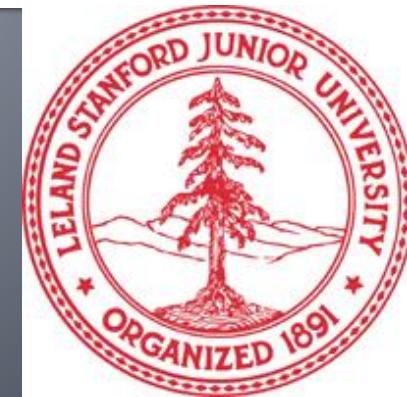


Stanford CS224W: Graph Neural Networks

CS224W: Machine Learning with Graphs

Jure Leskovec, Stanford University

<http://cs224w.stanford.edu>



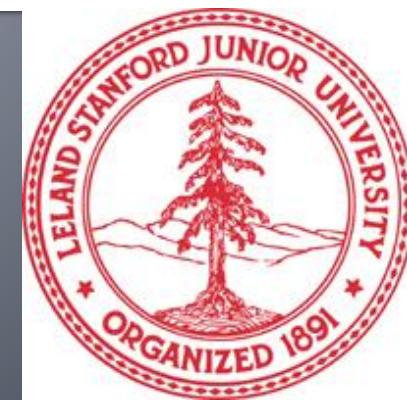
ANNOUNCEMENTS

- **Today (01/19):** HW 1 out
- **Monday (01/23):** Recitation session for HW 1
- **Next Thursday (01/26):** Colab 1 due, Colab 2 out

CS224W: Machine Learning with Graphs

Jure Leskovec, Stanford University

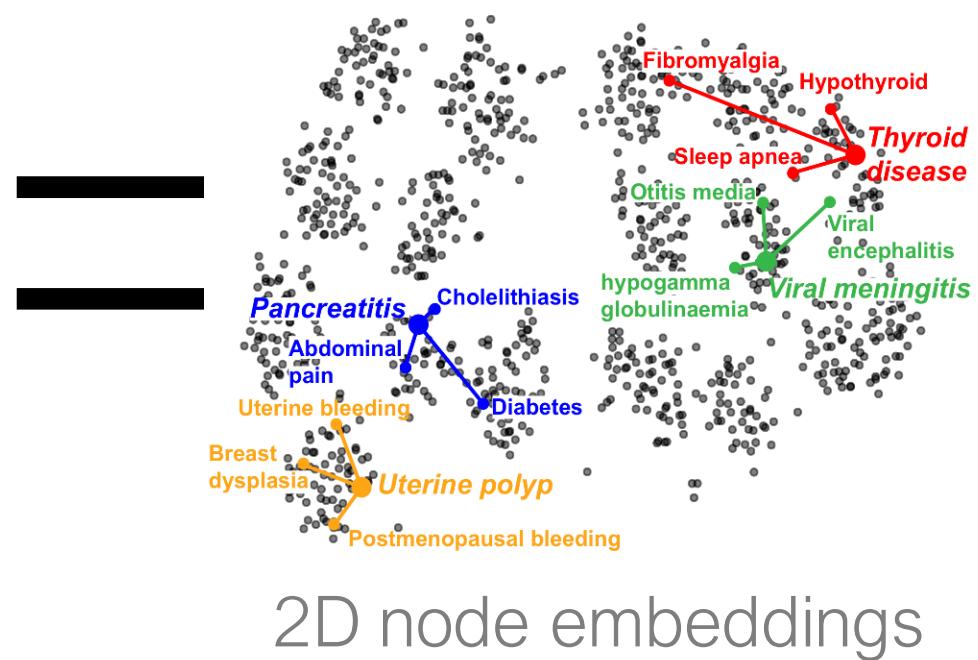
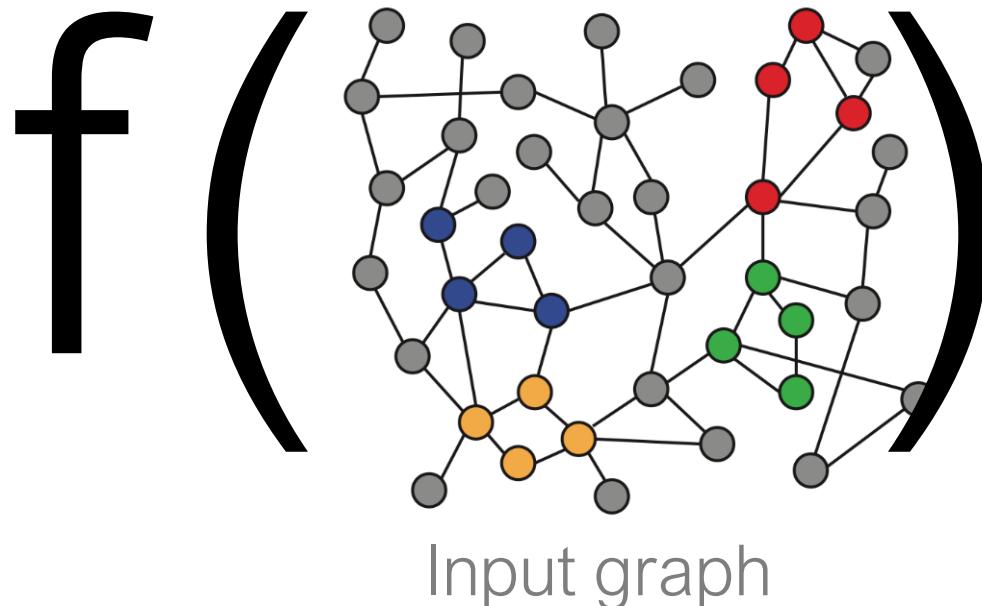
<http://cs224w.stanford.edu>



Recap: Node Embeddings

图嵌入表示学习：

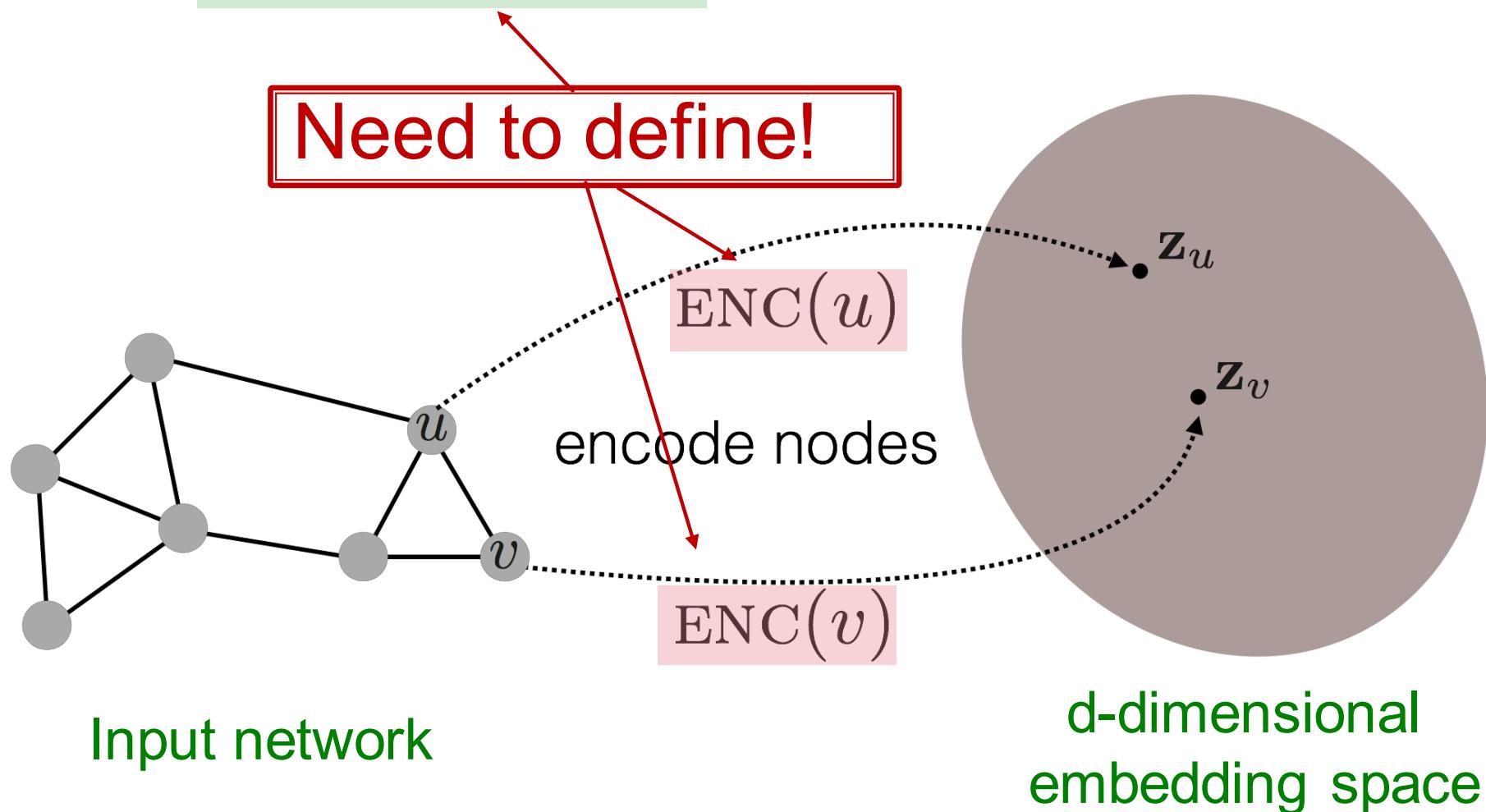
- **Intuition:** Map nodes to d -dimensional embeddings such that similar nodes in the graph are embedded close together



How to learn mapping function f ?

Recap: Node Embeddings

Goal: $\text{similarity}(u, v) \approx \mathbf{z}_v^T \mathbf{z}_u$



Recap: Two Key Components

- **Encoder:** Maps each node to a **low-dimensional vector**

$$\text{ENC}(v) = \mathbf{z}_v$$

node in the input graph

\mathbf{z}_v is a d -dimensional embedding

- **Similarity function:** Specifies how the relationships in vector space map to the relationships in the original network

$$\text{similarity}(u, v) \approx \mathbf{z}_v^T \mathbf{z}_u$$

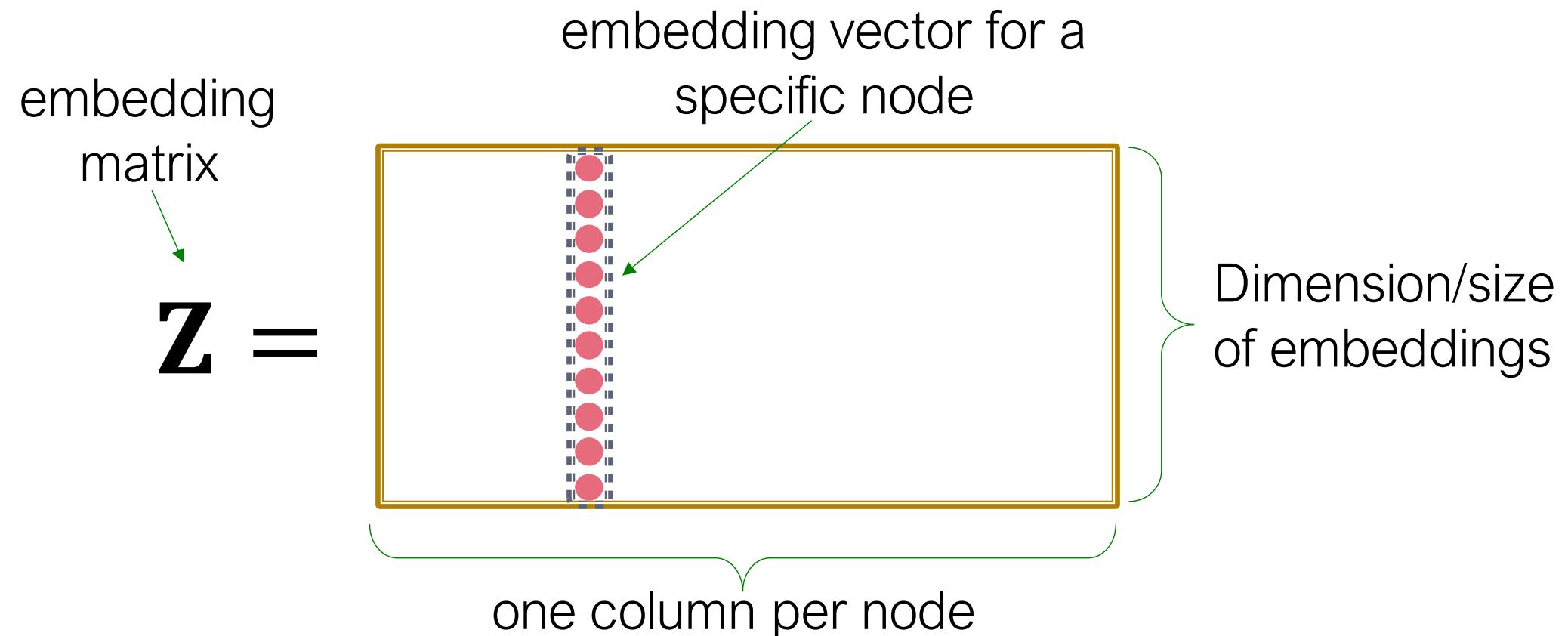
Similarity of u and v in the original network

$\mathbf{z}_v^T \mathbf{z}_u$ is a dot product between node embeddings

Decoder

Recap: “Shallow” Encoding

Simplest encoding approach: **Encoder is just an embedding-lookup**



Recap: Shallow Encoders

Deepwalk, Node2vec, LINE

■ Limitations of shallow embedding methods:

- **$O(|V|d)$ parameters are needed:** 每个节点的嵌入向量都要
参数共享 独立训练.
 - No sharing of parameters between nodes
 - Every node has its own unique embedding
- **Inherently “transductive”:** 直推式: 无法泛化到新图/新节点
 - Cannot generate embeddings for nodes that are not seen during training
- **Do not incorporate node features:** 没有用到节点属性特征标注
 - Nodes in many graphs have features that we can and should leverage
利用

Today: Deep Graph Encoders

图深度学习, 图神经网络

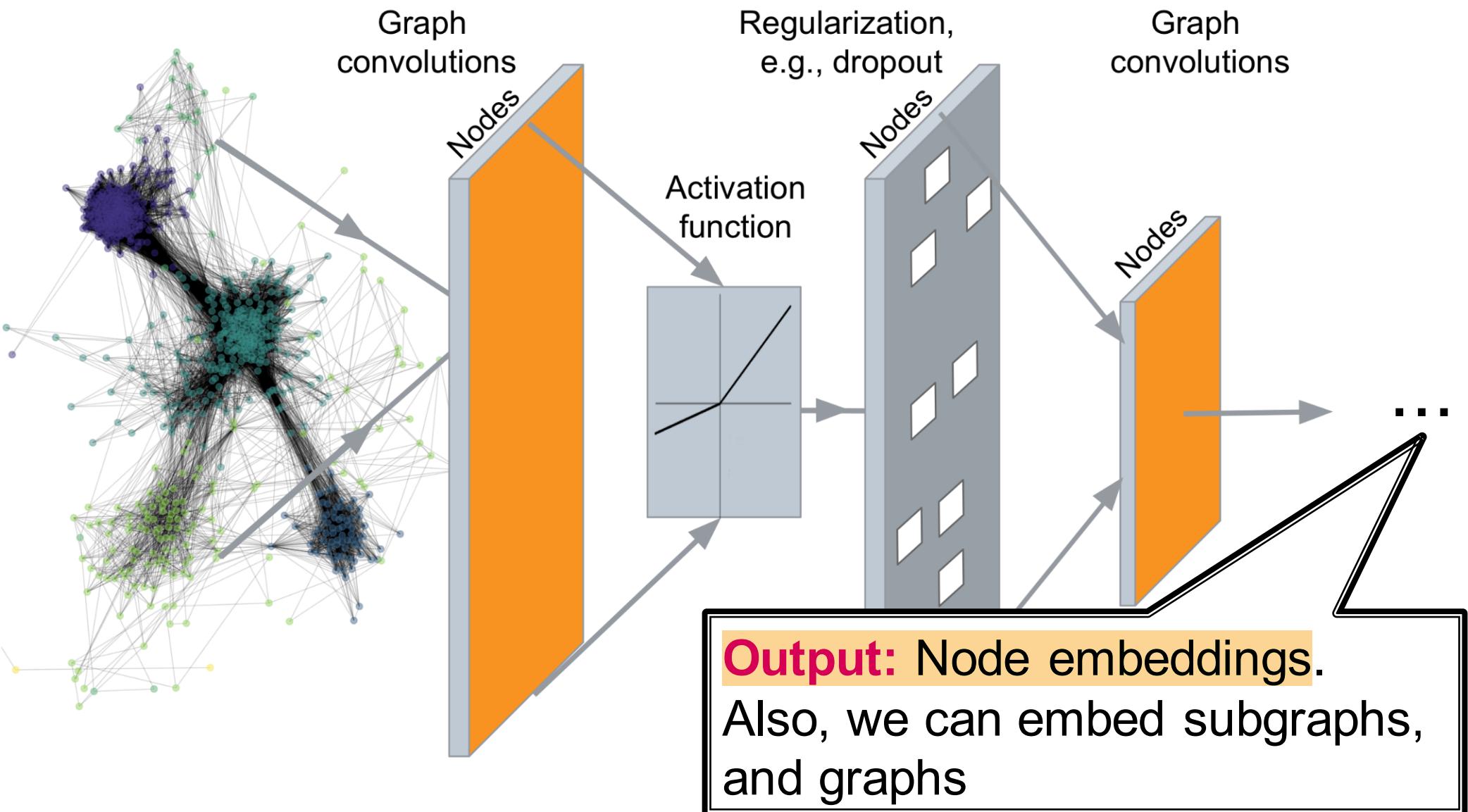
- **Today:** We will now discuss deep learning methods based on graph neural networks (GNNs):

$$\text{ENC}(v) =$$

multiple layers of
non-linear transformations
based on graph structure

- **Note:** All these deep encoders can be **combined with node similarity functions** defined in the Lecture 3.

Deep Graph Encoders

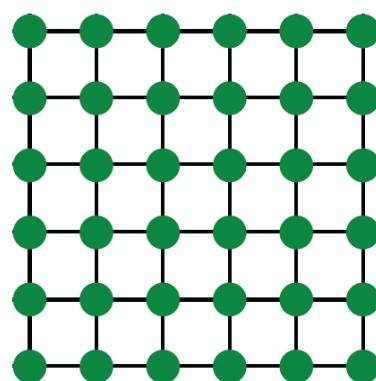


Tasks on Networks

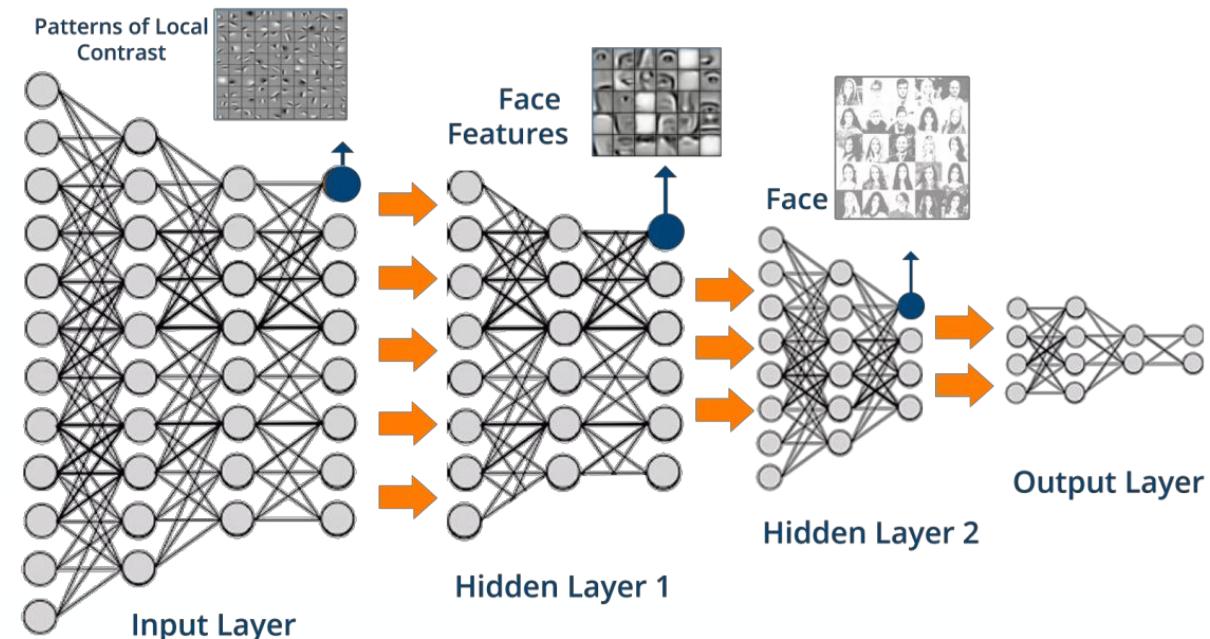
Tasks we will be able to solve:

- Node classification
 - Predict the type of a given node
- Link prediction
 - Predict whether two nodes are linked
- Community detection
 - Identify densely linked clusters of nodes
- Network similarity
 - How similar are two (sub)networks

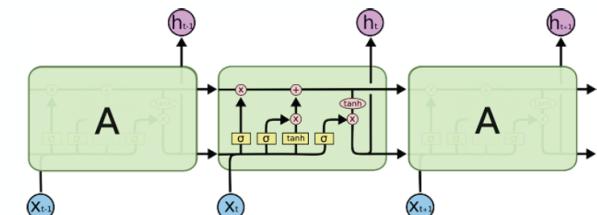
Modern ML Toolbox



Images



Text/Speech

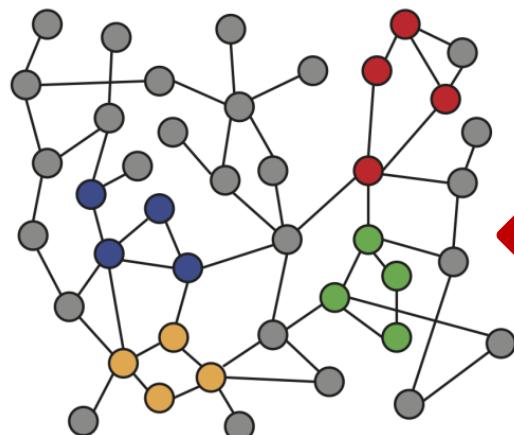


Modern deep learning toolbox is designed for simple sequences & grids

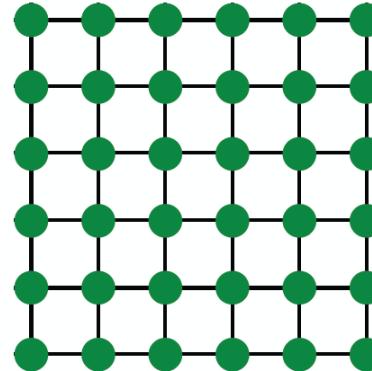
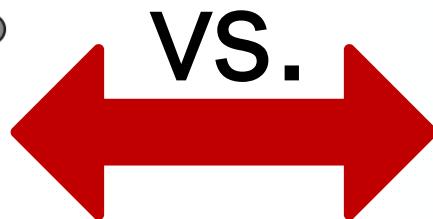
Why is it Hard?

But networks are far more complex!

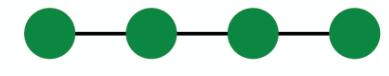
- Arbitrary size and complex topological structure (i.e., no spatial locality like grids)



Networks



Images



Text

- No fixed node ordering or reference point
- Often dynamic and have multimodal features

Outline of Today's Lecture

1. Basics of deep learning



2. Deep learning for graphs

3. Graph Convolutional Networks

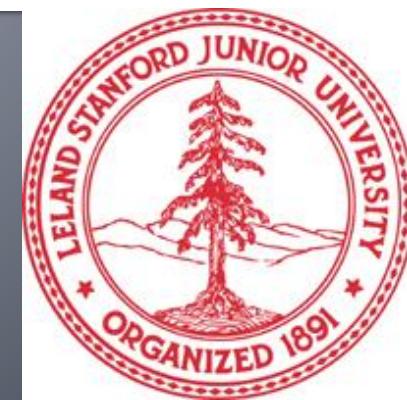
4. GNNs subsume CNNs

Stanford CS224W: Basics of Deep Learning

CS224W: Machine Learning with Graphs

Jure Leskovec, Stanford University

<http://cs224w.stanford.edu>



Machine Learning as Optimization

- **Supervised learning:** we are given input x , and the goal is to predict label y .
- **Input x can be:**
 - Vectors of real numbers
 - Sequences (natural language)
 - Matrices (images)
 - Graphs (potentially with node and edge features)
- **We formulate the task as an optimization problem.**

Machine Learning as Optimization

- **Formulate the task as an optimization problem:**

$$\min_{\Theta} \mathcal{L}(y, f(x))$$

Objective function

- Θ : a set of **parameters** we optimize
 - Could contain one or more scalars, vectors, matrices ...
 - E.g. $\Theta = \{Z\}$ in the shallow encoder (the embedding lookup)

- **\mathcal{L} : loss function.** Example: L2 loss

$$\mathcal{L}(y, f(x)) = \|y - f(x)\|_2$$

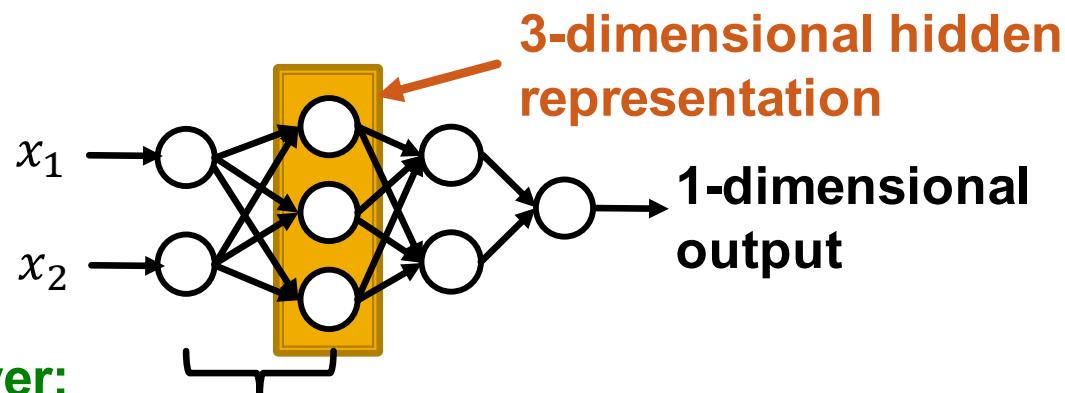
- Other common loss functions:
 - L1 loss, huber loss, max margin (hinge loss), cross entropy ...
 - See <https://pytorch.org/docs/stable/nn.html#loss-functions>

Multi-layer Perceptron (MLP)

- Each layer of MLP combines linear transformation and non-linearity:

$$\mathbf{x}^{(l+1)} = \sigma(W_l \mathbf{x}^{(l)} + b^l)$$

- where W_l is weight matrix that transforms hidden representation at layer l to layer $l + 1$
- b^l is bias at layer l , and is added to the linear transformation of $\mathbf{x}^{(l)}$
- σ is non-linearity function (e.g., sigmoid)
- Suppose \mathbf{x} is 2-dimensional, with entries x_1 and x_2



Summary

- **Objective function:**

$$\min_{\Theta} \mathcal{L}(y, f(\mathbf{x}))$$

- f can be a simple linear layer, an MLP, or other neural networks (e.g., a GNN later)
- Sample a minibatch of input \mathbf{x}
- **Forward propagation:** Compute \mathcal{L} given \mathbf{x}
- **Back-propagation:** Obtain gradient $\nabla_{\mathbf{W}} \mathcal{L}$ using a chain rule.
- Use **stochastic gradient descent (SGD)** to optimize for Θ over many iterations.

Outline of Today's Lecture

1. Basics of deep learning



2. Deep learning for graphs



3. Graph Convolutional Networks

4. GNNs subsume CNNs

Stanford CS224W: Deep Learning for Graphs

CS224W: Machine Learning with Graphs

Jure Leskovec, Stanford University

<http://cs224w.stanford.edu>



Content

- **Local network neighborhoods:**
 - Describe aggregation strategies
 - Define computation graphs
- **Stacking multiple layers:**
 - Describe the model, parameters, training
 - How to fit the model?
 - Simple example for unsupervised and supervised training

Setup

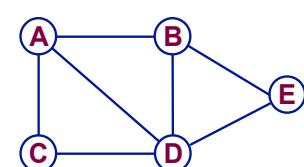
■ Assume we have a graph G :

- V is the **vertex set**
- A is the **adjacency matrix** (assume binary)
- $X \in \mathbb{R}^{\frac{|V| \times d}{\text{节点个数}}}$ is a matrix of **node features**
- v : a node in V ; $N(v)$: the set of neighbors of v .
- **Node features:** 节点的邻域: 与节点 v 相连的所有节点集合
 - Social networks: User profile, User image
 - Biological networks: Gene expression profiles, gene functional information
 - When there is no node feature in the graph dataset:
 - Indicator vectors (one-hot encoding of a node)
 - Vector of constant 1: $[1, 1, \dots, 1]$

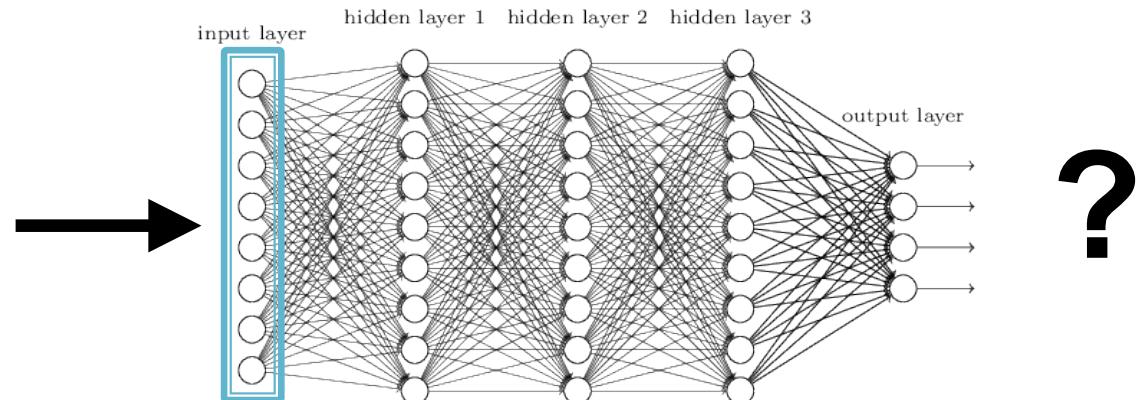
A Naïve Approach

朴素想法：直接输入邻接矩阵A

- Join adjacency matrix and features
- Feed them into a deep neural net:



	A	B	C	D	E	Feat
A	0	1	1	1	0	1 0
B	1	0	0	1	1	0 0
C	1	0	0	1	0	0 1
D	1	1	1	0	1	1 1
E	0	1	0	1	0	1 0



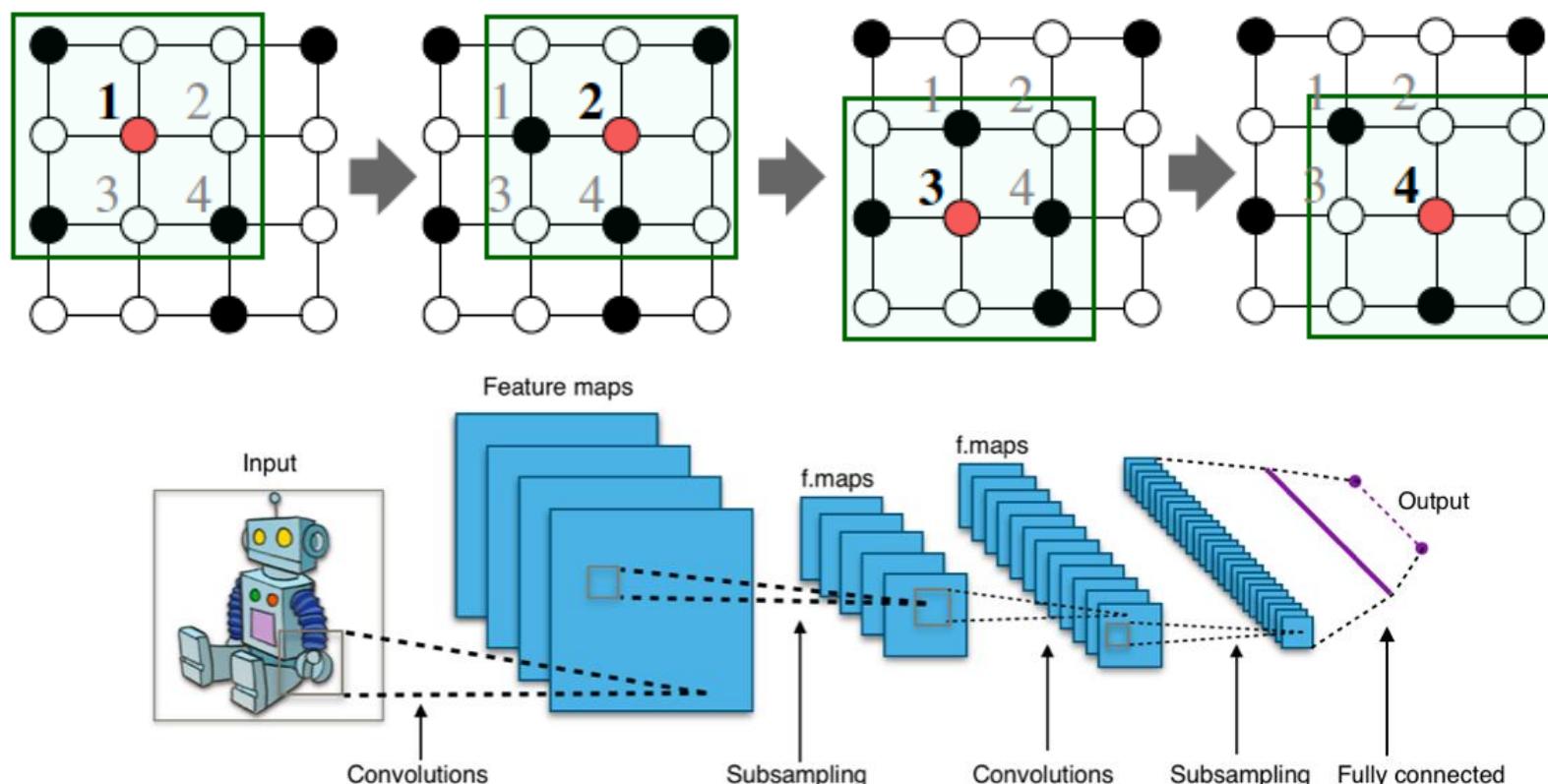
- Issues with this idea:

全连接NN要求输入维度固定

- $O(|V|)$ parameters 参数量过大 (过拟合)
- Not applicable to graphs of different sizes 无法泛化
- Sensitive to node ordering 不具备“变换不变性”

Idea: Convolutional Networks

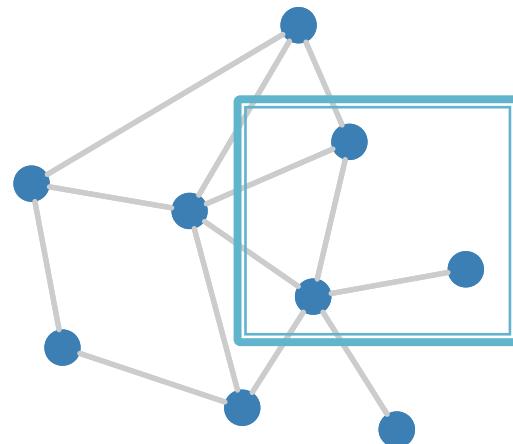
CNN on an image:



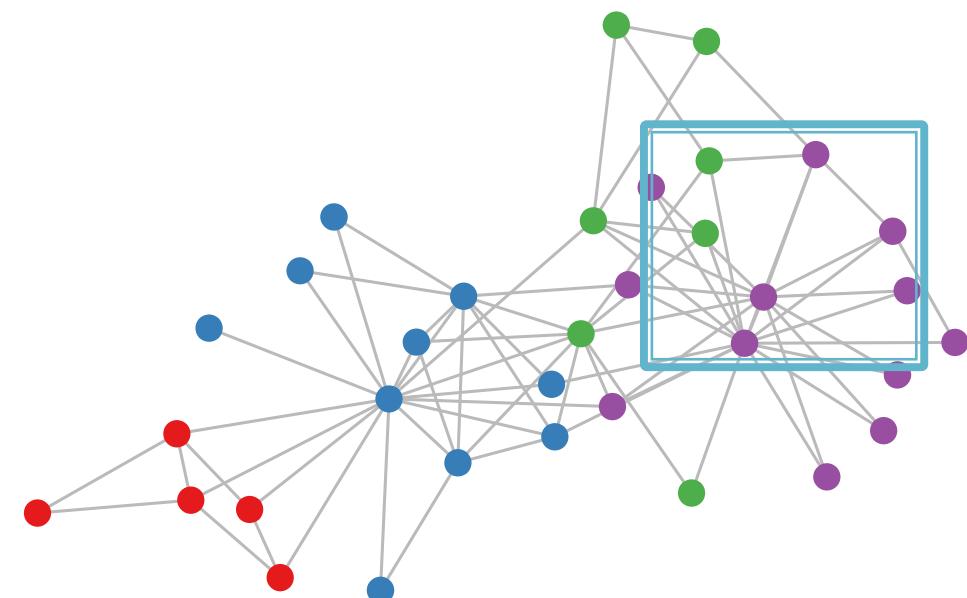
Goal is to generalize convolutions beyond simple lattices
Leverage node features/attributes (e.g., text, images)

Real-World Graphs

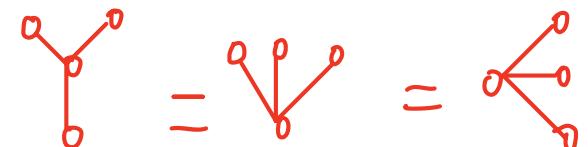
But our graphs look like this:



or this:



- There is no fixed notion of locality or sliding window on the graph
- Graph is permutation invariant



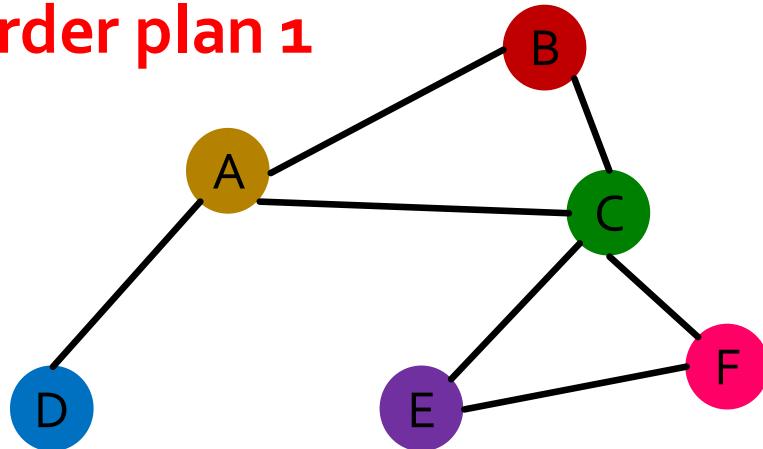
Permutation Invariance

- **Graph does not have a canonical order of the nodes!**
- We can have many different order plans.

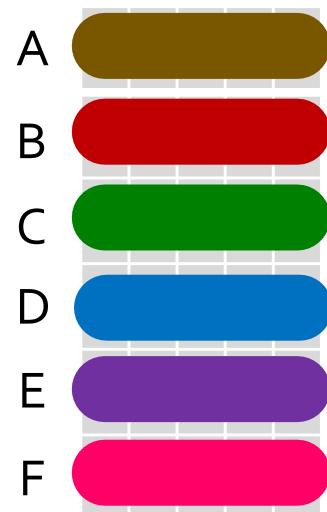
Permutation Invariance

- Graph does not have a canonical order of the nodes!

Order plan 1



Node features X_1



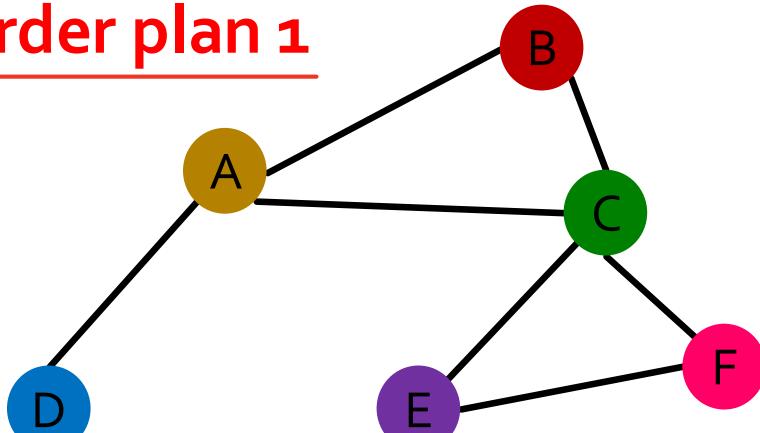
Adjacency matrix A_1

	A	B	C	D	E	F
A	1	0	1	0	0	0
B	0	1	1	0	0	0
C	1	1	0	1	1	0
D	0	0	0	1	0	0
E	0	0	1	0	0	1
F	0	0	0	0	1	1

Permutation Invariance

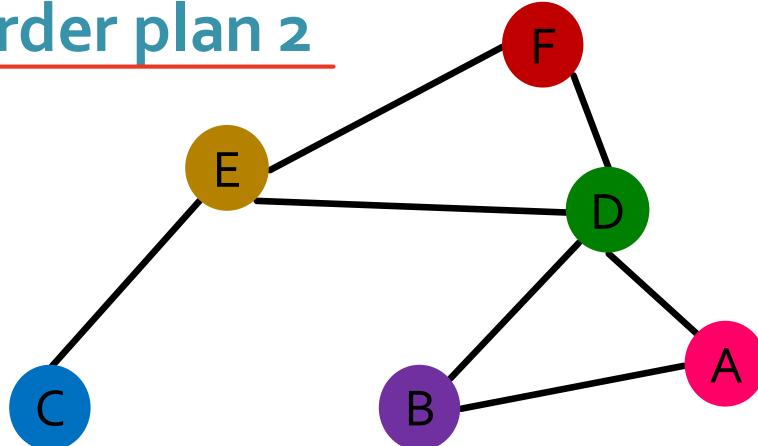
- Graph does not have a canonical order of the nodes!

Order plan 1

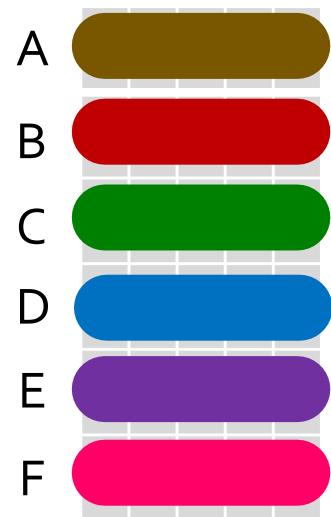


同一张图, 邻接矩阵却不同.

Order plan 2



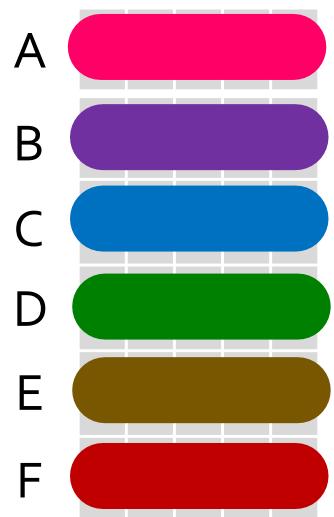
Node features X_1



Adjacency matrix A_1

	A	B	C	D	E	F
A	1	0	1	1	0	0
B	0	1	0	0	0	0
C	1	0	1	1	1	1
D	1	0	0	1	0	0
E	0	0	1	0	1	1
F	0	0	1	0	1	1

Node features X_2



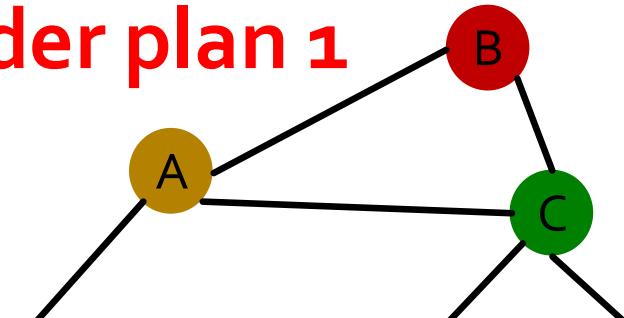
Adjacency matrix A_2

	A	B	C	D	E	F
A	1	0	0	0	0	0
B	0	1	0	0	0	0
C	0	0	1	0	1	1
D	0	0	1	1	0	0
E	0	0	1	0	1	1
F	0	0	1	0	1	1

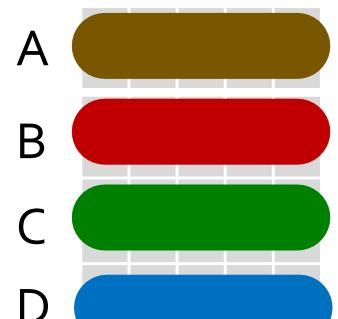
Permutation Invariance

- Graph does not have a canonical order of the nodes!

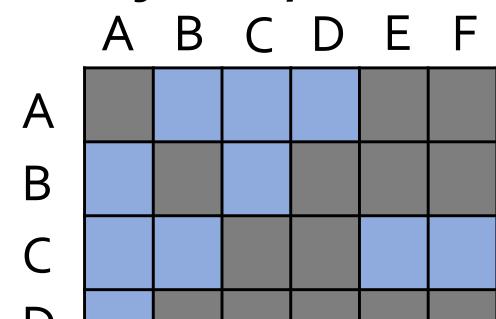
Order plan 1



Node feature X_1



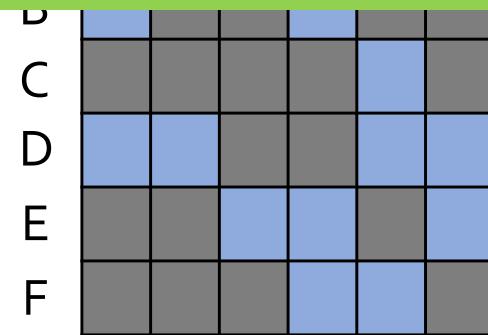
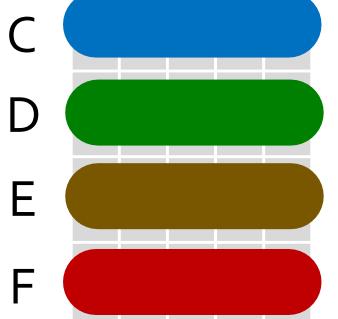
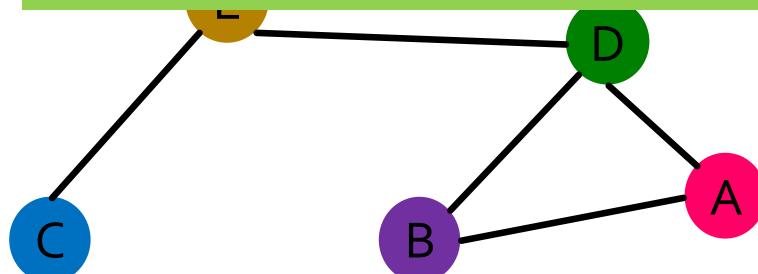
Adjacency matrix A_1



全图和节点的嵌入表示应与节点编号顺序无关

Graph and node representations

should be the same for Order plan 1
and Order plan 2



Permutation Invariance

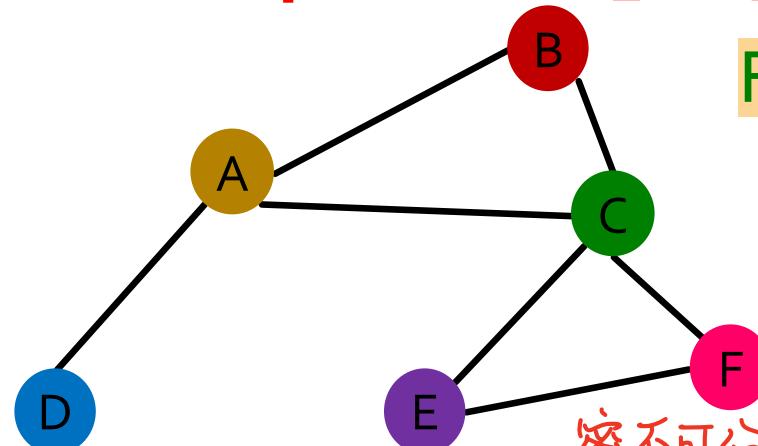
What does it mean by “graph representation is same for two order plans”?

- Consider we learn a function f that maps a graph $G = (A, X)$ to a vector \mathbb{R}^d then

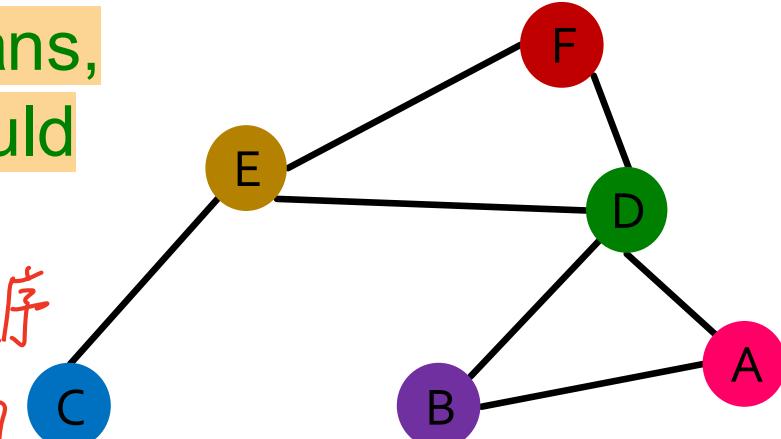
$$f(A_1, X_1) = f(A_2, X_2)$$

A is the adjacency matrix
 X is the node feature matrix

Order plan 1: A_1, X_1



Order plan 2: A_2, X_2



For two order plans,
output of f should
be the same!

邻接矩阵与编号顺序
密不可分, 用邻接矩阵行不通

Permutation Invariance

What does it mean by “graph representation is same for two order plans”?

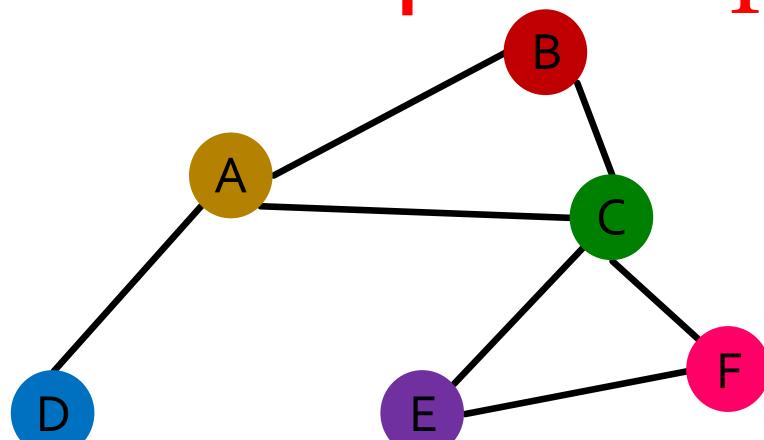
- Consider we learn a function f that maps a graph $G = (A, X)$ to a vector \mathbb{R}^d .
 A is the adjacency matrix
 X is the node feature matrix
- Then, if $f(A_i, X_i) = f(A_j, X_j)$ for any order 任意顺序 plan i and j , we formally say f is a **permutation invariant function**.
For a graph with $|V|$ nodes, there are $|V|!$ different order plans.
- Definition:** For any **graph** function $f: \mathbb{R}^{|V| \times m} \times \mathbb{R}^{|V| \times |V|} \rightarrow \mathbb{R}^d$, f is **permutation-invariant** if $f(A, X) = f(PAP^T, PX)$ for any permutation P .

Permutation P : a shuffle of the node order
Example: (A,B,C)->(B,C,A)

Permutation Equivariance

For node representation: We learn a function f that maps nodes of G to a matrix $\mathbb{R}^{m \times d}$.

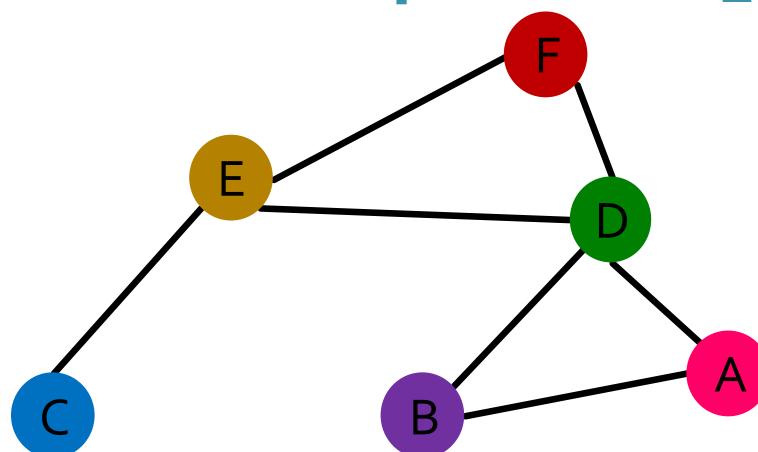
Order plan 1: A_1, X_1



$$f(A_1, X_1) =$$

A		
B		
C		
D		
E		
F		

Order plan 2: A_2, X_2



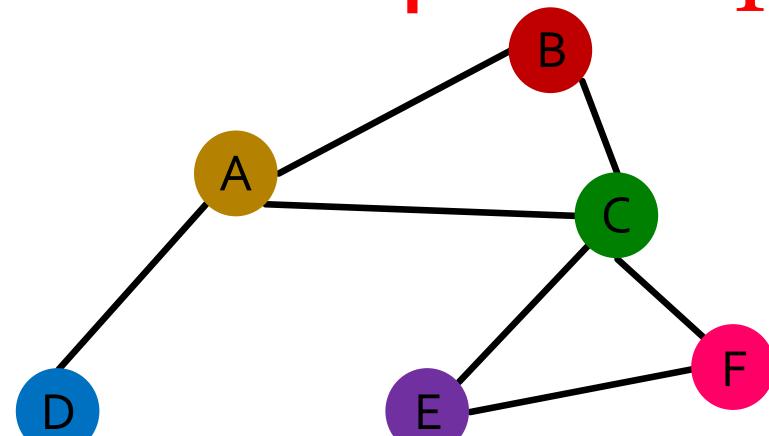
$$f(A_2, X_2) =$$

A		
B		
C		
D		
E		
F		

Permutation Equivariance

For node representation: We learn a function f that maps nodes of G to a matrix $\mathbb{R}^{m \times d}$.

Order plan 1: A_1, X_1

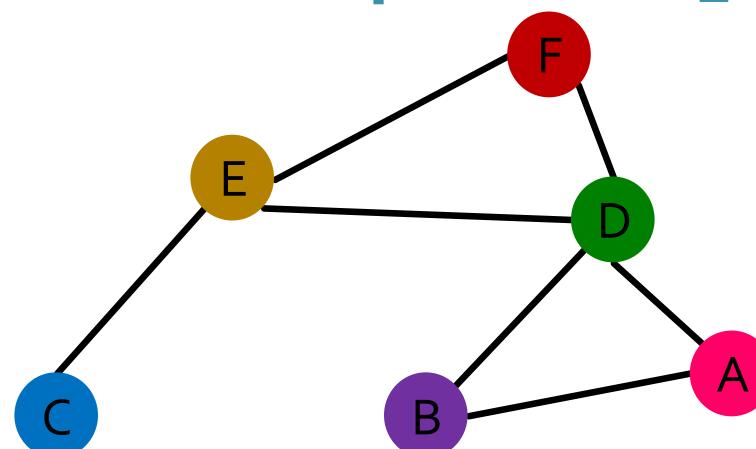


Representation vector
of the brown node A

A	■	■
B	■	■
C	■	■
D	■	■
E	■	■
F	■	■

$$f(A_1, X_1) =$$

Order plan 2: A_2, X_2



$$f(A_2, X_2) =$$

For two order plans, the vector of node at the same position in the graph is the same!

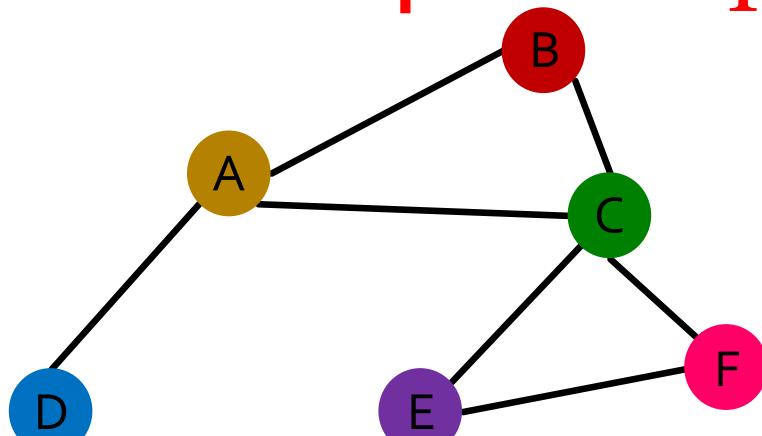
A	■	■
B	■	■
C	■	■
D	■	■
E	■	■
F	■	■

Representation vector
of the brown node E

Permutation Equivariance

For node representation: We learn a function f that maps nodes of G to a matrix $\mathbb{R}^{m \times d}$.

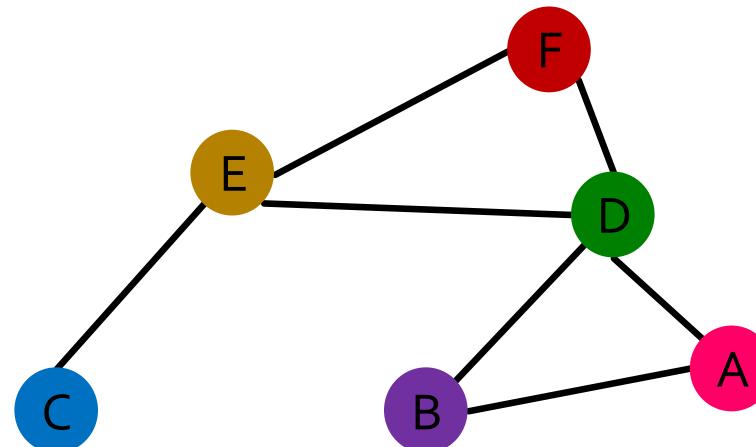
Order plan 1: A_1, X_1



$$f(A_1, X_1) = \begin{matrix} & \text{A} & \text{B} & \text{C} & \text{D} & \text{E} & \text{F} \\ \text{A} & \text{brown} & & & & & \\ \text{B} & & \text{red} & & & & \\ \text{C} & & & \text{green} & & & \\ \text{D} & & & & \text{blue} & & \\ \text{E} & & & & & \text{purple} & \\ \text{F} & & & & & & \text{pink} \end{matrix}$$

Representation vector of the green node C

Order plan 2: A_2, X_2



$$f(A_2, X_2) = \begin{matrix} & \text{A} & \text{B} & \text{C} & \text{D} & \text{E} & \text{F} \\ \text{A} & \text{red} & & & & & \\ \text{B} & & \text{purple} & & & & \\ \text{C} & & & \text{blue} & & & \\ \text{D} & & & & \text{green} & & \\ \text{E} & & & & & \text{brown} & \\ \text{F} & & & & & & \text{red} \end{matrix}$$

Representation vector of the green node D

For two order plans, the vector of node at the same position in the graph is the same! F

Permutation Equivariance

For node representation

- Consider we learn a function f that maps a graph $G = (A, X)$ to a matrix $\mathbb{R}^{m \times d}$
- If the output vector of a node at the same position in the graph remains unchanged for any order plan, we say f is **permutation equivariant**.
- **Definition:** For any **node** function $f: \mathbb{R}^{|V| \times m} \times \mathbb{R}^{|V| \times |V|} \rightarrow \mathbb{R}^{|V| \times m}$, f is **permutation-equivariant** if $Pf(A, X) = f(PAP^T, PX)$ for any permutation P .

Summary: Invariance and Equivariance

■ Permutation-invariant

$$f(A, X) = f(PAP^T, PX)$$

Permute the input, the output stays the same.
(map a graph to a vector)

■ Permutation-equivariant

$$Pf(A, X) = f(PAP^T, PX)$$

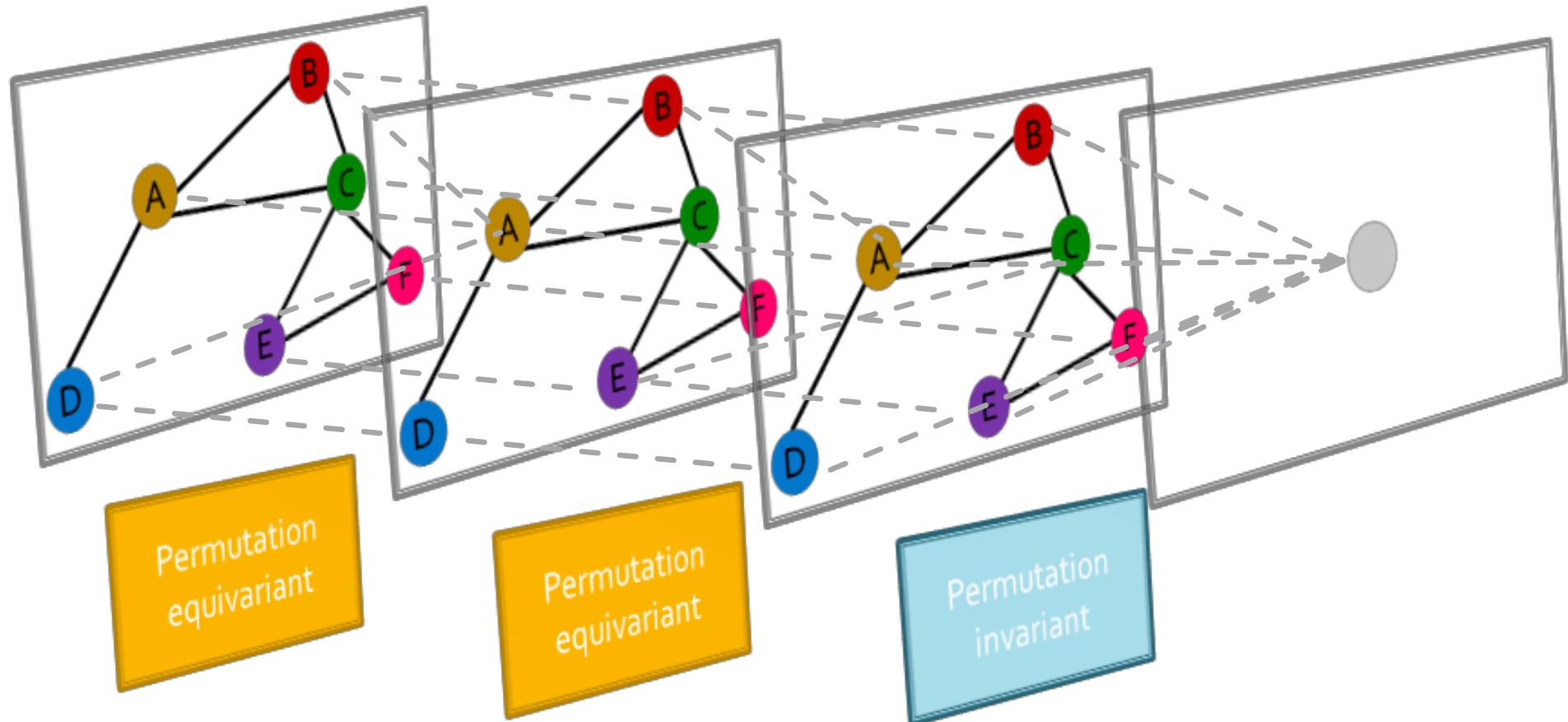
Permute the input, output also permutes accordingly.
(map a graph to a matrix)

■ Examples:

- $f(A, X) = \mathbf{1}^T X$: Permutation-**invariant**
 - Reason: $f(PAP^T, PX) = \mathbf{1}^T \mathbf{P}X = \mathbf{1}^T X = f(A, X)$
- $f(A, X) = X$: Permutation-**equivariant**
 - Reason: $f(PAP^T, PX) = \mathbf{P}X = Pf(A, X)$
- $f(A, X) = AX$: Permutation-**equivariant**
 - Reason: $f(PAP^T, PX) = PAP^T PX = \mathbf{P}AX = Pf(A, X)$

Graph Neural Network Overview

- Graph neural networks consist of multiple permutation equivariant / invariant functions.



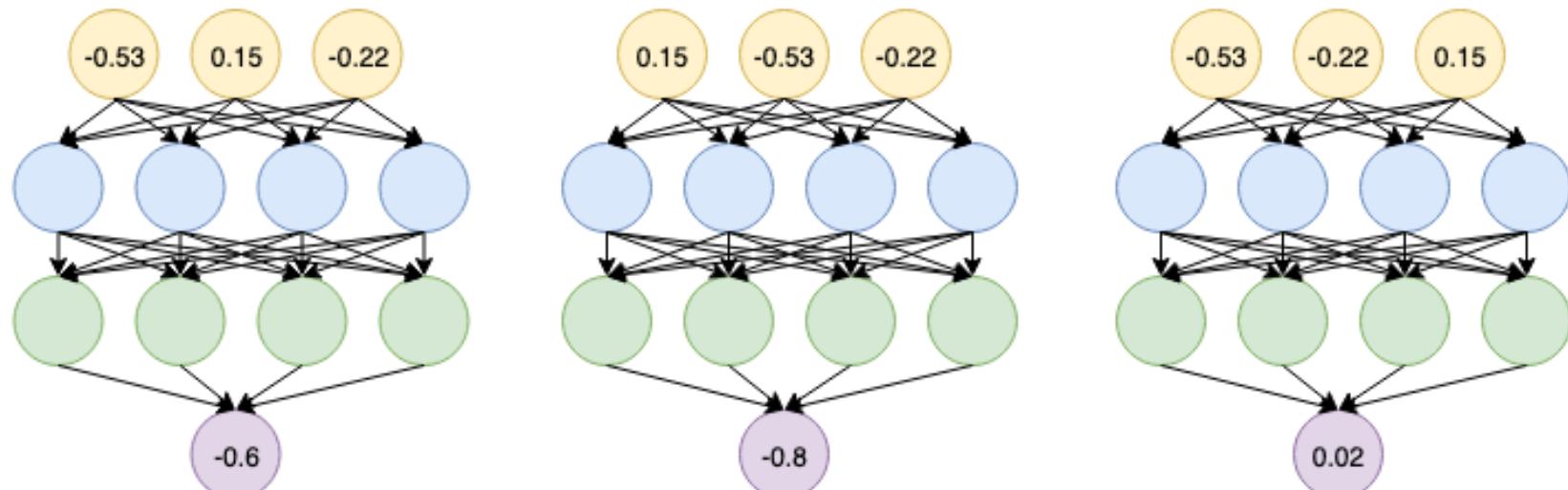
Graph Neural Network Overview

Are other neural network architectures, e.g.,
MLPs, permutation invariant / equivariant?

- No.

MLP 没有顺序不变性

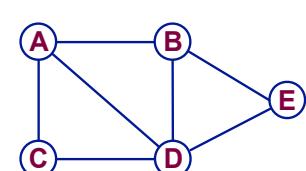
Switching the order of the
input leads to different
outputs!



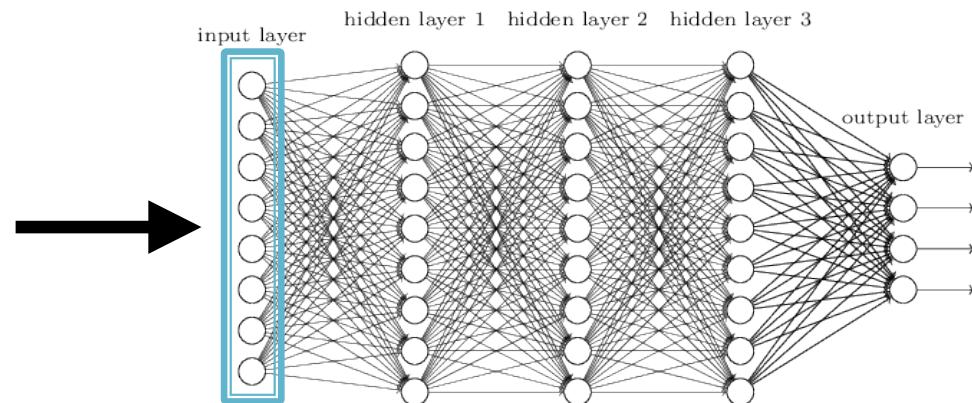
Graph Neural Network Overview

Are other neural network architectures, e.g.,
MLPs, permutation invariant / equivariant?

■ No.



	A	B	C	D	E	Feat
A	0	1	1	1	0	1 0
B	1	0	0	1	1	0 0
C	1	0	0	1	0	0 1
D	1	1	1	0	1	1 1
E	0	1	0	1	0	1 0



直接把邻接矩阵输入 NN

This explains why the naïve MLP approach
fails for graphs! 不可行

Graph Neural Network Overview

- Are any neural network architecture, e.g.,

Next: Design graph neural networks that are permutation invariant / equivariant by passing and aggregating information from neighbors!

消息传递图神经网络

?

Outline of Today's Lecture

1. Basics of deep learning



2. Deep learning for graphs



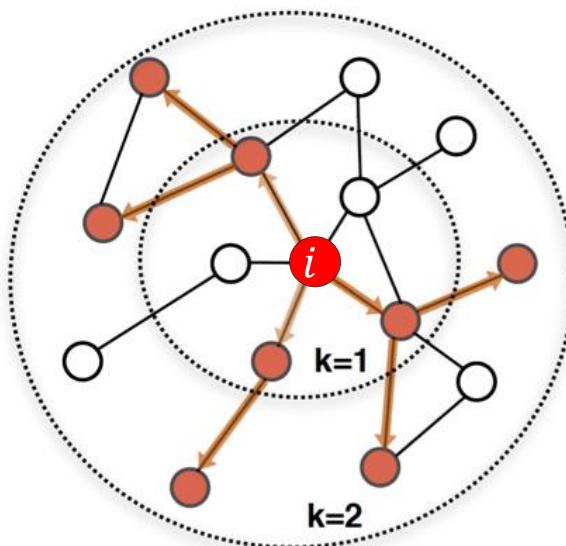
3. Graph Convolutional Networks



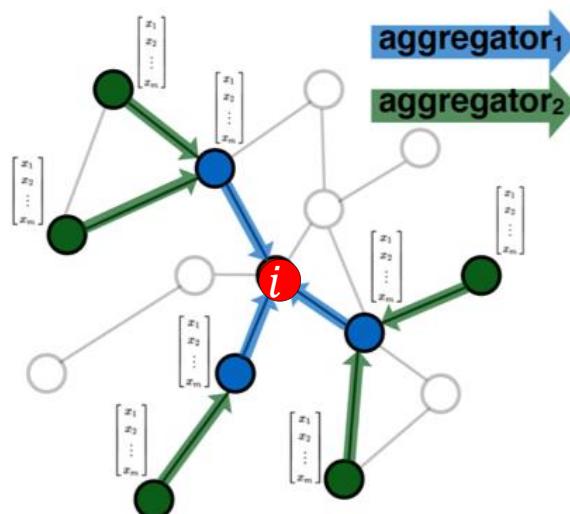
4. GNNs subsume CNNs

Graph Convolutional Networks

Idea: Node's neighborhood defines a computation graph



Determine node computation graph



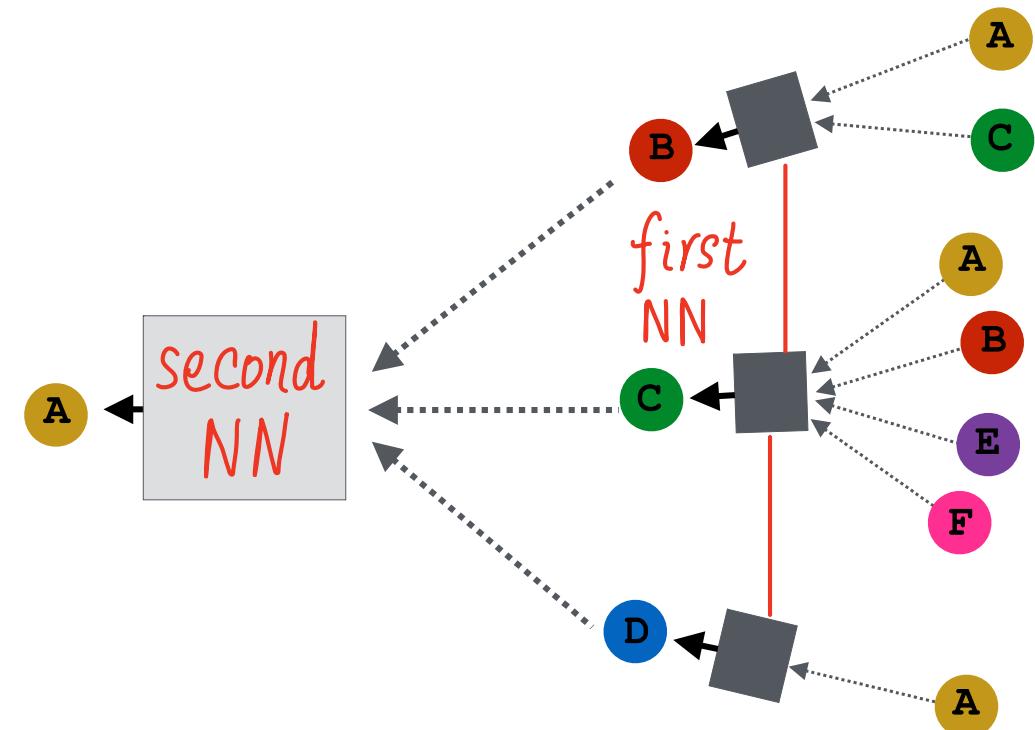
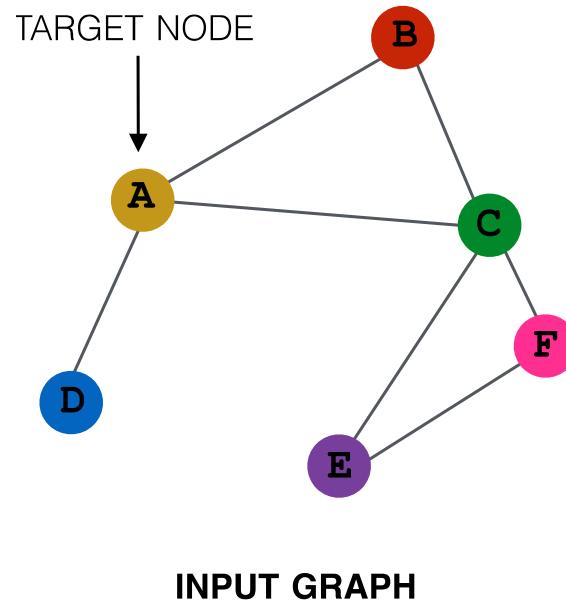
Propagate and transform information

Learn how to propagate information across the graph to compute node features

Idea: Aggregate Neighbors

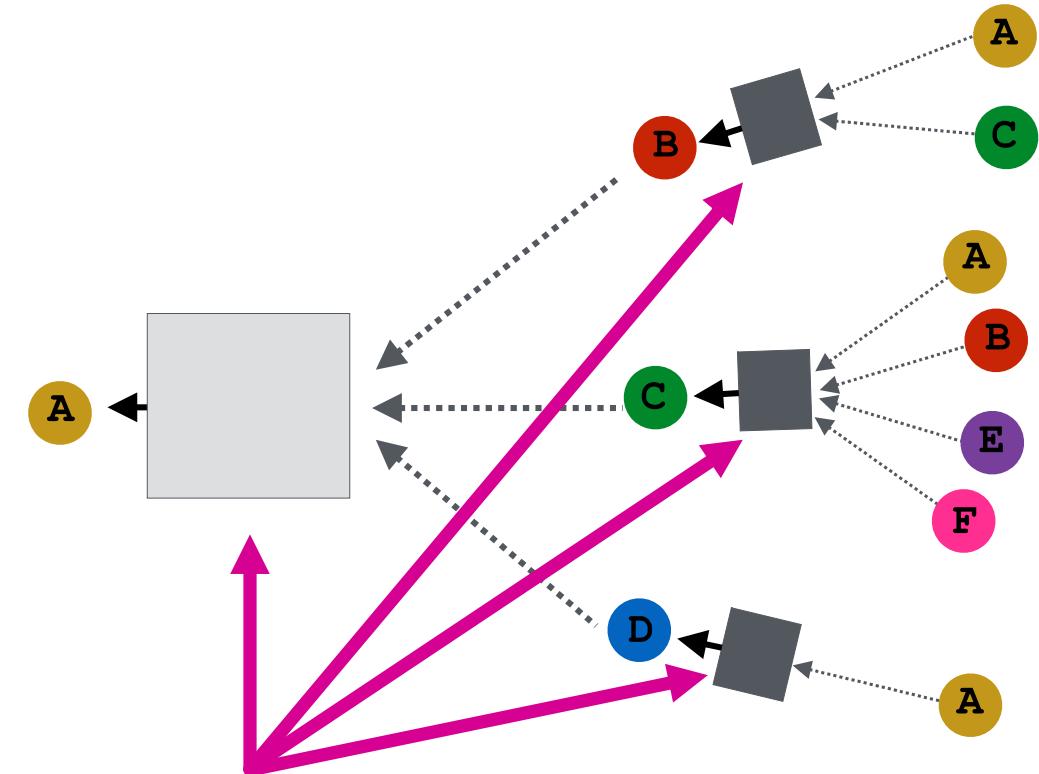
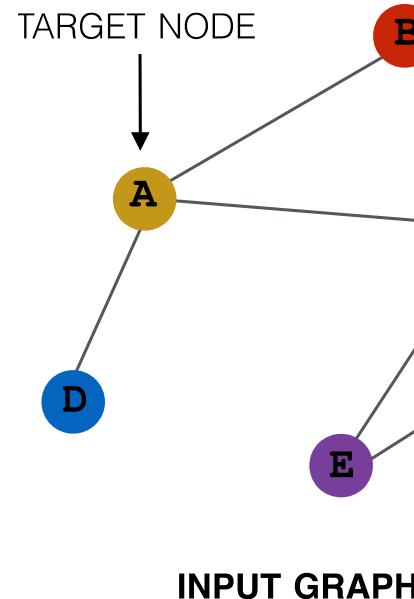
- **Key idea:** Generate node embeddings based on local network neighborhoods

通过局部邻域构建计算图



Idea: Aggregate Neighbors

- **Intuition:** Nodes aggregate information from their neighbors using neural networks

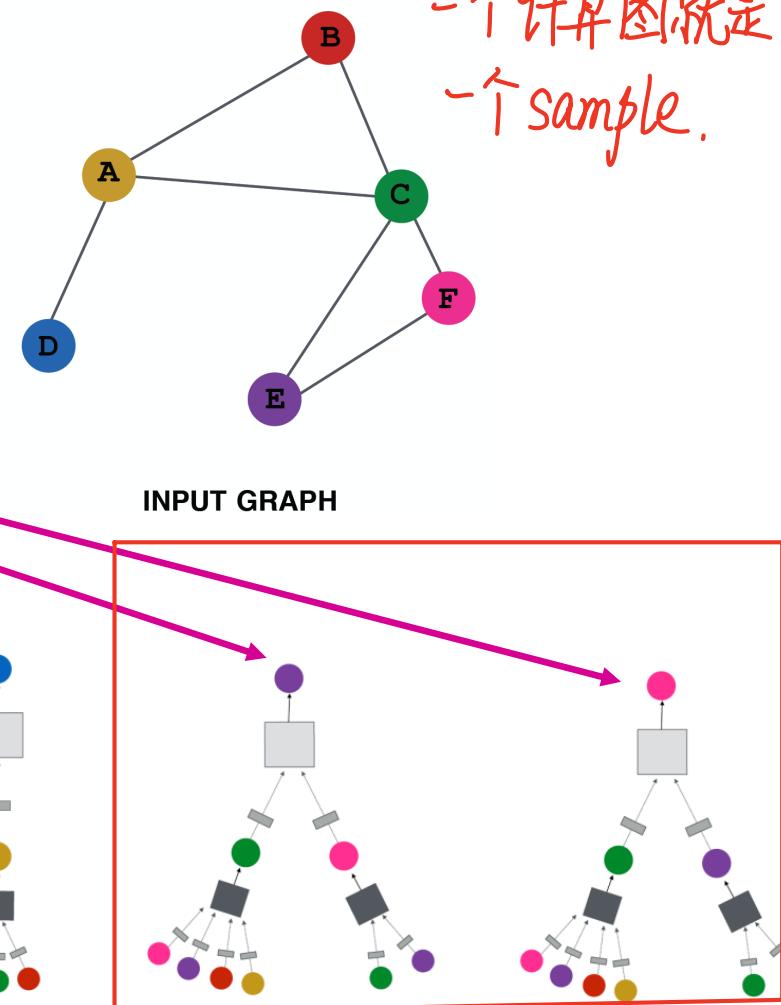
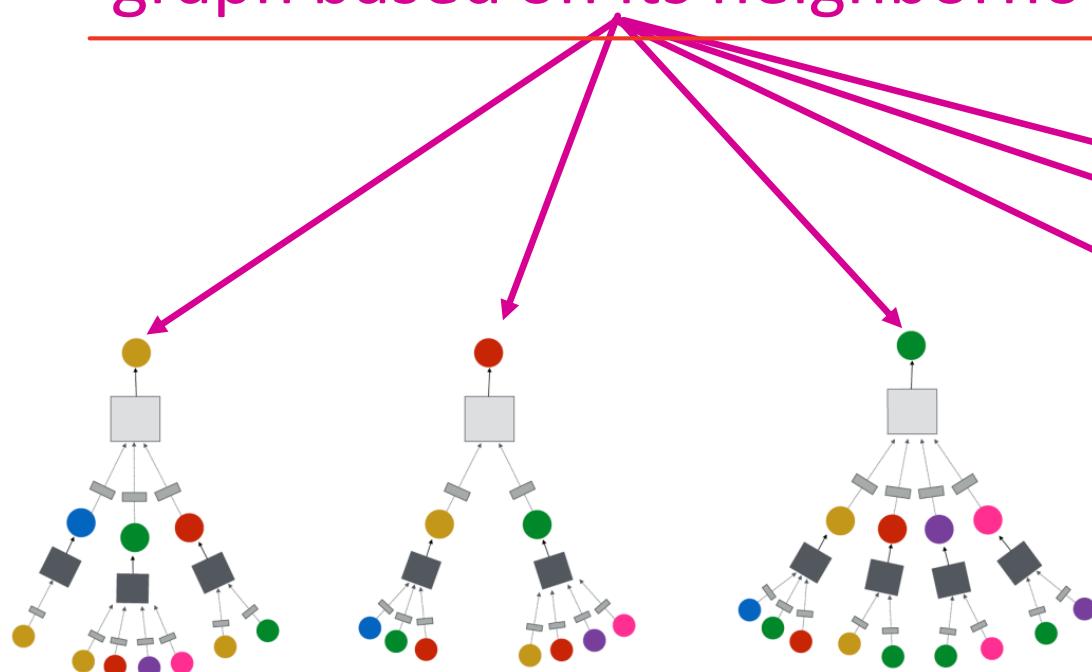


Idea: Aggregate Neighbors

- **Intuition:** Network neighborhood defines a computation graph

每个节点分别构建自己的计算图

Every node defines a computation graph based on its neighborhood!

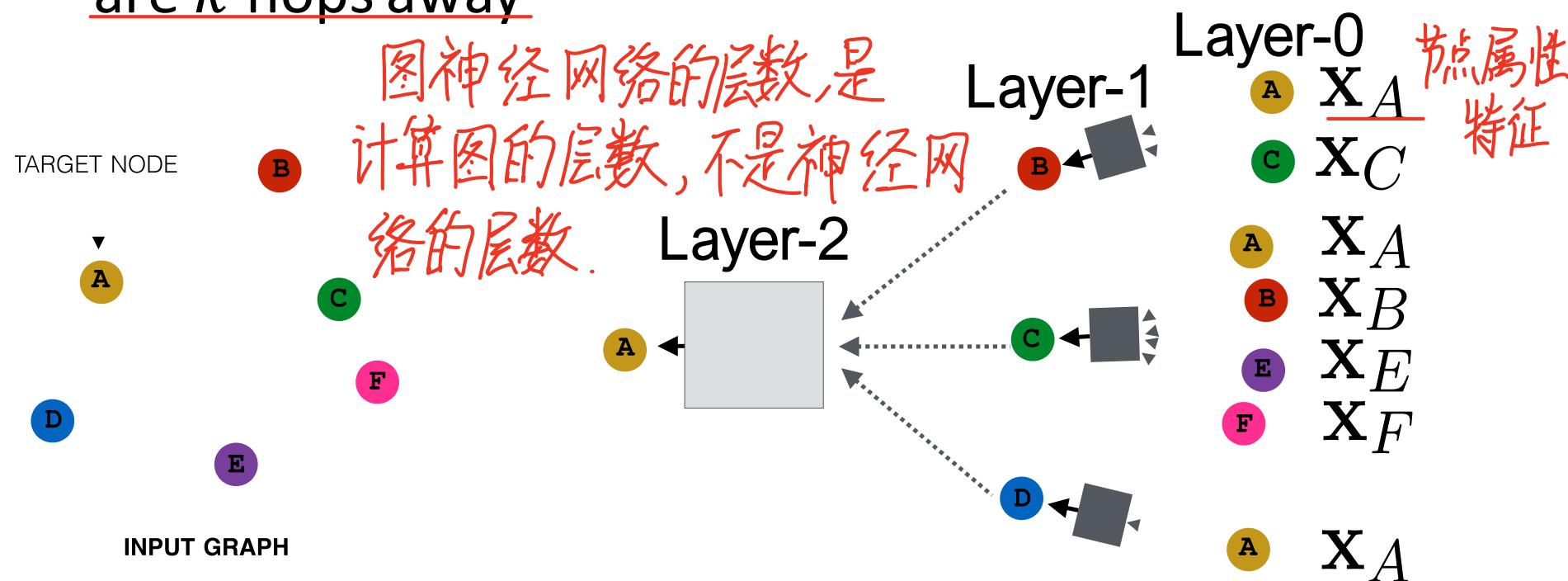


一个计算图就是
一个 sample.

计算图结构相同 (结构, 功能, 角色)

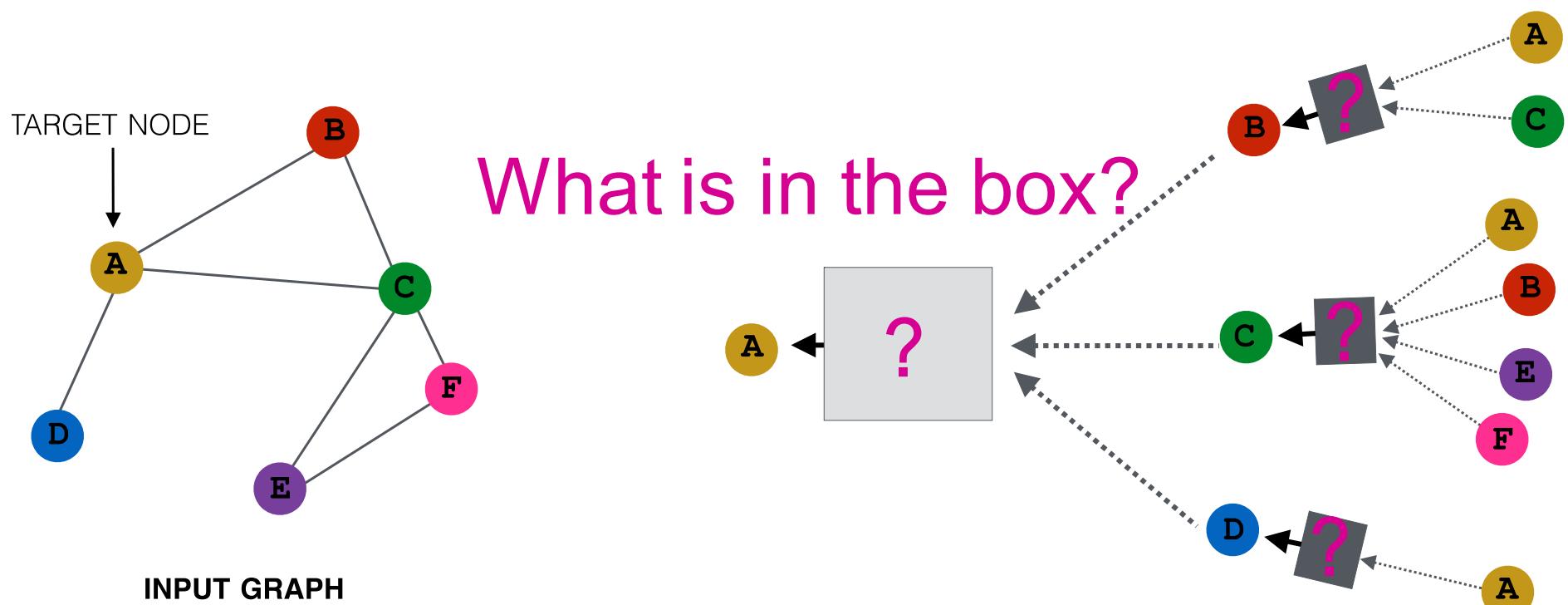
Deep Model: Many Layers

- 理论上任意深度
- Model can be of arbitrary depth:
 - Nodes have embeddings at each layer
 - Layer-0 embedding of node v is its input feature, x_v
 - Layer- k embedding gets information from nodes that are k hops away
- “六度空间”理论: 深度最多6层够用
(一般2层就好)



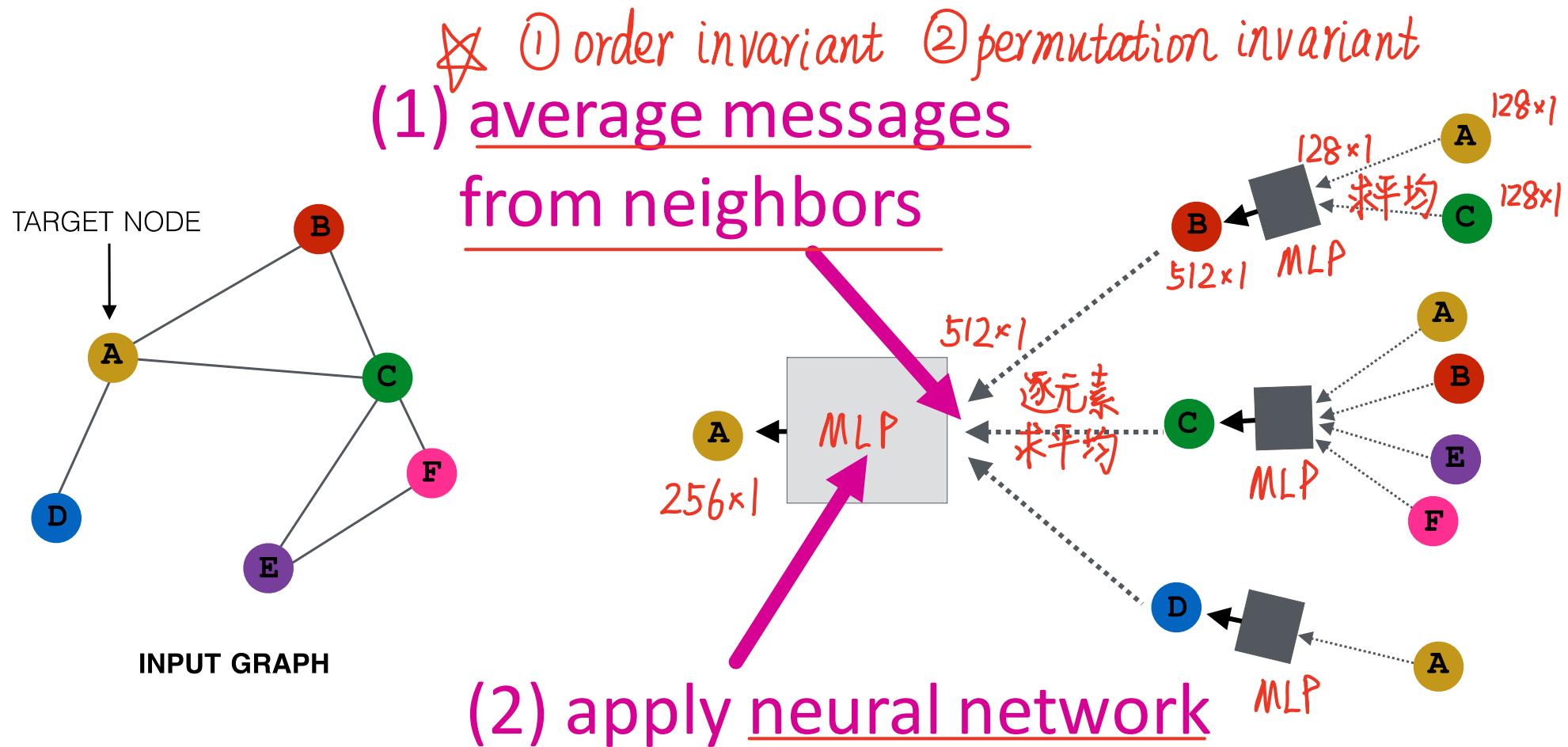
Neighborhood Aggregation

- **Neighborhood aggregation:** Key distinctions are in how different approaches aggregate information across the layers



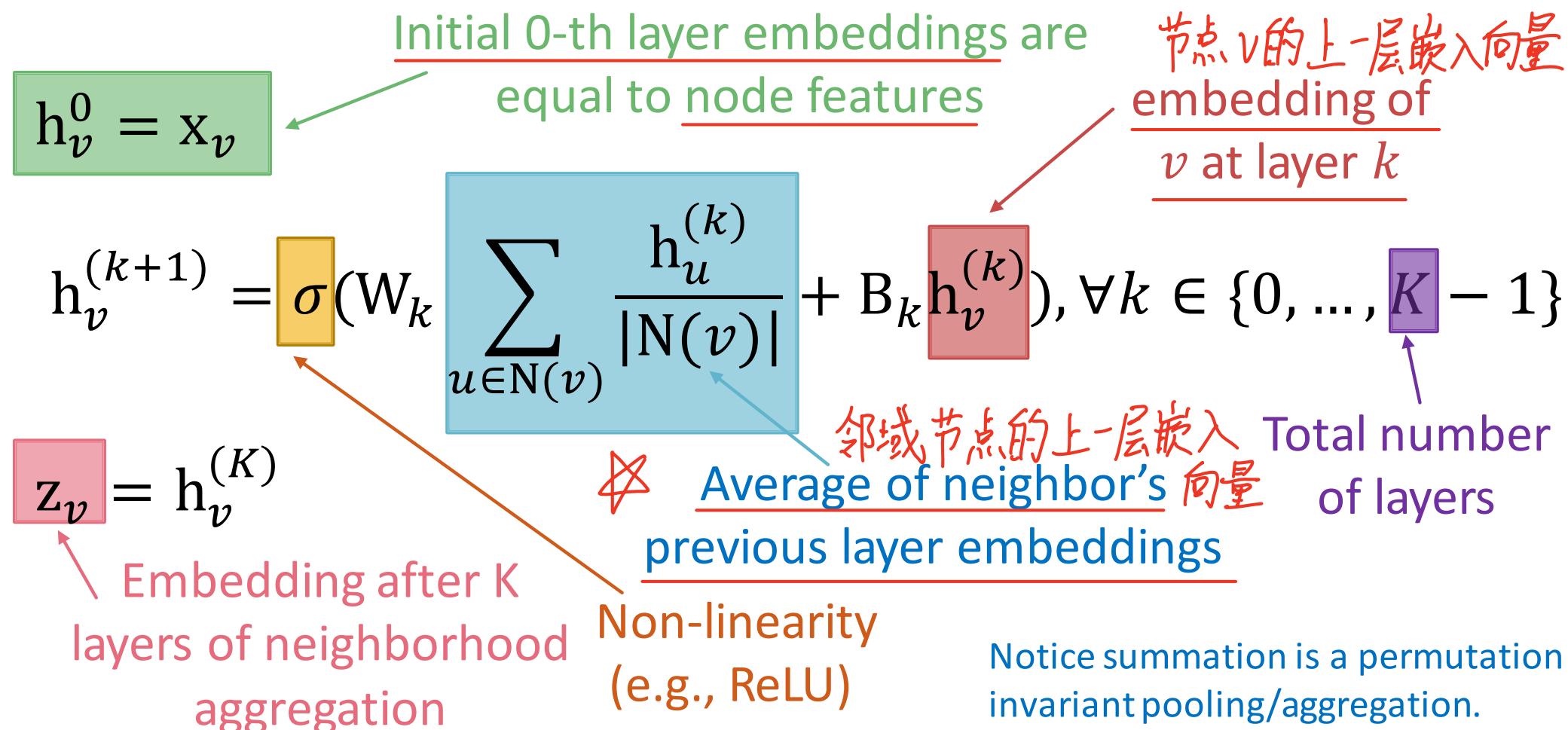
Neighborhood Aggregation

- **Basic approach:** Average information from neighbors and apply a neural network



The Math: Deep Encoder

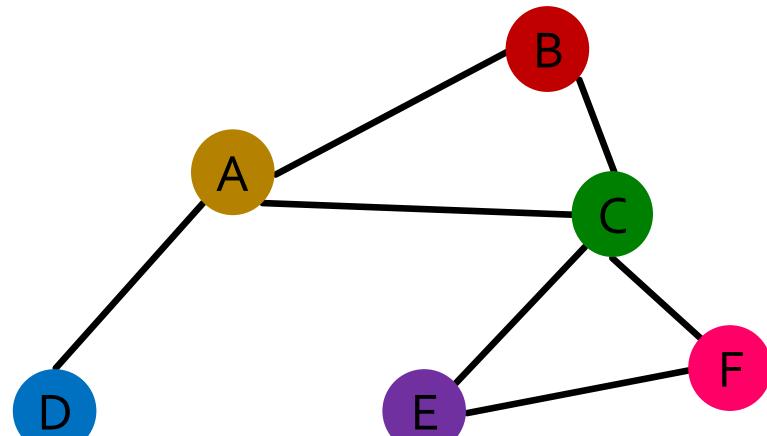
- **Basic approach:** Average neighbor messages and apply a neural network



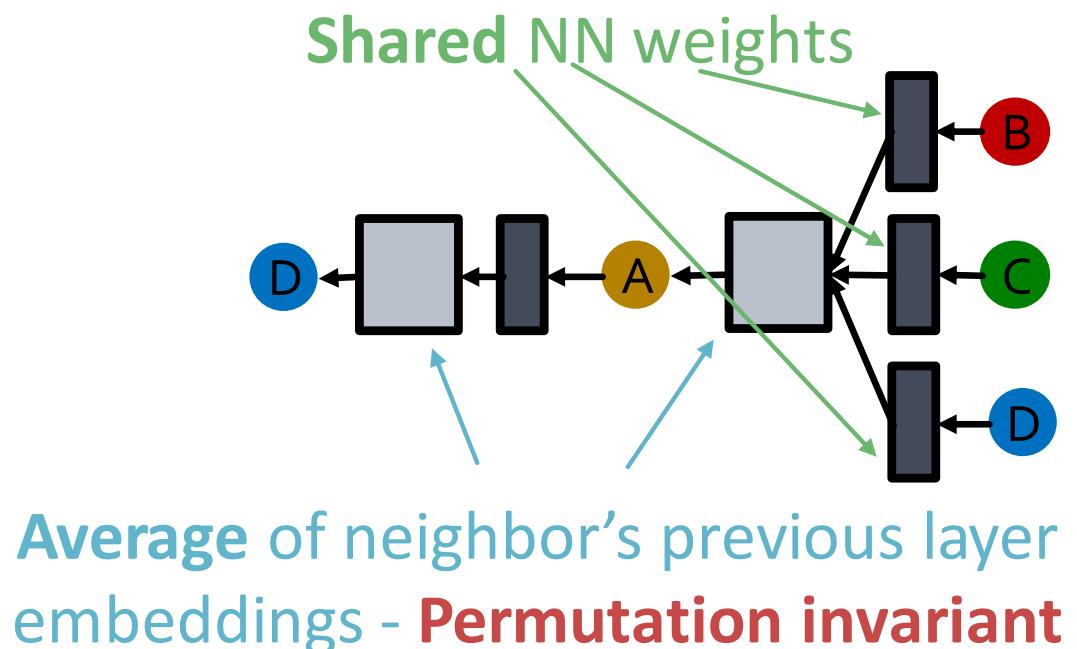
GCN: Invariance and Equivariance

What are the **invariance** and **equivariance** properties for a GCN?

- Given a node, the GCN that computes its embedding is **permutation invariant**



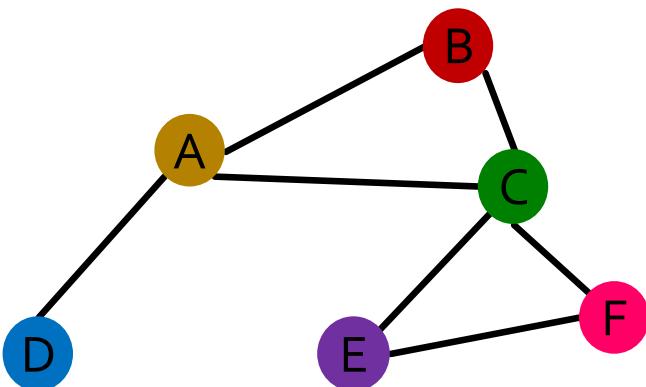
Target Node



GCN: Invariance and Equivariance

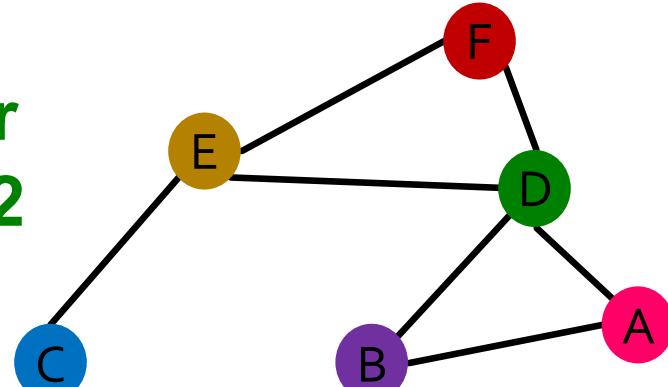
- Considering all nodes in a graph, GCN computation is **permutation equivariant**

Order plan 1



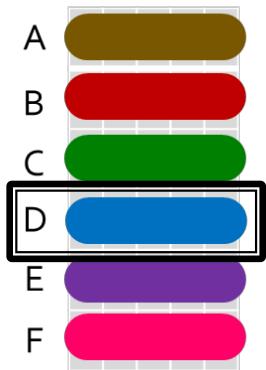
Target Node

Order plan 2

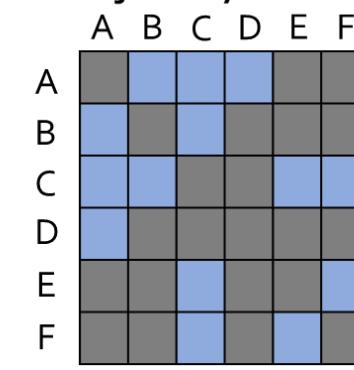


Target Node

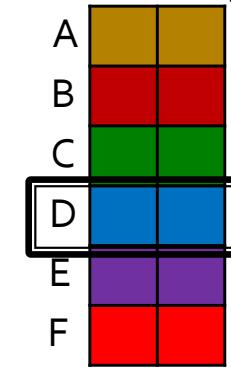
Node feature X_1



Adjacency matrix A_1

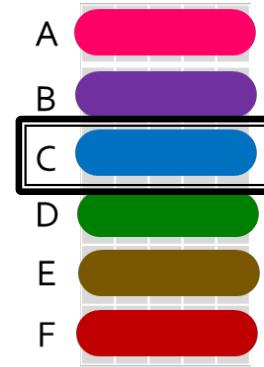


Embeddings H_1

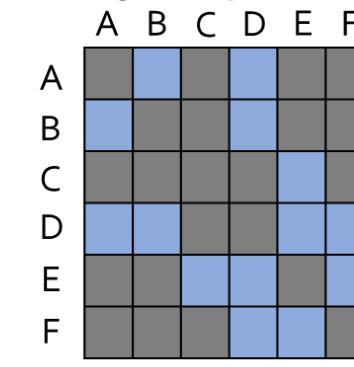


Permute the input, the output also permutes accordingly - **permutation equivariant**

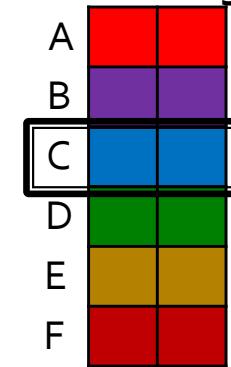
Node feature X_2



Adjacency matrix A_2



Embeddings H_2



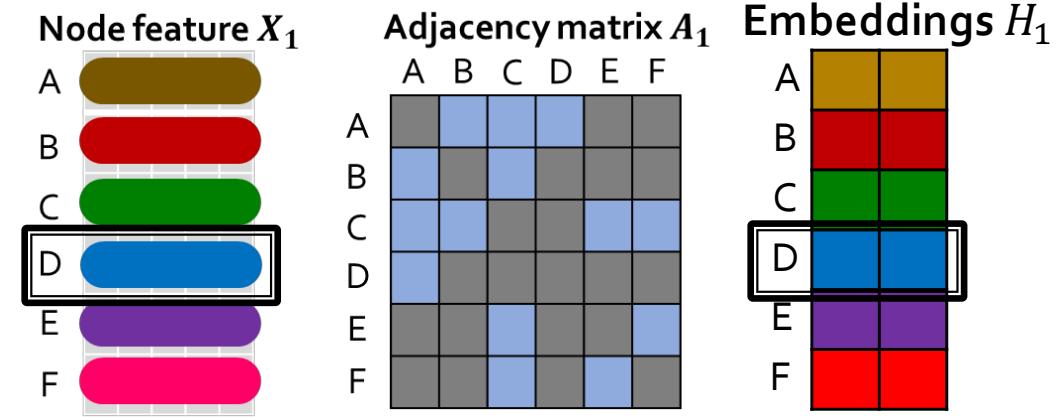
GCN: Invariance and Equivariance

- Considering all nodes in a graph, GCN computation is **permutation equivariant**

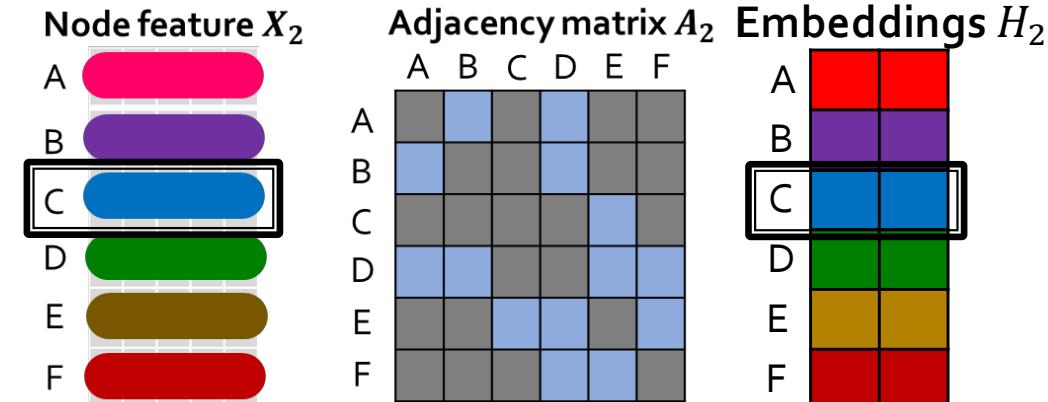
Detailed reasoning:

1. The rows of **input node features** and **output embeddings** are aligned
2. We know computing the embedding of a **given node** with GCN is **invariant**.
3. So, after permutation, the **location of a given node in the input node feature matrix** is changed, and the **the output embedding of a given node stays the same** (the colors of node feature and embedding are **matched**)

This is permutation equivariant

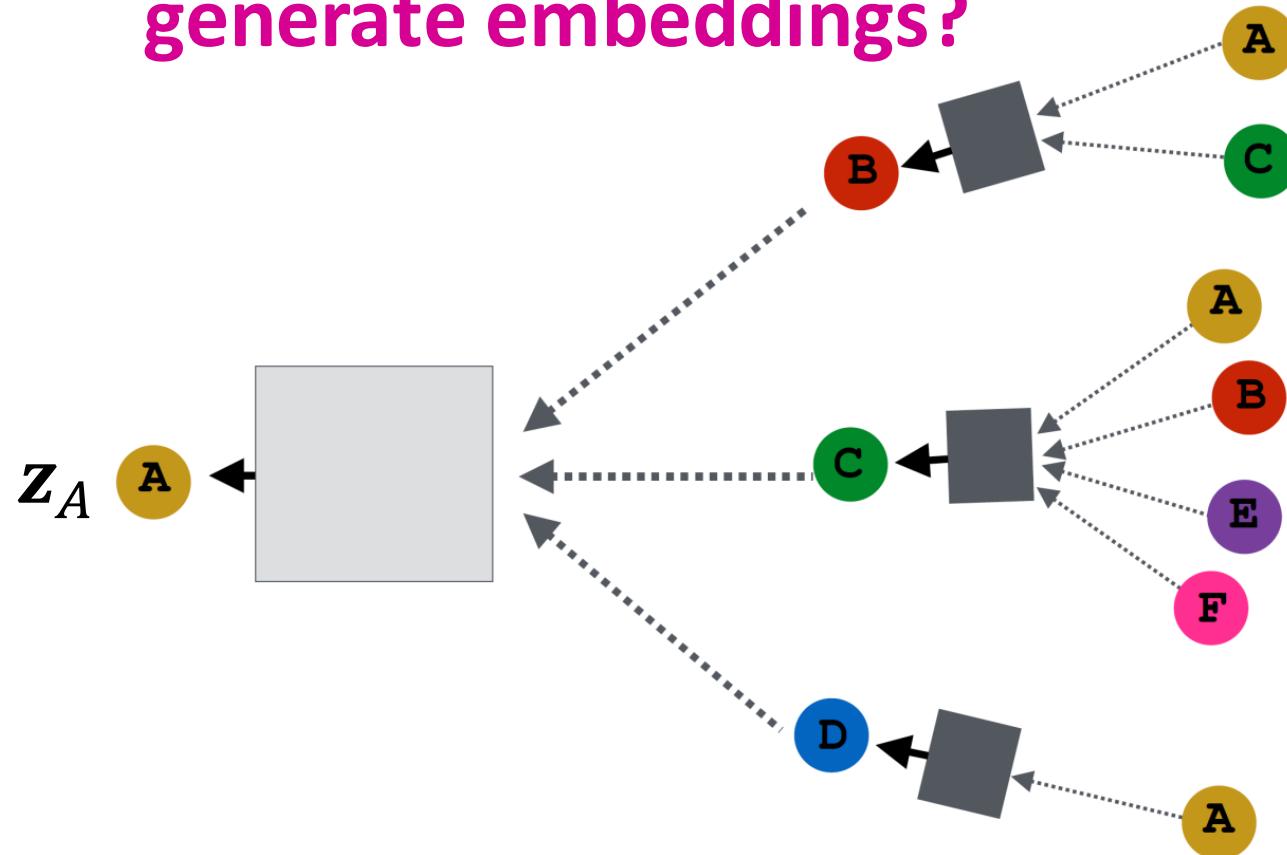


Permute the input, the output also permutes accordingly - permutation equivariant



Training the Model

How do we train the GCN to generate embeddings?



Need to define a loss function on the embeddings.

Model Parameters

Trainable weight matrices
(i.e., what we learn)

$$\begin{aligned} h_v^{(0)} &= x_v \\ h_v^{(k+1)} &= \sigma \left(\sum_{u \in N(v)} \frac{h_u^{(k)}}{|N(v)|} + B_k h_v^{(k)} \right), \forall k \in \{0..K-1\} \\ z_v &= h_v^{(K)} \end{aligned}$$

Final node embedding

We can feed these **embeddings into any loss function** and run SGD to **train the weight parameters**

- h_v^k : the hidden representation of node v at layer k
- W_k : weight matrix for neighborhood aggregation
 - B_k : weight matrix for transforming hidden vector of self

Matrix Formulation (1)

- Many aggregations can be performed efficiently by (sparse) matrix operations

- Let $H^{(k)} = [h_1^{(k)} \dots h_{|V|}^{(k)}]^T$ Matrix of hidden embeddings $H^{(k-1)}$

- Then: $\sum_{u \in N_v} h_u^{(k)} = A_{v,:} H^{(k)}$ 选出 v 节点邻域 的 embedding

- Let D be diagonal matrix where

$$D_{v,v} = \text{Deg}(v) = |N(v)|$$

- The inverse of D : D^{-1} is also diagonal:

$$D_{v,v}^{-1} = 1/|N(v)|$$

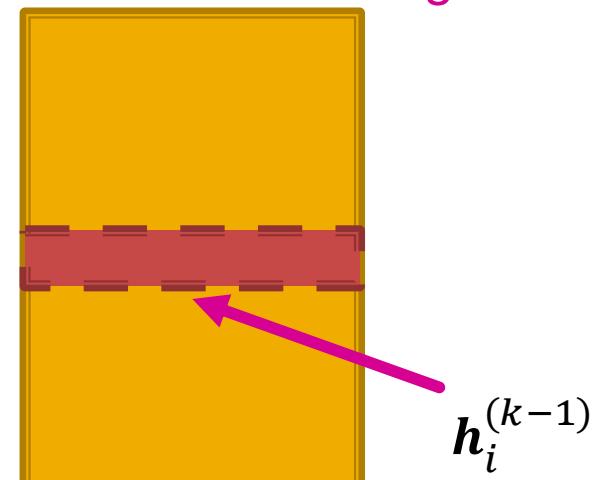
- Therefore,

$$\sum_{u \in N(v)} \frac{h_u^{(k-1)}}{|N(v)|}$$

$$D = \begin{bmatrix} d(1) & 0 & \dots \\ 0 & d(2) & \dots \\ \dots & \dots & \dots \end{bmatrix}$$

$$H^{(k+1)} = D^{-1} A H^{(k)}$$

求和
求平均



$$h_v^{(0)} = \chi_v$$

$$h_v^{(k+1)} = \delta (W_k \sum_{u \in N(v)} \frac{h_u^{(k)}}{|N(v)|})$$

$$z_v = h_v^{(k)}$$

$$\sum_{u \in N(v)} \frac{h_u^{(k-1)}}{|N(v)|} \rightarrow D^{-1} A H^{(k)}$$



Row Normalized Matrix
最大特征值为1

⇒ 向量经过矩阵变换，幅值不会变小

* 只按自己的 degree，对所有渠道来的信息求平均。
没有考虑对方的连接数

改进 \rightarrow $A_{naive} = D^{-1} A D^{-1}$

$\lambda \in (-1, 1) \rightarrow$ 输入向量幅值变小

改进 \rightarrow $A_{sym} = D^{-\frac{1}{2}} A D^{-\frac{1}{2}}$

最大特征值为1 \rightarrow 输入向量幅值不变

$$\sum_{u \in N(v)} \frac{h_u^{(k-1)}}{|N(v)|} \longrightarrow D^{-1} A H^{(k)}$$

当图给定时, \tilde{A} 就给定了.

$$\longrightarrow \boxed{D^{-\frac{1}{2}} A D^{-\frac{1}{2}}} H^{(k)}$$

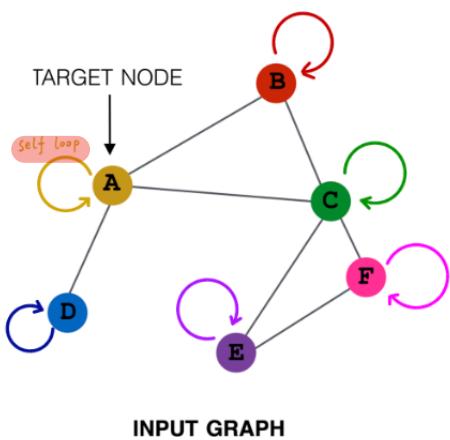
Normalized Adjacency Matrix
 Normalized Diffusion Matrix \tilde{A}
 Symmetric Normalized Matrix

图卷积神经网络GCN

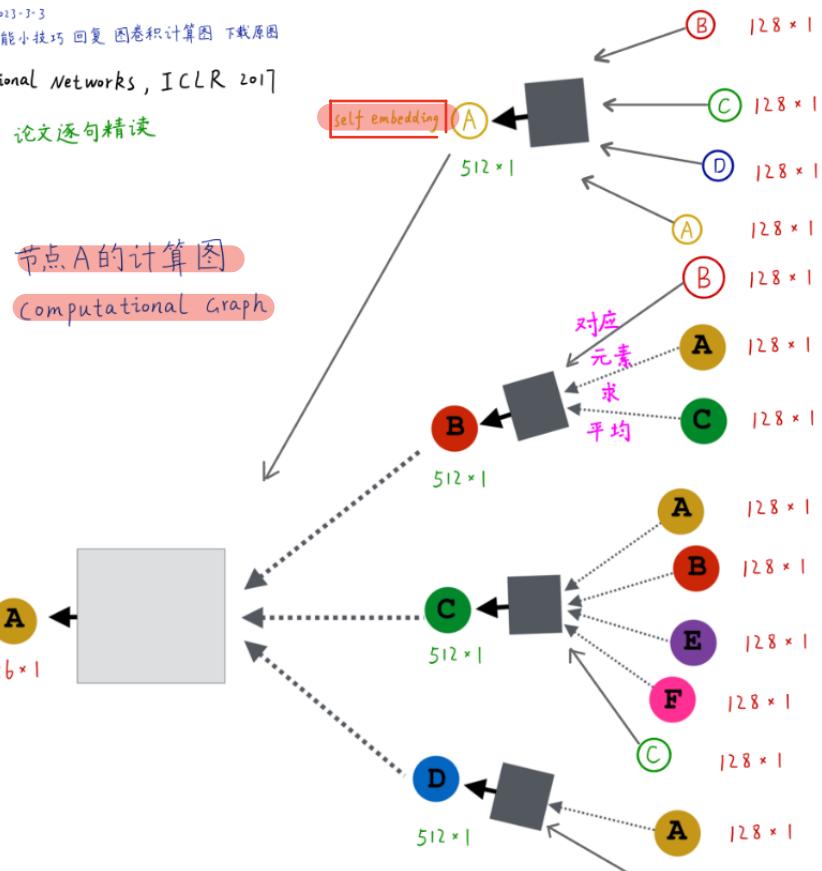
同济子豪兄 2023-3-3
公众号 人工智能小技巧 回复 因卷积计算图 下载原因

Semi-Supervised classification with Graph Convolutional Networks, ICLR 2017

斯坦福CS224W图机器学习公开课 中文精讲 论文逐句精读



	A	B	C	D	E	F
A	1	1	1	1	0	0
B	1	1	1	0	0	0
C	1	1	1	0	1	1
D	1	0	0	1	0	0
E	0	0	1	0	1	1
F	0	0	1	0	1	1



图神经网络优点

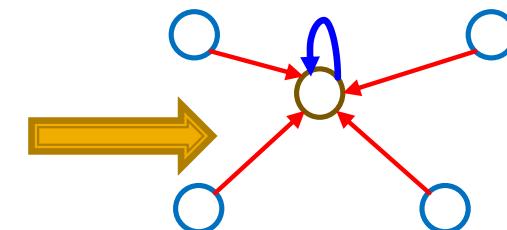
- 深度学习拟合学习能力强，表示学习得到的嵌入向量质量高
- 站在深度学习巨人肩膀上
- 归纳式学习能力 Inductive Learning: 泛化到新节点、新图
- 参数量少、所有计算图共享神经网络
- 利用节点属性特征
- 利用节点标注类别
- 区分节点结构功能角色(桥接、中枢、外围边缘)
- 只需寥寥几层，就可以让任意两个节点相互影响

Matrix Formulation (2)

- Re-writing update function in matrix form:

$$H^{(k+1)} = \sigma(\tilde{A}H^{(k)}W_k^T + H^{(k)}B_k^T)$$

where $\tilde{A} = D^{-1}A$



$$H^{(k)} = [h_1^{(k)} \dots h_{|V|}^{(k)}]^T$$

- Red: neighborhood aggregation
- Blue: self transformation

- In practice, this implies that efficient sparse matrix multiplication can be used (\tilde{A} is sparse)
- **Note:** not all GNNs can be expressed in matrix form, when aggregation function is complex

How to Train A GNN

- Node embedding \mathbf{z}_v is a function of input graph
- **Supervised setting:** we want to minimize the loss

\mathcal{L} (see also Slide 15): f : 分类或回归预测头

$$\min_{\Theta} \mathcal{L}(y, f(\mathbf{z}_v))$$

- y : node label
- \mathcal{L} could be L2 if y is real number, or cross entropy if y is categorical

- **Unsupervised setting:**

- No node label available

自监督

- Use the graph structure as the supervision!

原图中直接相连的2个点学习到的embedding尽可能相似。

Unsupervised Training

- “Similar” nodes have similar embeddings

$$\mathcal{L} = \sum_{z_u, z_v} \text{CE}(y_{u,v}, \text{DEC}(z_u, z_v))$$

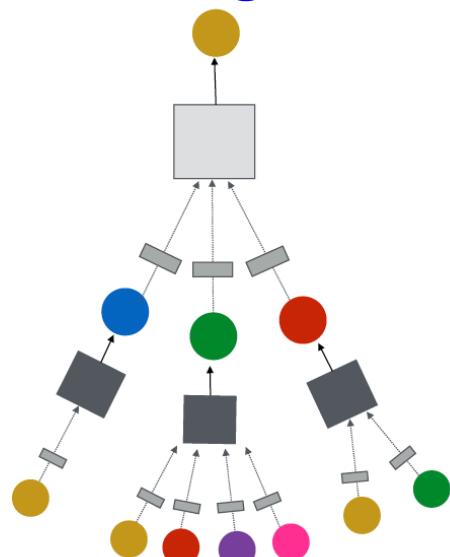
node u 和 v 的 embedding

- Where $\underline{y_{u,v} = 1}$ when node u and v are similar
- CE is the cross entropy (Slide 16)
- DEC is the decoder such as inner product (Lecture 4)
- Node similarity can be anything from Lecture 3, e.g., a loss based on:
 - Random walks (node2vec, DeepWalk, struc2vec)
 - Matrix factorization
 - Node proximity in the graph

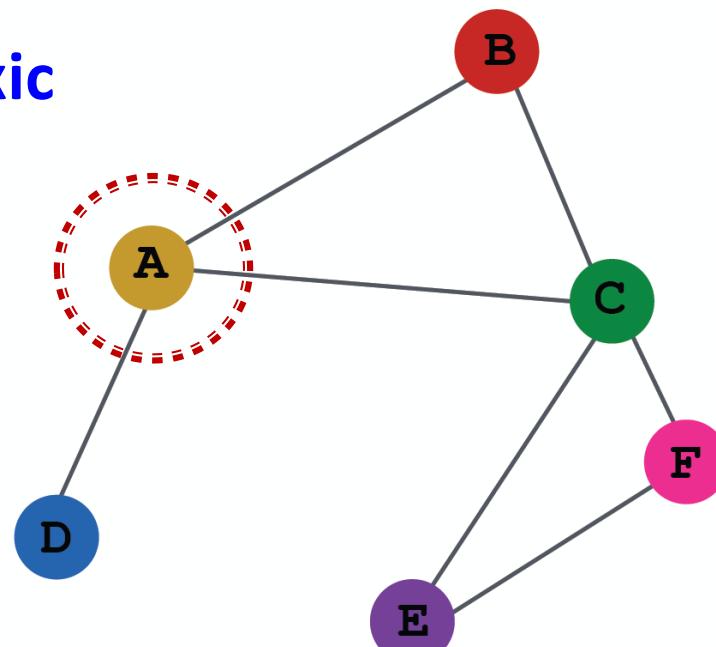
Supervised Training

Directly train the model for a supervised task
(e.g., node classification)

Safe or toxic
drug?



Safe or toxic
drug?

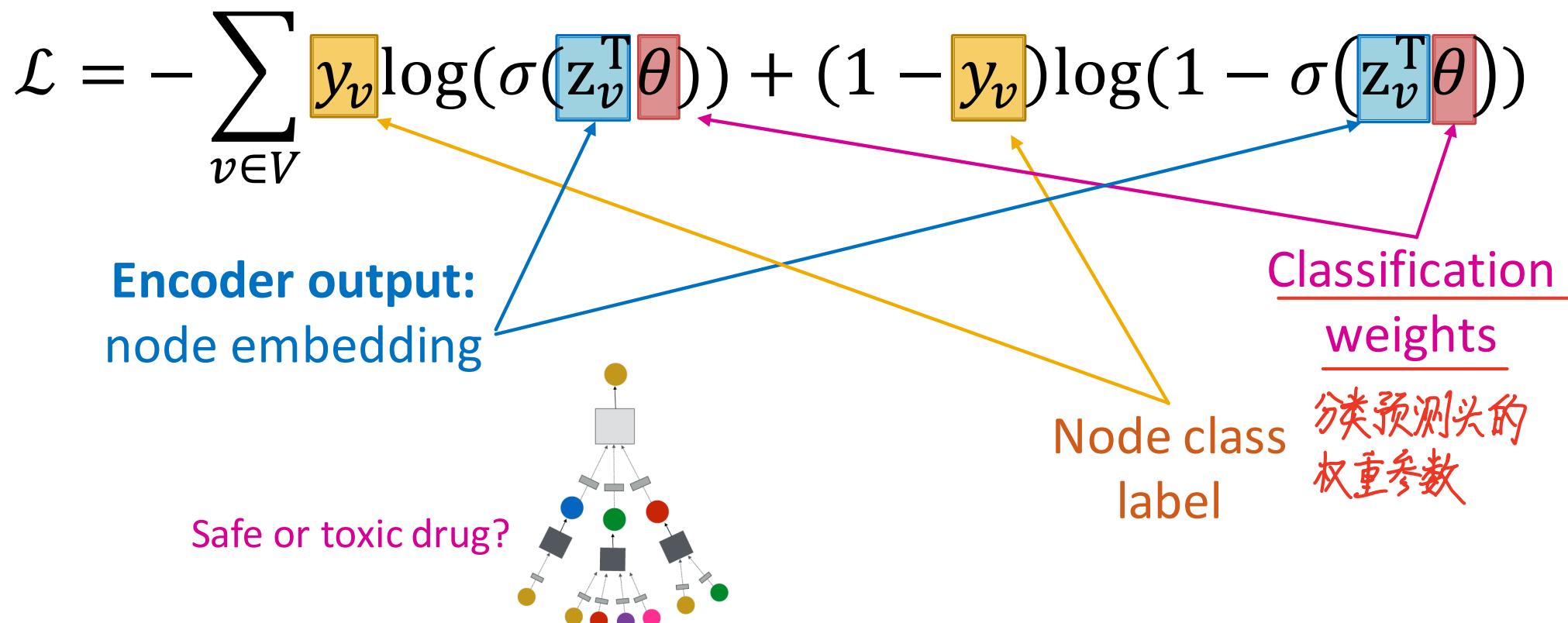


E.g., a drug-drug
interaction network

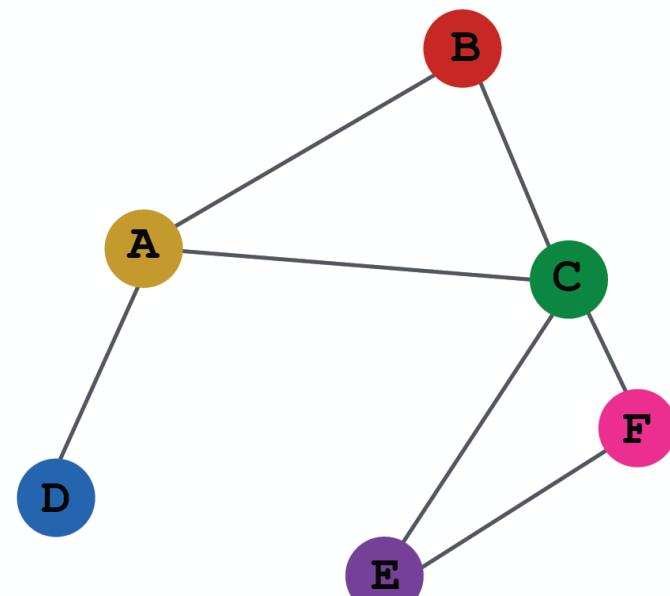
Supervised Training

Directly train the model for a supervised task
(e.g., **node classification**)

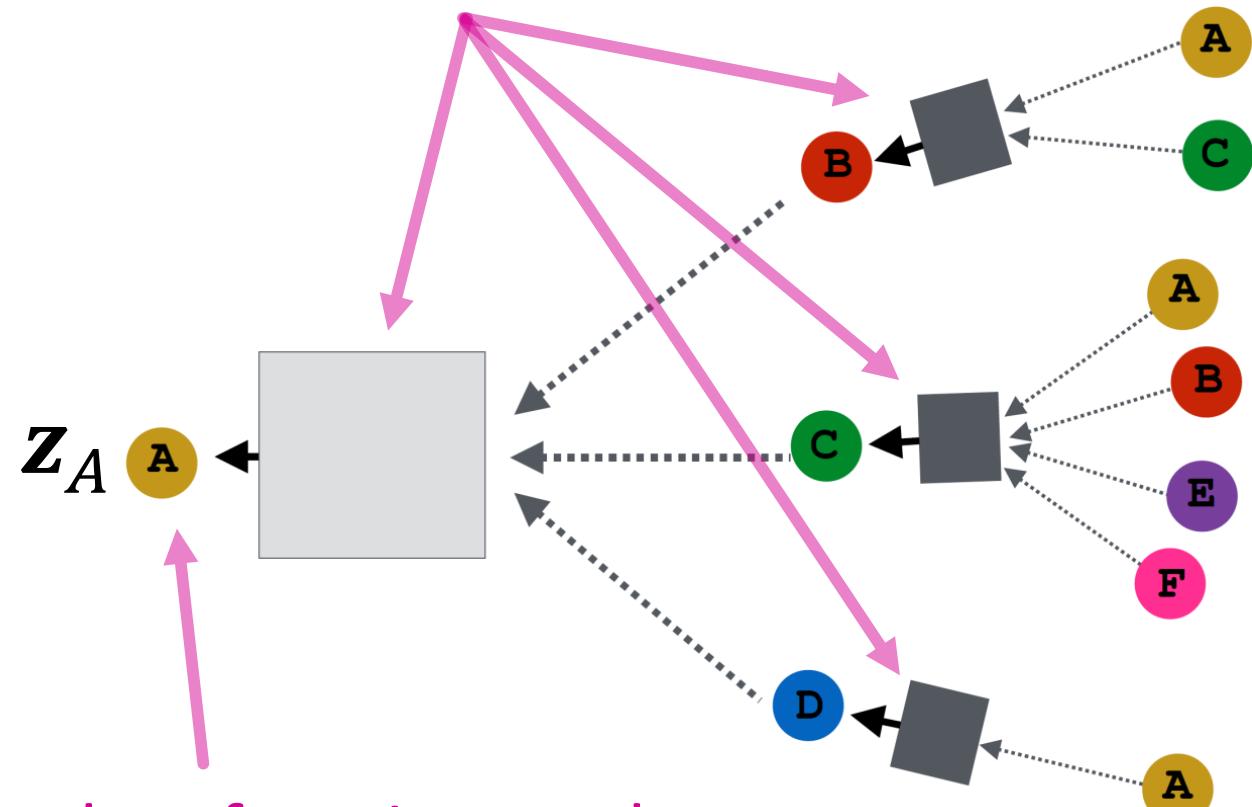
- Use cross entropy loss (Slide 16)



Model Design: Overview

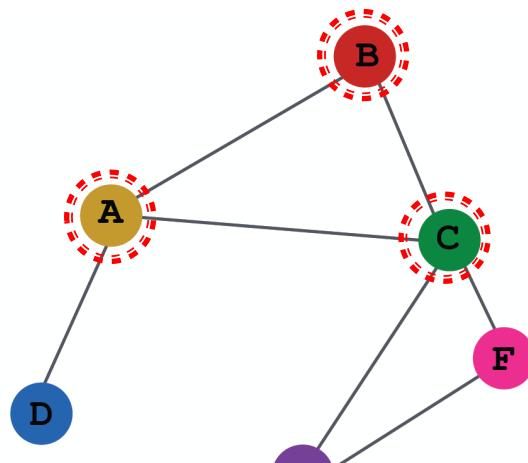


(1) Define a neighborhood aggregation function



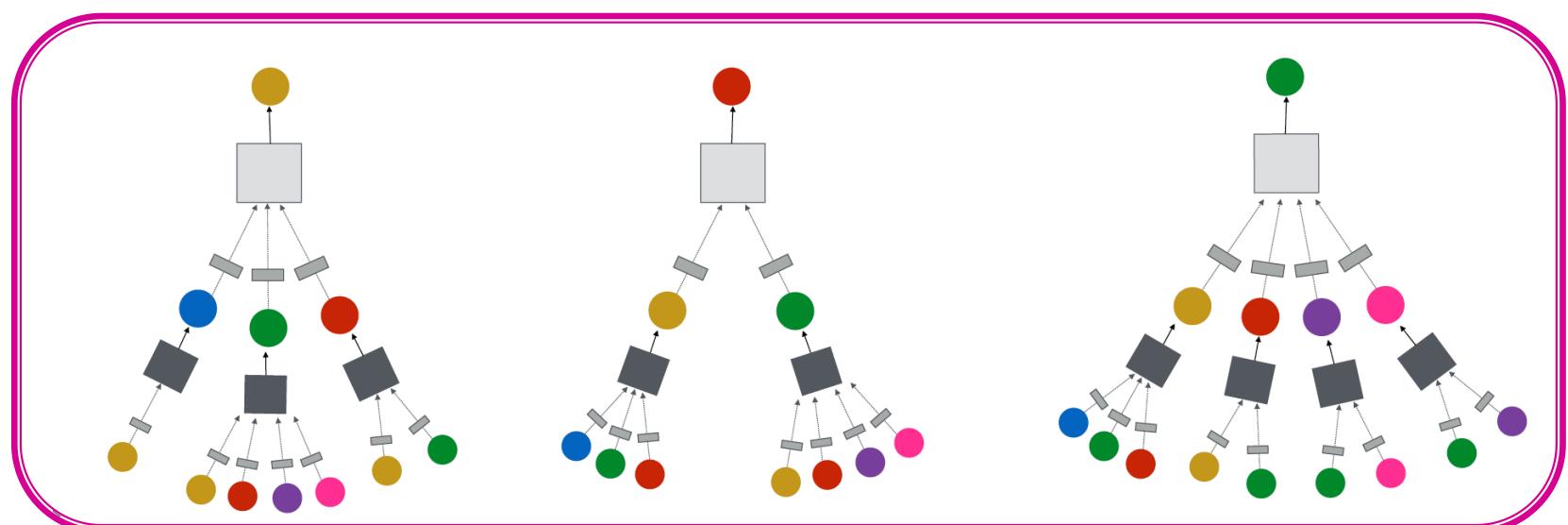
(2) Define a loss function on the embeddings

Model Design: Overview

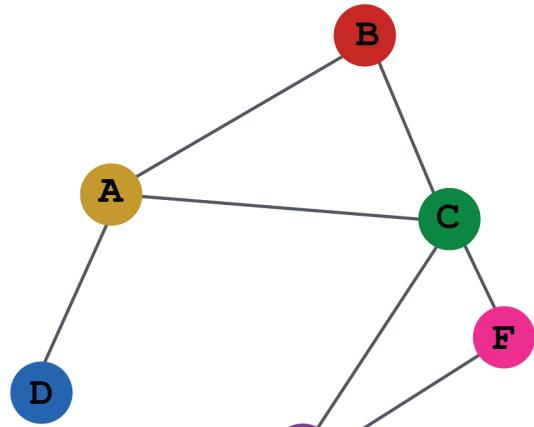


INPUT GRAPH

**(3) Train on a set of nodes, i.e.,
a batch of compute graphs**



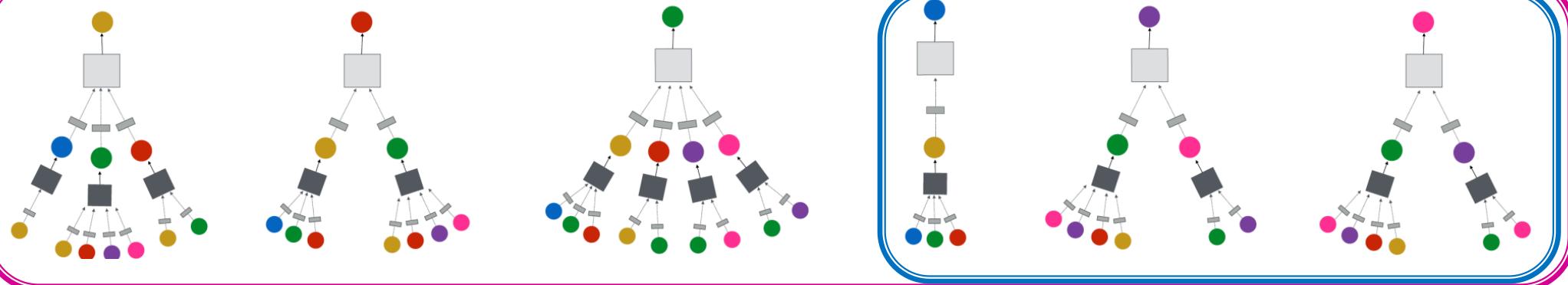
Model Design: Overview



INPUT GRAPH

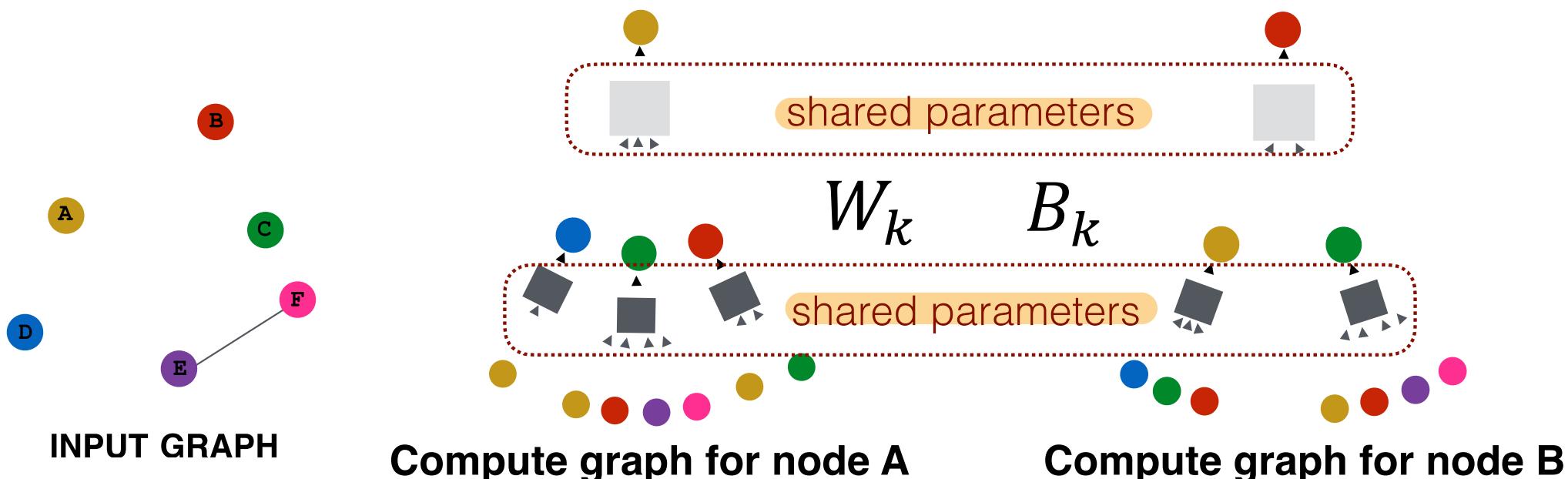
**(4) Generate embeddings
for nodes as needed**

**Even for nodes we never
trained on!**



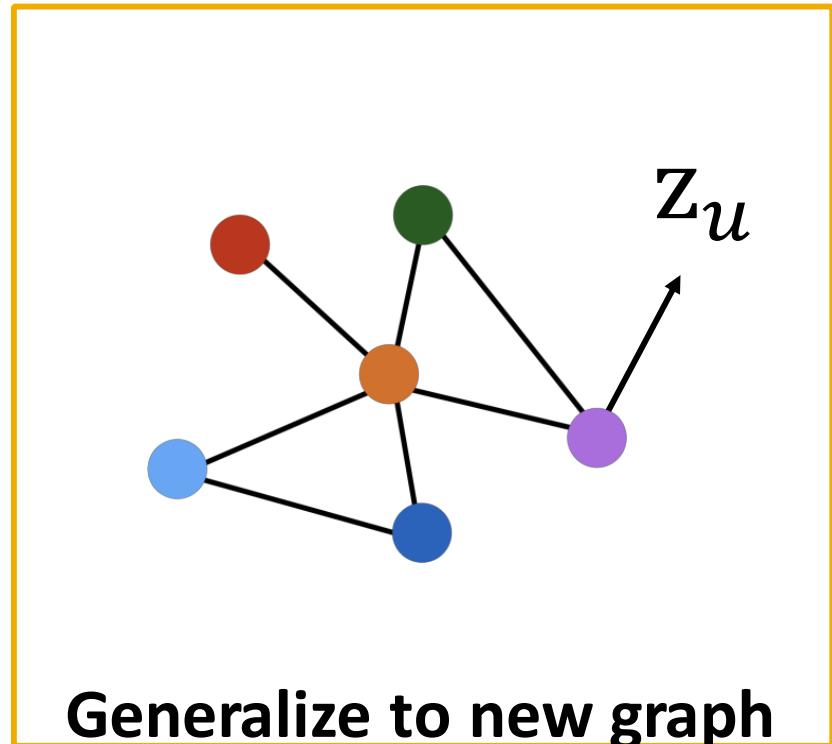
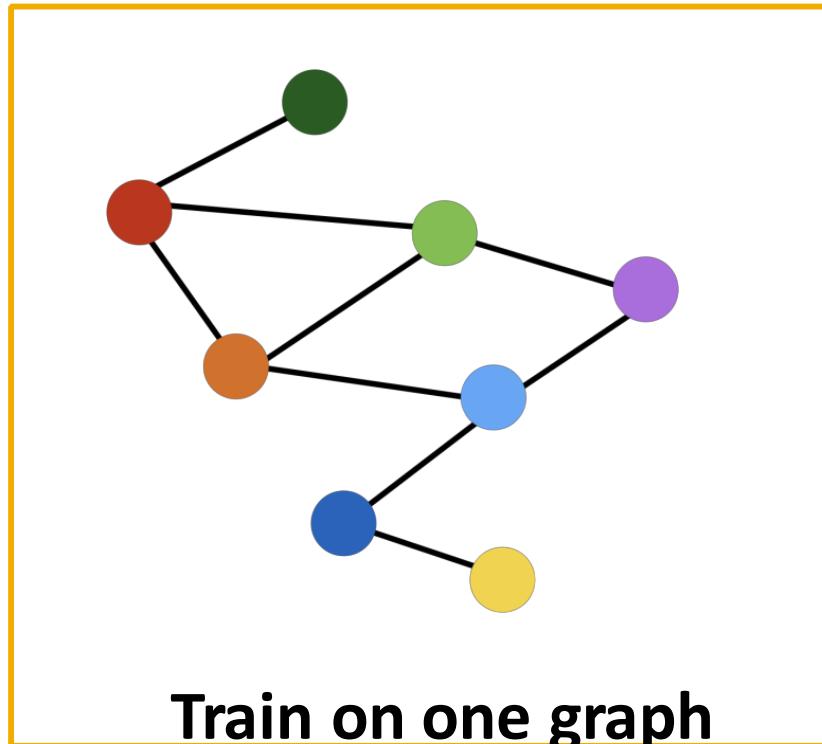
Inductive Capability

- The same aggregation parameters are shared for all nodes:
 - The number of model parameters is sublinear in $|V|$ and we can **generalize to unseen nodes!**



Inductive Capability: New Graphs

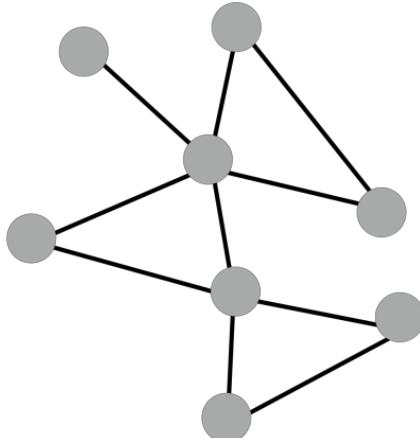
Inductive learning 归纳式学习: 泛化到新节点甚至新图 (迁移学习)



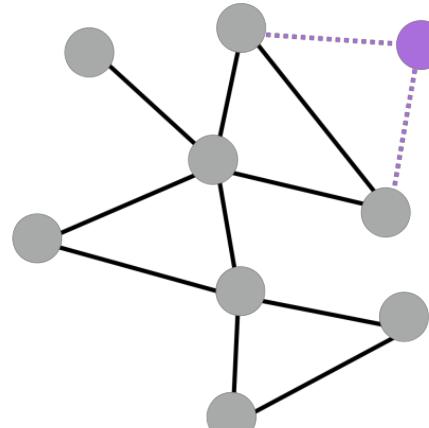
Inductive node embedding → Generalize to entirely unseen graphs

E.g., train on protein interaction graph from model organism A and generate embeddings on newly collected data about organism B

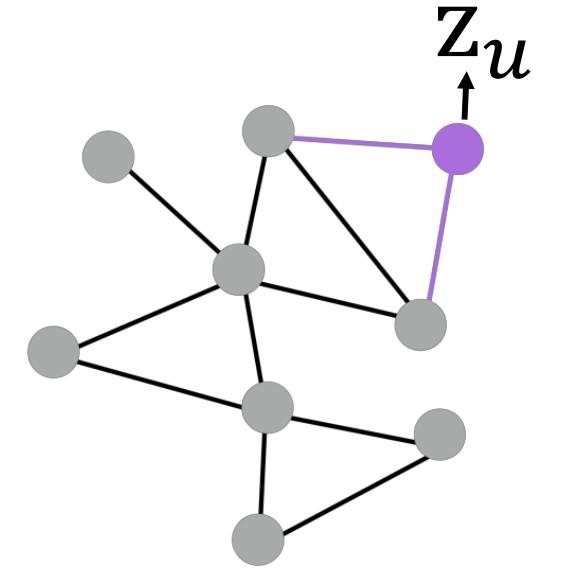
Inductive Capability: New Nodes



Train with snapshot



New node arrives



Generate embedding
for new node

- Many application settings constantly encounter previously unseen nodes:
 - E.g., Reddit, YouTube, Google Scholar
- Need to generate new embeddings "on the fly" 随用随启 冷启动

Outline of Today's Lecture

1. Basics of deep learning



2. Deep learning for graphs



3. Graph Convolutional Networks



4. GNNs subsume CNNs

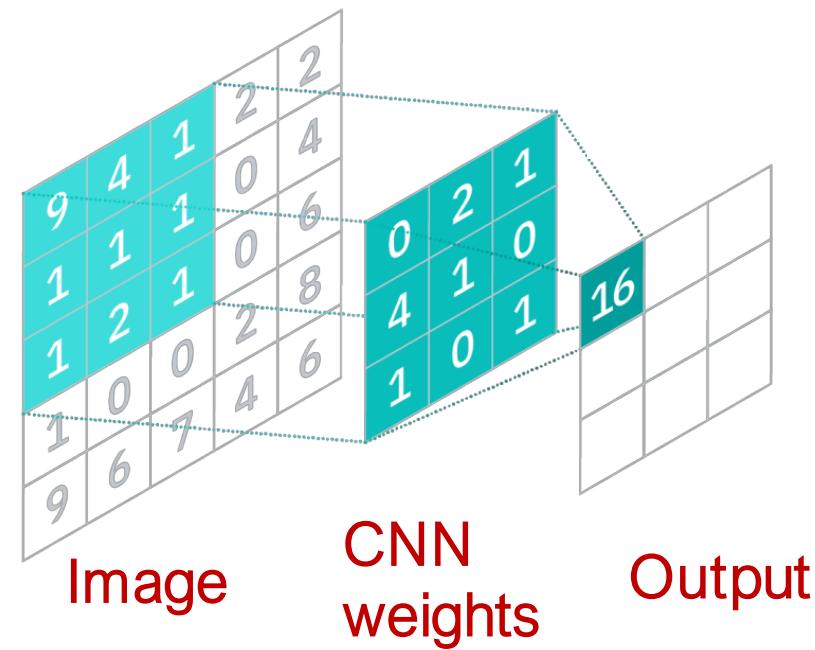
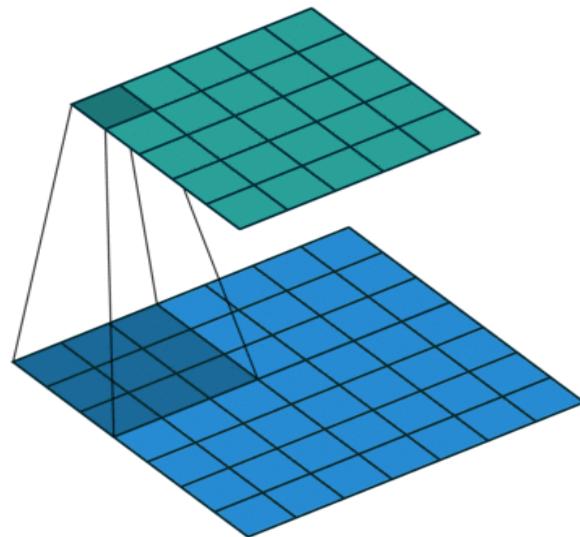


Architecture Comparison

- How do GNNs compare to prominent architectures such as Convolutional Neural Nets?

Convolutional Neural Network

Convolutional neural network (CNN) layer with 3x3 filter:

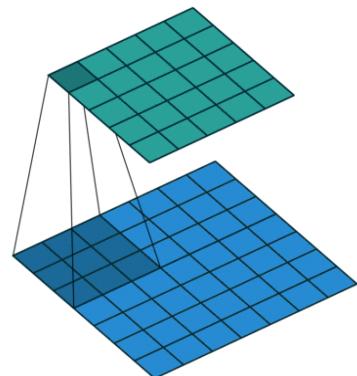


$$\text{CNN formulation: } h_v^{(l+1)} = \sigma(\sum_{u \in N(v) \cup \{v\}} W_l^u h_u^{(l)}), \quad \forall l \in \{0, \dots, L-1\}$$

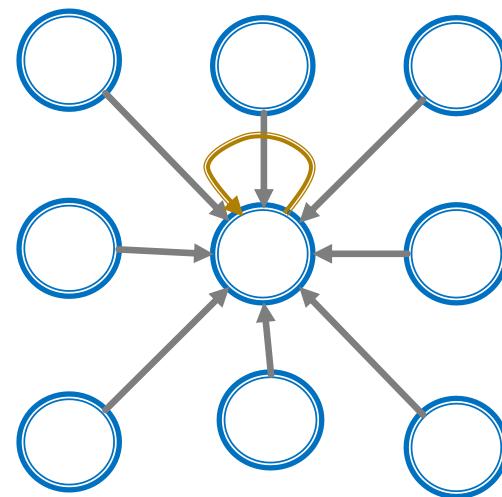
$N(v)$ represents the 8 neighbor pixels of v .

GNN vs. CNN

Convolutional neural network (CNN) layer with 3x3 filter:



Image

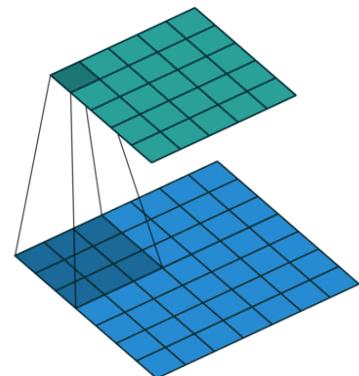


Graph

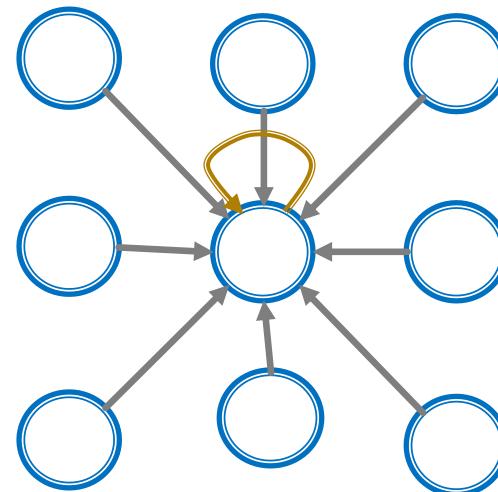
- GNN formulation: $h_v^{(l+1)} = \sigma(\underbrace{\mathbf{W}_l}_{\text{matrix}} \sum_{u \in N(v)} \frac{h_u^{(l)}}{|N(v)|} + B_l h_v^{(l)}), \forall l \in \{0, \dots, L-1\}$
- CNN formulation: (previous slide) $h_v^{(l+1)} = \sigma(\sum_{u \in N(v) \cup \{v\}} W_l^u h_u^{(l)}), \forall l \in \{0, \dots, L-1\}$
if we rewrite: $h_v^{(l+1)} = \sigma(\sum_{u \in N(v)} \underbrace{\mathbf{W}_l^u}_{\text{matrix}} h_u^{(l)} + B_l h_v^{(l)}), \forall l \in \{0, \dots, L-1\}$

GNN vs. CNN

Convolutional neural network (CNN) layer with 3x3 filter:



Image



Graph

$$\text{GNN formulation: } h_v^{(l+1)} = \sigma \left(\underbrace{\mathbf{W}_l \sum_{u \in N(v)} \frac{h_u^{(l)}}{|N(v)|}}_{\text{Learn different } W_l \text{ for different neighbors}} + B_l h_v^{(l)} \right), \forall l \in \{0, \dots, L-1\}$$

$$\text{CNN formulation: } h_v^{(l+1)} = \sigma \left(\sum_{u \in N(v)} \mathbf{W}_l^u h_u^{(l)} + B_l h_v^{(l)} \right), \forall l \in \{0, \dots, L-1\}$$

不同位置的邻域分别学习不同权重

Key difference: We can learn different W_l^u for different “neighbor” u for pixel v on the image. The reason is we can pick an order for the 9 neighbors using **relative position** to the center pixel: $\{(-1, -1), (-1, 0), (-1, 1), \dots, (1, 1)\}$

GNN vs. CNN

Convolutional neural network (CNN) layer with
3x3 filter: 
CNN 卷积核权重需学习得到
GCN 卷积核权重由 \hat{A} 预定义

- CNN can be seen as a special GNN with fixed neighbor size and ordering: 固定的邻域和固定的顺序
- The size of the filter is pre-defined for a CNN.
- The advantage of GNN is it processes arbitrary graphs with different degrees for each node.
- CNN is not permutation invariant/equivariant.
- Switching the order of pixels will leads to different outputs.

Key difference: We can learn different W_l^u for different “neighbor” u for pixel v on the image. The reason is we can pick an order for the 9 neighbors using **relative position** to the center pixel: $\{(-1, -1), (-1, 0), (-1, 1), \dots, (1, 1)\}$

Summary

■ In this lecture, we introduced

- Basics of neural networks
 - Loss, Optimization, Gradient, SGD, non-linearity, MLP
- Idea for Deep Learning for Graphs
 - Multiple layers of embedding transformation
 - At every layer, use the embedding at previous layer as the input
 - Aggregation of neighbors and self-embeddings
- Graph Convolutional Network
 - Mean aggregation; can be expressed in matrix form
 - GNN is a general architecture
 - CNN can be viewed as a special GNN

①

计算图

{

②