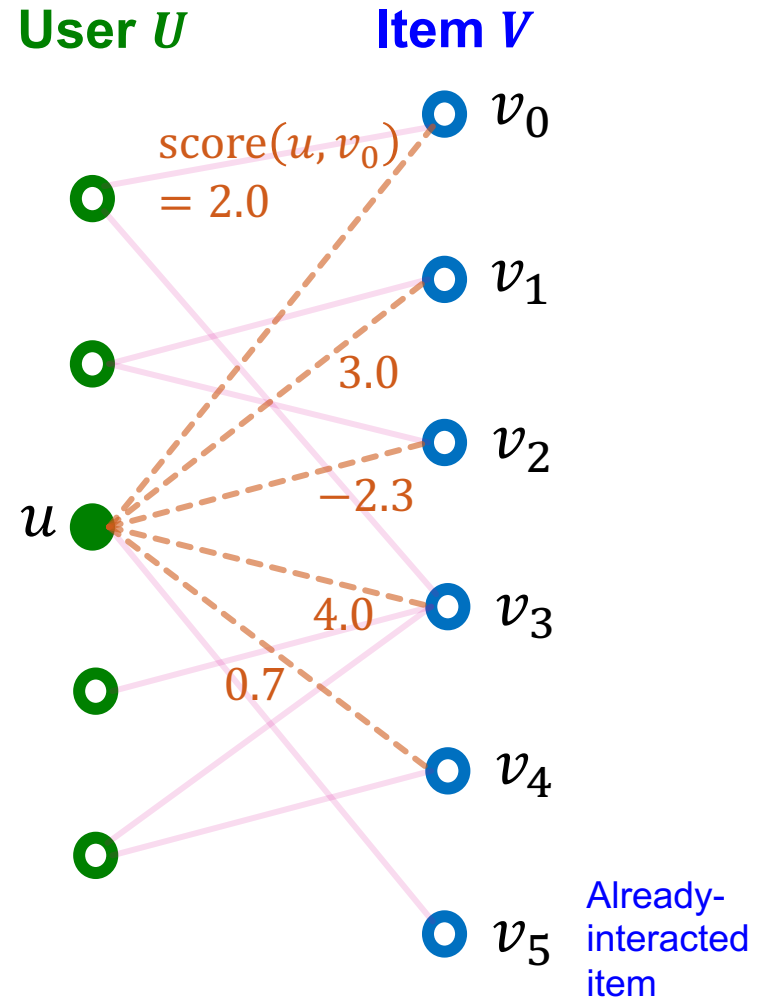Jure Leskovec, Stanford University
http://cs224w.stanford.edu

# Notation

- **Notation**:

  - $U$: A set of all users

  - $V$: A set of all items

  - $E$: A set of observed user-item interactions

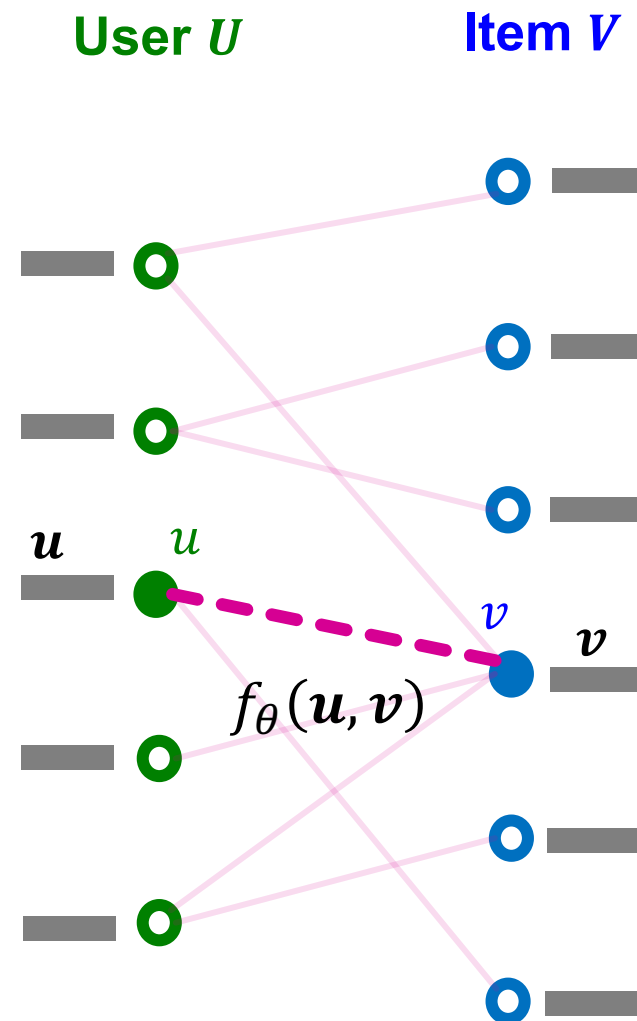    - $E = \{(u, v) \mid u \in U, v \in V, u \text{ interacted with } v\}$

# Score Function

- To get the top-$K$ items, we need a score function for user-item interaction:

  - For $u \in U, v \in V$, we need to get a real-valued scalar **score**$(u, v)$.

  - **$K$ items with the largest scores for a given user $u$** (excluding already-interacted items) are then recommended.



**User $U$**          **Item $V$**

score$(u, v_0)$ = 2.0

3.0

−2.3

4.0

0.7

Already-interacted item

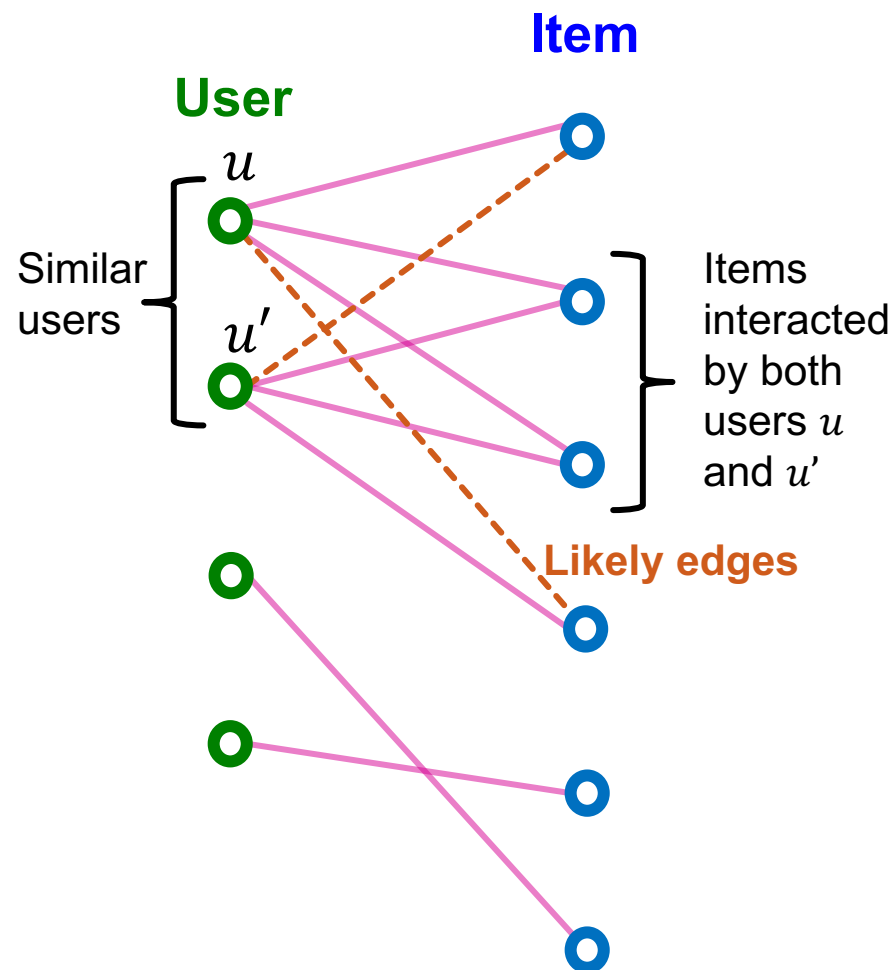For $K = 2$, recommended items for user $u$ would be $\{v_1, v_3\}$.

# Embedding-Based Models

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

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

  - For each item $v \in V$, let $\boldsymbol{v} \in \mathbb{R}^D$ be its $D$-dimensional embedding.

  - Let $f_\theta(\cdot,\cdot): \mathbb{R}^D \times \mathbb{R}^D \to \mathbb{R}$ be a parametrized function.

  - Then, $\boxed{\text{score}(u,v) \equiv f_\theta(\boldsymbol{u}, \boldsymbol{v})}$

**User $U$**    **Item $V$**



$f_\theta(\boldsymbol{u}, \boldsymbol{v})$

# Why Embedding Models Work?

- **Underlying idea: Collaborative filtering**

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

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

- **Key question: How to capture similarity between users/items?**

# Why Embedding Models Work?

- 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.
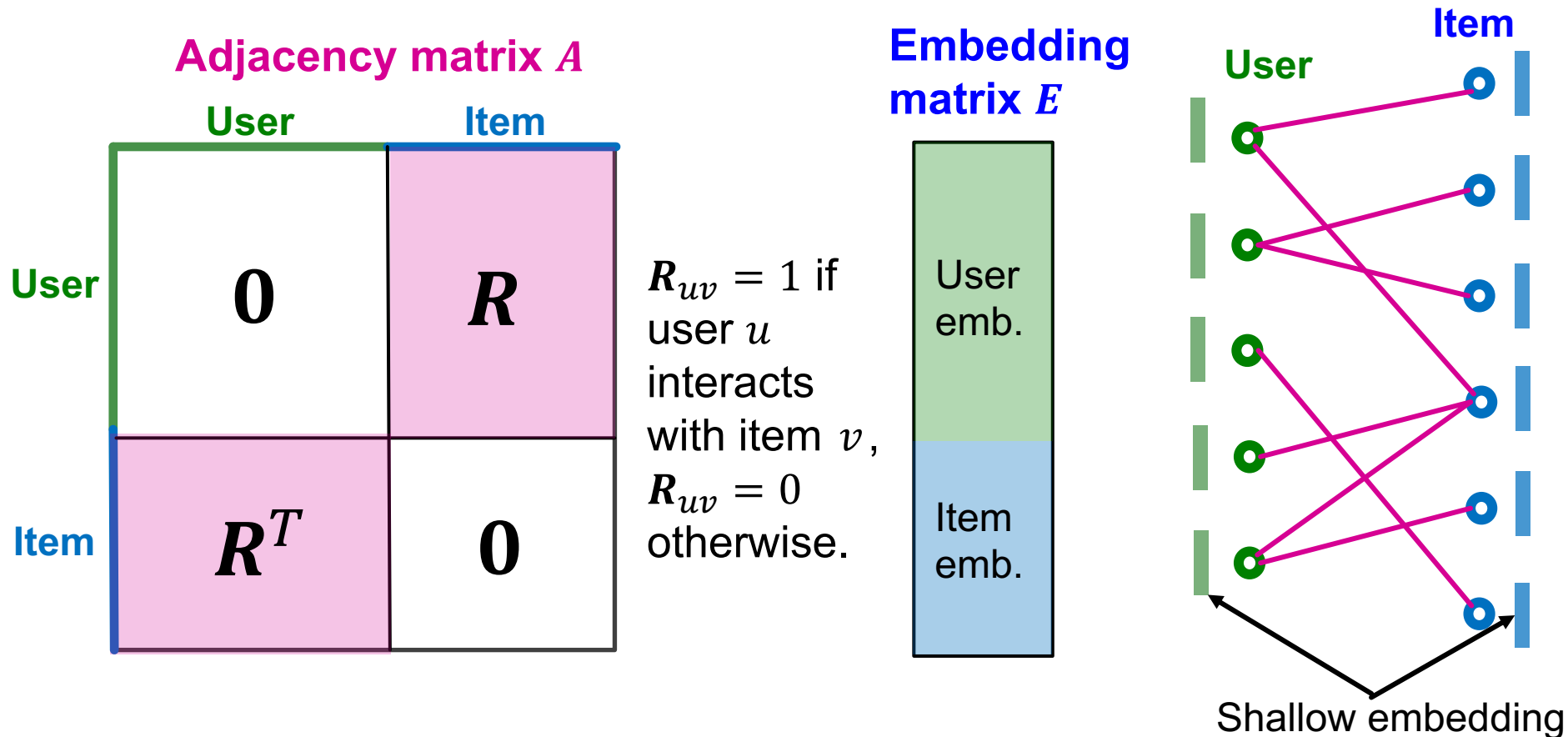
# Stanford CS224W: LightGCN

CS224W: Machine Learning with Graphs
Jure Leskovec, Stanford University
http://cs224w.stanford.edu

# Adjacency and Embedding Matrices

- **Adjacency matrix** of a (undirected) bipartite graph.
- **Shallow embedding matrix**.



Adjacency matrix $A$

| | User | Item |
|---|---|---|
| **User** | $\mathbf{0}$ | $R$ |
| **Item** | $R^T$ | $\mathbf{0}$ |

$R_{uv} = 1$ if user $u$ interacts with item $v$, $R_{uv} = 0$ otherwise.

Embedding matrix $E$

User emb.

Item emb.

User    Item

Shallow embedding

# Matrix Formulation of GCN

- **Recall**: Diffusion matrix of C&S.
- Let $D$ be the degree matrix of $A$.
- Define the normalized adjacency matrix $\widetilde{A}$ as
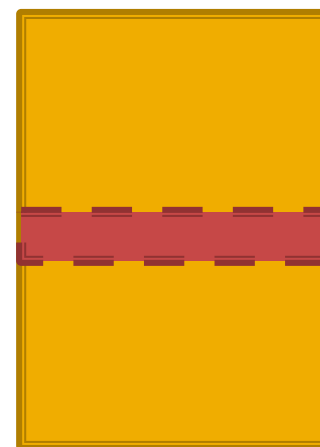
$$\widetilde{A} \equiv D^{-1/2}AD^{-1/2}$$

- Let $E^{(k)}$ be the embedding matrix at $k$-th layer.
- Each layer of GCN's aggregation can be written in a matrix form:

$$E^{(k+1)} = \text{ReLU}\left(\widetilde{A}E^{(k)}\,W^{(k)}\right)$$

**Neighbor aggregation**   **Learnable linear transformation**

Matrix of node embeddings $E^{(k)}$

**Note**: Different from the original GCN, self-connection is omitted here.



**Each row stores node embedding**

# Simplifying GCN (1)

- **Simplify GCN by removing ReLU non-linearity:**

$$E^{(k+1)} = \widetilde{A} E^{(k)} W^{(k)}$$

<span style="color:darkred">**Original idea from SGC [Wu et al. 2019]**</span>

- The final node embedding matrix is given as

$$E^{(K)} = \widetilde{A} \, E^{(K-1)} W^{(K-1)}$$

$$= \widetilde{A} \big( \widetilde{A} E^{(K-2)} W^{(K-2)} \big) W^{(K-1)}$$

$$= \widetilde{A} \big( \widetilde{A} ( \cdots ( \widetilde{A} E^{(0)} W^{(0)} ) \cdots ) W^{(K-2)} \big) W^{(K-1)}$$

Set $E$ as input embedding $E^{(0)}$

$$= \widetilde{A}^{K} \, E \, \big( W^{(0)} \cdots W^{(K-1)} \big)$$

# Simplifying GCN (2)

- Removing ReLU significantly simplifies GCN!

$$E^{(K)} = \boxed{\widetilde{A}^K \, E} \, W$$

**Diffusing node embeddings along the graph**
(similar to C&S that diffuses soft labels along the graph)

$$W \equiv W^{(0)} \cdots W^{(K-1)}$$

- **Algorithm**: Apply $E \leftarrow \widetilde{A} \, E$ for $K$ times.

  - Each matrix multiplication diffuses the current embeddings to their one-hop neighbors.

  - **Note:** $\widetilde{A}^K$ is dense and never gets materialized. Instead, the above iterative matrix-vector product is used to compute $\widetilde{A}^K \, E$.
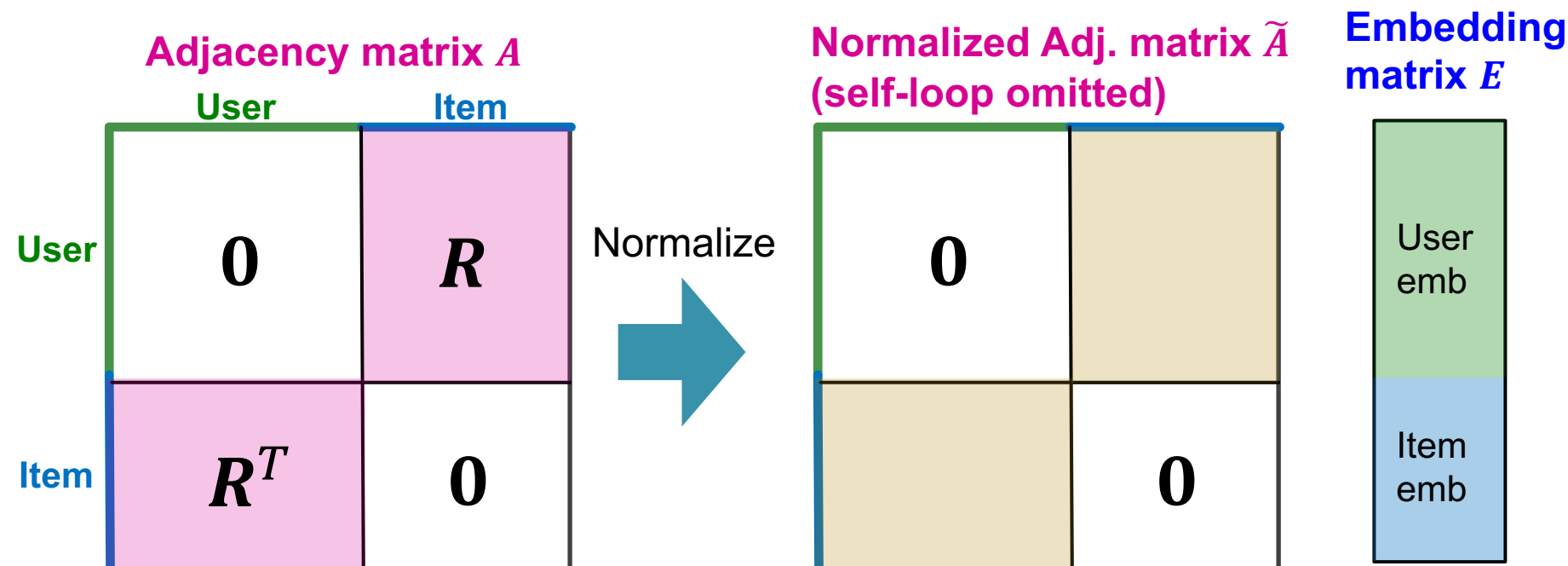
# Multi-Scale Diffusion

- We can consider **multi-scale diffusion**

$$\alpha_0 E^{(0)} + \alpha_1 E^{(1)} + \alpha_2 E^{(2)} + \cdots + \alpha_K E^{(K)}$$

  - The above includes embeddings diffused at multiple hop scales.

  - $\alpha_0 E^{(0)} = \alpha_0 \widetilde{A}^0 E^{(0)}$ acts as a self-connection (that is omitted in the definition $\widetilde{A}$)

  - The coefficients, $\alpha_0, \dots, \alpha_K$, are hyper-parameters.

- For simplicity, LightGCN uses the uniform coefficient, i.e., $\alpha_k = \frac{1}{K+1}$ for $k = 0, \dots, K$.

# LightGCN: Model Overview (1)

- **Given**:
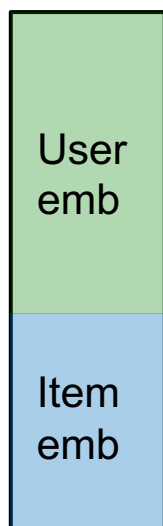  - **Adjacency matrix A**
  - **Initial learnable embedding matrix $E$**

**Adjacency matrix $A$**

|  | User | Item |
|---|---|---|
| **User** | $0$ | $R$ |
| **Item** | $R^T$ | $0$ |

Normalize →

**Normalized Adj. matrix $\widetilde{A}$ (self-loop omitted)**

|  |  |
|---|---|
| $0$ |  |
|  | $0$ |

**Embedding matrix $E$**

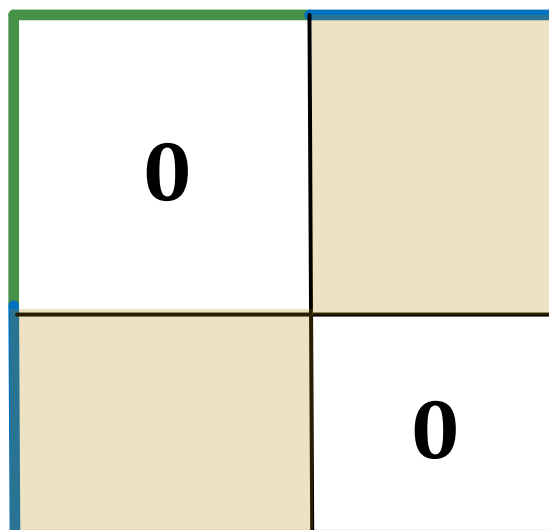| User emb |
|---|
| Item emb |

- Iteratively diffuse embedding matrix $E$ using $\widetilde{A}$
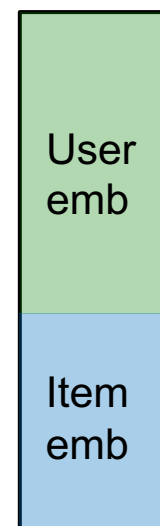
For $k = 0 \ldots K - 1$,



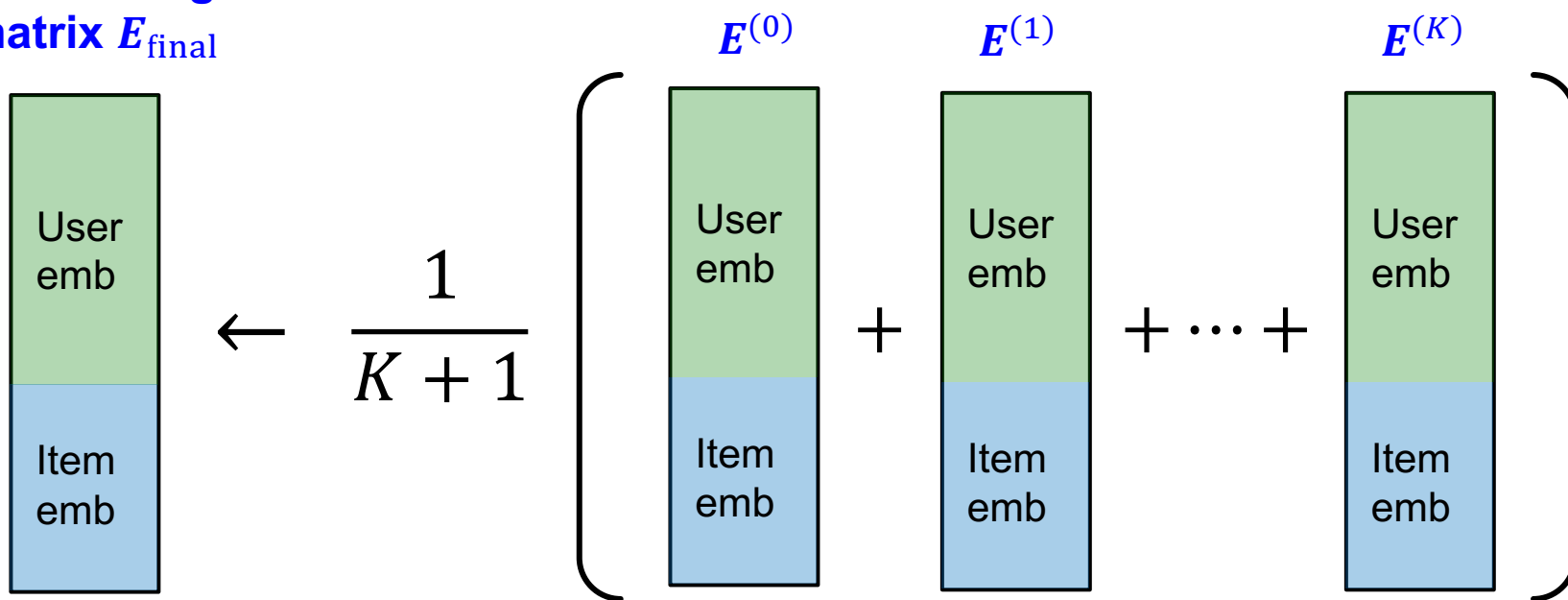**Embedding matrix $E^{(k+1)}$**

**Normalized Adj. matrix $\widetilde{A}$ (self-loop omitted)**

**Embedding matrix $E^{(k)}$ ($E^{(0)}$ is set to $E$)**

User emb

Item emb

0

0

$\cdot$

User emb

Item emb

# LightGCN: Model Overview (3)

- Average the embedding matrices at different scales.

- **Score function**:
  - Use user/item vectors from $E_{\text{final}}$ to score user-item interaction

# LightGCN: Intuition

- **Question**: **Why does the simple diffusion propagation work well?**
- **Answer**: The diffusion directly encourages the embeddings of similar users/items to be similar.

  - Similar users share many common neighbors (items) and are expected to have similar future preferences (interact with similar items).

# LightGCN and GCN/C&S

- The embedding propagation of LightGCN is closely related to GCN/C&S.
- **Recall**: GCN/C&S (neighbor aggregation part)

$$h_v^{(k+1)} = \sum_{u \in N(v)} \frac{1}{\sqrt{d_u}\sqrt{d_v}} \cdot h_u^{(k)}$$

**Node degree**

  - Self-loop is added in the neighborhood definition.
- LightGCN uses the same equation except that
  - Self-loop is *not* added in the neighborhood definition.
  - Final embedding takes the average of embeddings from all the layers: $h_v = \frac{1}{K+1}\sum_{k=0}^{K} h_v^{(k)}$.

# LightGCN and MF: Comparison

- Both LightGCN and Matrix Factorization (MF) **learn a unique embedding for each user/item.**
- The difference is that
  - MF directly uses the shallow user/item embeddings for scoring.
  - LightGCN uses the *diffused* user/item embeddings for scoring.
- LightGCN performs better than MF but are also more computationally expensive due to the additional diffusion step.
  - The final embedding of a user/item is obtained by aggregating embeddings of its multi-hop neighboring nodes.

# LightGCN: Summary

- LightGCN simplifies NGCF by **removing the learnable parameters of GNNs.**

- **Learnable parameters are all in the shallow input node embeddings.**

  - Diffusion propagation only involves matrix-vector multiplication.

  - The simplification leads to better empirical performance than NGCF.