

Contents

I	Graph Convolutional Neural Network (GCN)	3
1	Problem	3
2	Method	3
3	Formula	3
3.1	Matrix form	3
4	Experimental result and Takeaways	4
4.1	Some side results	4
II	Simplifying Graph Convolutional Network (SGCN)	4
1	Problem	4
2	Method	4
3	Formula	4
4	Experimental result	5
III	Graph Attention Network (GAT)	5
1	Problem	5
2	Method	5
3	Formula	6
4	Experimental results and Takeaways	7
IV	LightGCN: Simplifying and Powering GCN for Recommendation	7
1	Problem	7
2	Method	7

3	Formula	8
3.1	Matrix form	8
4	Experimental Results and Takeaways	8
V	Graph Neural Networks for Social Recommendation	8
1	Problem	9
2	Method	9
3	Formula	10
3.1	User Modeling	10
3.1.1	Item Aggregation	10
3.1.2	Social Aggregation	10
3.1.3	Learning User Latent Factor	11
3.2	Item Modeling	11
3.2.1	User Aggregation	11
3.3	Rating prediction	11
3.4	Loss function	11
4	Experiment Results and Takeaways	11

Part I

Graph Convolutional Neural Network (GCN)

Paper

1 Problem

Given dataset $X \in \mathbb{R}^{n \times d}$ where n is the number of nodes, d is the number of features. We want to classify each node into one of C categories.

2 Method

Similar to CNNs, GCNs learn new feature representations of X over multiple layers. In k th-layer, the node representation is a matrix $H^{(k)} \in \mathbb{R}^{n \times d}$. The core idea of GCN is the node representation of a node i at k th-layer depends on the node representations of all its neighboring nodes at $(k-1)$ -th layer. This is called **feature propagation**. After stacking L layers, the final node representations of node i may contain information propagated from all neighboring nodes which are at most k -hop away from i .

3 Formula

3.1 Matrix form

- In one line:

$$X = H^{(0)} \rightarrow H^{(1)} \rightarrow \dots \rightarrow H^L \rightarrow \hat{Y} \quad (1)$$

- Matrix form:

$$H^{(l+1)} = \sigma(\tilde{D}^{-0.5} \tilde{A} \tilde{D}^{-0.5} H^{(l)} W^{(l)}) \quad (2)$$

$$\hat{Y} = \text{softmax}(H^{(L)}) \quad (3)$$

where $H^{(l)} \in \mathbb{R}^{n \times d}$, $\tilde{A} \in \mathbb{R}^{n \times n}$, $\tilde{D} \in \mathbb{R}^{n \times n}$, $W^{(l)} \in \mathbb{R}^{d \times d}$.

- Step-by-step:

Let $S = \tilde{D}^{-0.5} \tilde{A} \tilde{D}^{-0.5}$.

- Feature propagation:

$$\bar{H}^{(k)} = S H^{(k-1)} \quad (4)$$

- Feature transformation

$$H^{(k)} = \text{ReLU}(\bar{H}^{(k)} \Theta^{(k)}) \quad (5)$$

– Classifier

$$\hat{Y} = \text{softmax}(H^{(L)})$$

4 Experimental result and Takeaways

Run semi-supervised node classification on Citeseer, Cora, Pubmed, NELL. Achieved state-of-the-art performance at that time.

4.1 Some side results

- Renormalization trick made huge improvement over existing methods (like Chebyshev filter)
- Best results are obtained with 2-layer or 3-layer model. Beyond that, performance drops. By adding residual connection, deeper model can achieve nearly the same performance with shallow ones.

Part II

Simplifying Graph Convolutional Network (SGCN)

Paper

1 Problem

Given dataset $X \in \mathbb{R}^{n \times d}$ where n is the number of nodes, d is the number of features. We want to classify each node into one of C categories.

2 Method

Similar to GCN, but drop feature transformation step.

3 Formula

In one line:

$$\hat{Y} = \text{softmax}(S^K X \Theta) \tag{6}$$

4 Experimental result

Achieve comparable performance with GCN. Can be run on larger dataset (Reddit). Memory requirement and running time are better than GCN by constant factors.

Part III

Graph Attention Network (GAT)

Paper

1 Problem

Given dataset $X \in \mathbb{R}^{n \times d}$ where n is the number of nodes, d is the number of features. We want to classify each node into one of C categories.

2 Method

Learn new feature representations of X over multiple layers. Input to a layer is a set of node features $h = \{h_1, h_2, \dots, h_N\}$ where $h_i \in \mathbb{R}^F$ and F is the number of features in each node. Output of that layer is a set of node features $h' = \{h'_1, h'_2, \dots, h'_N\}$ where $h'_i \in \mathbb{R}^{F'}$ and F' is the number of features in each node. Information is propagated among all neighboring nodes via **attention mechanism**. Key point of the attention mechanism is that it allows for assigning different importances to nodes in the same neighborhood.

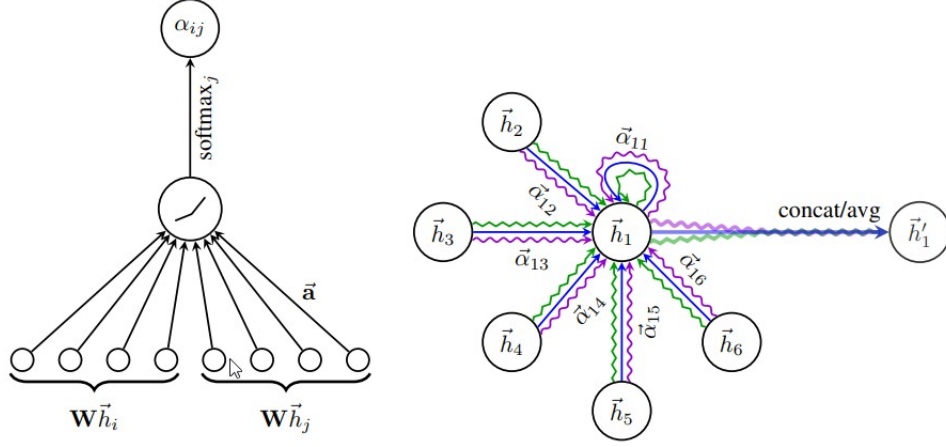


Figure 1: **Left:** The attention mechanism $a(\vec{W}h_i, \vec{W}h_j)$ employed by our model, parametrized by a weight vector $\vec{a} \in \mathbb{R}^{2F'}$, applying a LeakyReLU activation. **Right:** An illustration of multi-head attention (with $K = 3$ heads) by node 1 on its neighborhood. Different arrow styles and colors denote independent attention computations. The aggregated features from each head are concatenated or averaged to obtain \vec{h}'_1 .

3 Formula

- Attention of node i to node j:

$$e_{ij} = a(W h_i, W h_j) \quad (7)$$

- Attention, normalized:

$$\alpha_{ij} = \text{softmax}_j(e_{ij}) = \frac{\exp(e_{ij})}{\sum_{k \in N_i} \exp(e_{ik})} \quad (8)$$

- Attention, normalized, fully expanded equation:

$$\alpha_{ij} = \frac{\exp(\text{LeakyReLU}(a^T [W h_i || W h_j]))}{\sum_{k \in N_i} \exp(\text{LeakyReLU}(a^T [W h_i || W h_k]))} \quad (9)$$

where $||$ is concatenation and $a \in \mathbb{R}^{2F}$

- Feature propagation:

$$h'_i = \sigma \left(\sum_{j \in N_i} \alpha_{ij} W h_j \right) \quad (10)$$

- Feature propagation, with multi-head attention:

$$h'_i = ||_{k=1}^K \sigma \left(\sum_{j \in N_i} \alpha_{ij}^k W^k h_j \right) \quad (11)$$

- Last layer:

$$h'_i = \sigma \left(\frac{1}{K} \sum_{k=1}^K \sum_{j \in N_i} \alpha_{ij}^k W^k h_j \right) \quad (12)$$

4 Experimental results and Takeaways

- Run transductive learning semi-supervised node classification on Citeseer, Cora, Pubmed. Achieve significantly better result than GCN.
- Run inductive learning node classification on PPI. Achieve state-of-the-art and significantly better result than GraphSAGE-GCN.

Part IV

LightGCN: Simplifying and Powering GCN for Recommendation

Paper

1 Problem

Given N users, M items and a list $\{(u, i)\}$ of interactions between user u and item i , predict whether a user will interact with an item.

2 Method

- Learn **latent features (embedding)** to represent users and items, then perform prediction based on the embedding vectors
- Problem is represented as bipartite graph, one side are users and the other side are items.
- Adopt the feature propagation idea of GCN, however abandon the use of feature transformation and nonlinear activation. The k th-layer representation of a user u contain information propagated from all the nodes which are k -hop away from u . Final representation is a weighted sum of all the layer-presentations.
- Model is trained with Bayesian Personalized Ranking (BPR) loss.

3 Formula

3.1 Matrix form

- Adjacency matrix:

$$A = \begin{pmatrix} 0 & R \\ R^T & 0 \end{pmatrix} \quad (13)$$

where $R_{ui} = 1$ [user u interacted with item i].

- Feature propagation:

$$E^{(k+1)} = (D^{-0.5} A D^{-0.5}) E^{(k)} \quad (14)$$

- Final representation:

$$E = \alpha_0 E^{(0)} + \dots + \alpha_K \tilde{A}^K E^{(K)} \quad (15)$$

where $\tilde{A} = D^{-0.5} A D^{-0.5}$

- Prediction:

$$\hat{y}_{ui} = e_u^T e_i \quad (16)$$

- Loss function:

$$L_{BPR} = - \sum_{u=1}^M \sum_{i \in N_u} \sum_{j \notin N_u} \ln \sigma(\hat{y}_{ui} - \hat{y}_{uj}) + \lambda \|E^{(0)}\|^2 \quad (17)$$

4 Experimental Results and Takeaways

- Run on Gowalla, Yelp2018, Amazon-Book which are 3 datasets with the number of users 30,000 - 50,000; number of items 40,000 - 90,000; number of edges 1,00,000 - 3,000,000. Outperforms NGCF (which was state-of-the-art at collaborative filtering) by a large margin.
- The author argued that much of the effectiveness of LGCN was due to its successful embedding smoothness (i.e. two users overlapping in a large number of items should have similar embedding vectors).

Part V

Graph Neural Networks for Social Recommendation

Paper

1 Problem

Given $U = \{u_1, u_2, \dots, u_n\}$ and $V = \{v_1, v_2, \dots, v_m\}$ the set of users and items respectively; $R \in \mathbb{R}^{m \times n}$ the user-item rating matrix which have two parts: observed part $\mathbb{O} = \{(u_i, v_j) | r_{ij} \neq 0\}$ and missing part $\mathbb{T} = \{(u_i, v_j) | r_{ij} = 0\}$; and user-user social graph $T \in \mathbb{R}^{n \times n}$. We are to predict the missing part of R .

2 Method

- Learn user and item embeddings, use these embeddings to perform prediction
- User latent factors are obtained by combining information from both item space and item space
- Item latent factors are obtained via user aggregation
- Loss function is simply defined as discrepancy between predicted and ground-truth ratings for ratings in observed set

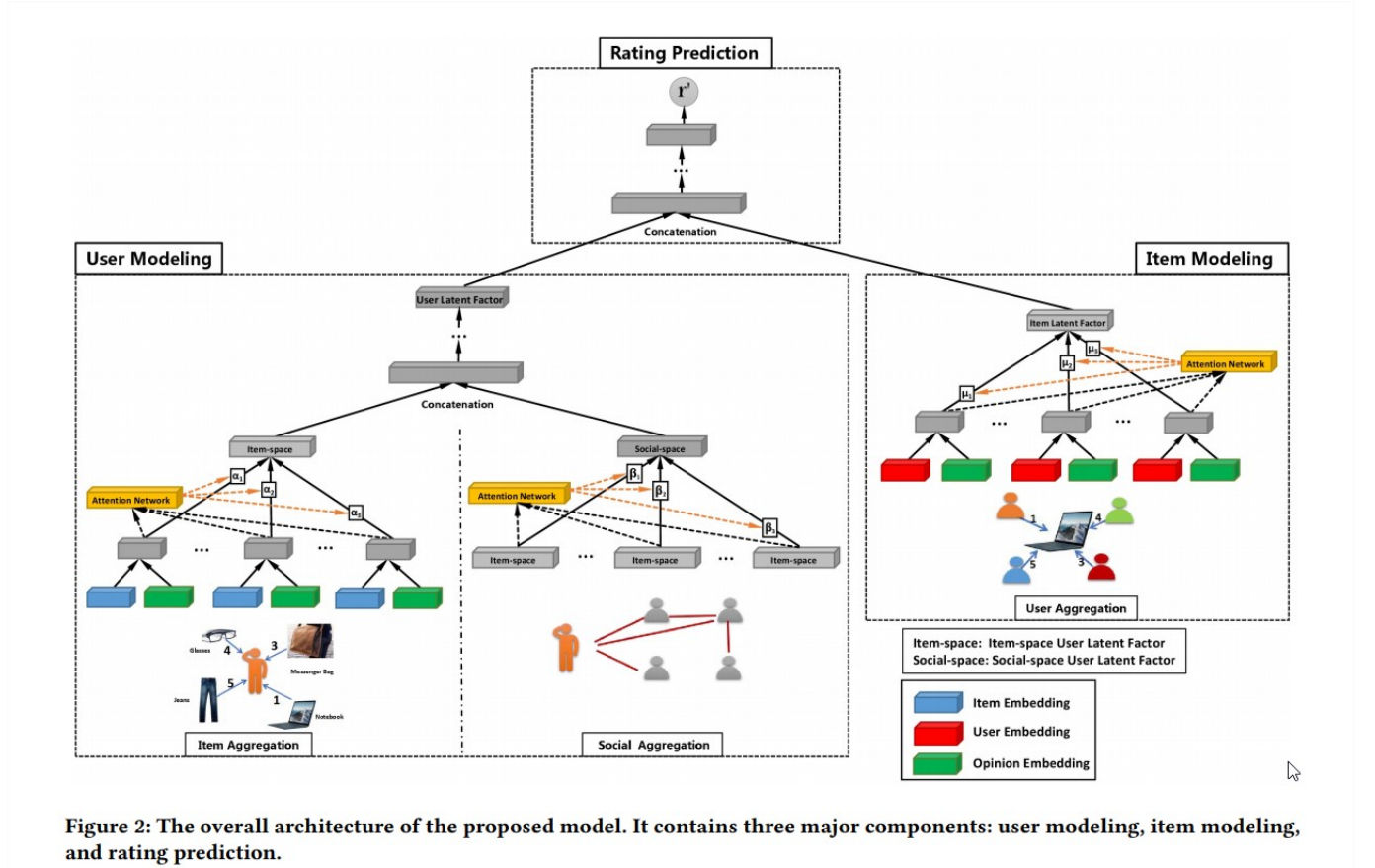


Figure 2: The overall architecture of the proposed model. It contains three major components: user modeling, item modeling, and rating prediction.

3 Formula

3.1 User Modeling

User modeling aims to learn user latent factors, denoted as $h_i \in R^d$ for user u_i . To do that, we try to learn item-space user latent factor $h_i^I \in R^d$ from the user-item graph and user-space user latent factor $h_i^S \in R^d$, then combine them together to obtain h_i .

3.1.1 Item Aggregation

- Item-space user latent factor:

$$h_i^I = \sigma \left(W \left\{ \sum_{a \in C(i)} \alpha_{ia} x_{ia} \right\} + b \right) \quad (18)$$

where x_{ia} is opinion-aware interaction between user i and item a, $C(i)$ is set of all items rated by user i, α_{ia} the attention weight of the iteration between user i and item a.

- Opinion-aware interaction:

$$x_{ia} = g_v([q_a || e_r]) \quad (19)$$

where g_v is a MLP, q_a is learnable item embedding of item a, e_r is the embedded rating of rating r where $r \in \{1, 2, 3, 4, 5\}$.

- Attention weight α_{ia} :

$$\alpha_{ia}^* = w_2^T \sigma \left(W_1 [x_{ia} || p_i] + b_1 \right) + b_2 \quad (20)$$

$$\alpha_{ia} = \text{softmax}(\alpha_{ia}^*) \quad (21)$$

where p_i is learnable user embedding of user i.

3.1.2 Social Aggregation

- Social-space user latent factor:

$$h_i^S = \sigma \left(W \left\{ \sum_{o \in N(i)} \beta_{io} h_o^I \right\} + b \right) \quad (22)$$

where β_{io} is the strength of connection between user i and user o

- Social attention β_{io}

$$\beta_{io}^* = w_2^T \sigma \left(W_1 [h_o^I || p_i] + b_1 \right) + b_2 \quad (23)$$

$$\beta_{io} = \text{softmax}(\beta_{io}^*) \quad (24)$$

where p_i is learnable user embedding of user i.

3.1.3 Learning User Latent Factor

$$h_i = FC^l([h_i^I || h_i^S]) \quad (25)$$

where FC is one fully connected layer of a neural network.

3.2 Item Modeling

3.2.1 User Aggregation

- Item latent factor:

$$z_j = \sigma\left(W\left\{\sum_{t \in B(j)} \mu_{jt} f_{jt}\right\} + b\right) \quad (26)$$

where $B(j)$ is the set of users who interacted with item j , f_{jt} the opinion-aware interaction representation of user t for item j , and μ_{jt} the user attention of user t in contributing to item j .

- User attention to item:

$$\mu_{jt}^* = w_2^T \sigma\left(W_1[f_{jt} || q_j] + b_1\right) + b_2 \quad (27)$$

$$\mu_{jt} = \text{softmax}(\mu_{jt}^*) \quad (28)$$

- Opinion-aware interaction:

$$f_{jt} = g_u([p_t || e_r]) \quad (29)$$

where g_u is a MLP, p_t is a learnable user embedding.

3.3 Rating prediction

$$r'_{ij} = w^T FC^{l-1}([h_i || z_j]) \quad (30)$$

where FC is one fully connected layer of a neural network.

3.4 Loss function

$$Loss = \frac{1}{2|\mathbb{O}|} \sum_{i,j \in \mathbb{O}} (r'_{ij} - r_{ij})^2 \quad (31)$$

4 Experiment Results and Takeaways

- Use two rating datasets: Ciao and Epinions
- Train-Validation-Test split: $(x, (1-x)/2, (1-x)/2)$, $x \in [60\%, 80\%]$. Achieve state-of-the-art for both datasets.
- Attention mechanism is an important factor in the success of this model.