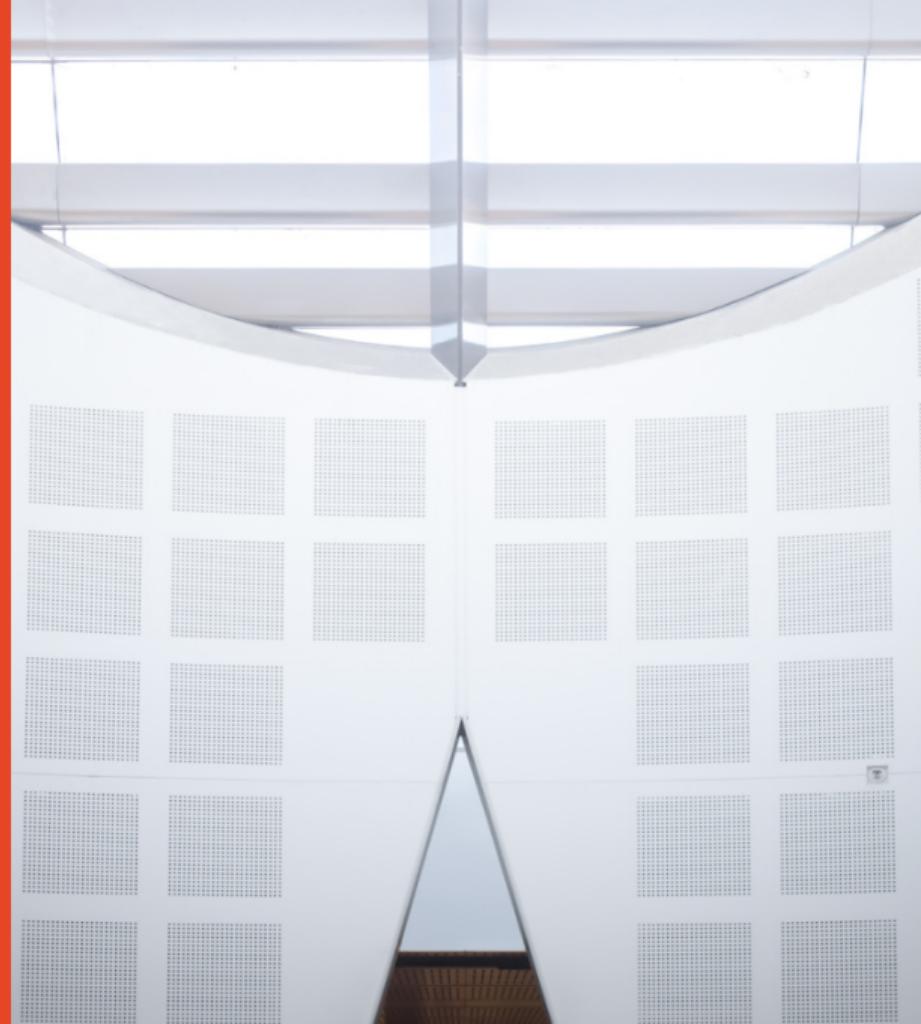


COMP5313/COMP4313 - Large Scale Networks

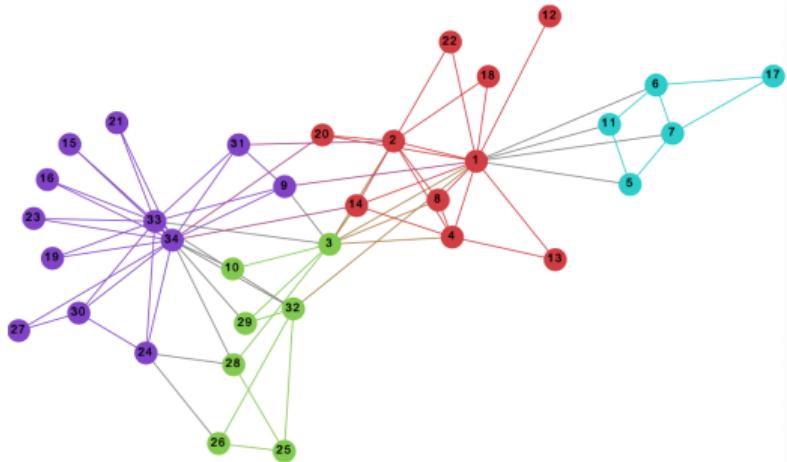
Week 7: Machine Learning on
Graphs (I)

Lijun Chang

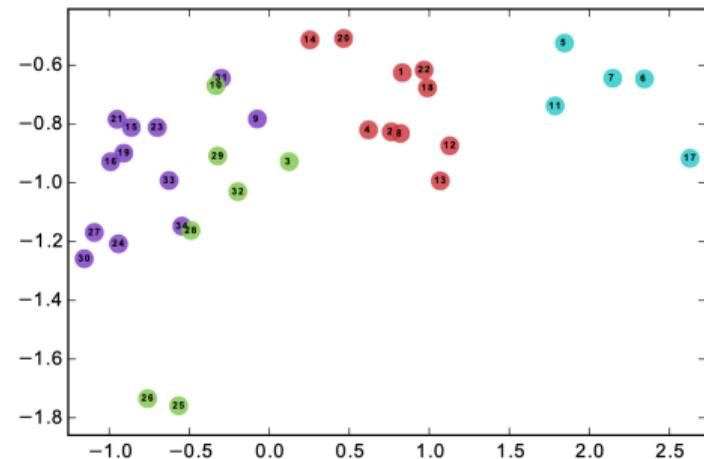
April 10, 2025



Graph Representation Learning (Node Embedding)



(a) Zachary's karate club network



(b) Embedding in \mathbb{R}^2 (by DeepWalk)

- Node color represents a modularity-based clustering of the graph

Bryan Perozzi, Rami Al-Rfou, Steven Skiena: DeepWalk: online learning of social representations. KDD 2014: 701-710

Outline

Background on Machine/Deep Learning (optional)

Machine Learning on Graphs

Graph Neural Networks (GNNs)

What is Machine Learning?

- ▶ Machine learning (ML) is the study of computer algorithms that can **improve automatically** through **data and experience**.
- ▶ Algorithms and processes that '**learn**' in the sense of being able to generalize **past data and experiences** in order to predict future outcomes.
- ▶ A computer program is said to learn form **experience E** with respect to some **task T** and some **performance measure P** , if its performance on T , as measured by P , improves with experience E .

Supervised Learning via Function Approximation

- ▶ Experience E is a **dataset**, containing measures of objects w.r.t a set of features
 - Features are partitioned into input features $\mathbf{x} \in \mathbb{R}^p$ and an output feature $y \in \mathbb{R}$
 - ▶ $\mathbf{x} = (\text{type of a property, area, no. of bedrooms})$, $y = \text{price}$.
 - ▶ $\mathbf{x} = (\text{pixels of an image at each position})$, $y = \text{cat or not.}$
 - Each example in the dataset is an input-output pair (\mathbf{x}_i, y_i)
 - ▶ the measures of one particular object
 - Training data: $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$
- ▶ Task T is **mapping** inputs \mathbf{x} to output y
 - Machine learning is to learn a function $f : \mathbf{x} \rightarrow y$
- ▶ Performance P is how well f **fits** the training data \mathcal{D} , and how well f **generalizes** to unseen (test) data (i.e., given arbitrary \mathbf{x} , predict y)
 - To avoid overfitting (to training data), there is an **assumption/restriction** on the function f (e.g., the degree of a polynomial function)

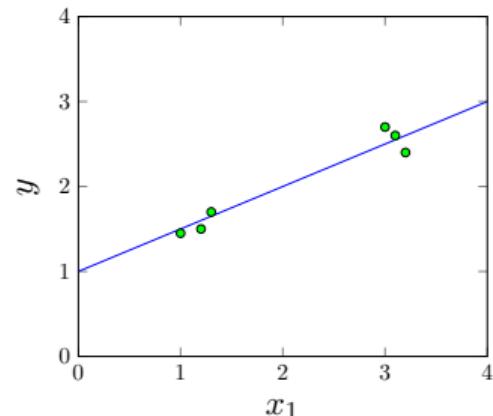
Linear Regression (for Quantitative Output)

- ▶ Problem setting:
 - Features: $x = (\text{type of a property, area, no. of bedrooms})$, $y = \text{price}$.
 - Learn a linear function $f(x) = b + \sum_{i=1}^p w_i x_i = b + \mathbf{x}^\top \mathbf{w}$ to approximate y
 - ▶ The learning algorithm is to estimate b, w_1, \dots, w_p .
 - ▶ We assume column vectors
 - Performance is measured by the mean square error on the training dataset
- MSE(\mathbf{w}) = $\frac{1}{N} \sum_{i=1}^N (y_i - f(\mathbf{x}_i))^2$.
 - ▶ N : the number of training examples
- ▶ Matrix form (by ignoring b):

$$\text{MSE}(\mathbf{w}) = \frac{1}{N} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2$$

$$-\quad \mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix} \text{ and } \mathbf{X} = \begin{bmatrix} \mathbf{x}_1^\top \\ \vdots \\ \mathbf{x}_N^\top \end{bmatrix}$$

- ▶ Assume the features are independent, then
- $$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \text{MSE}(\mathbf{w}) = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}.$$

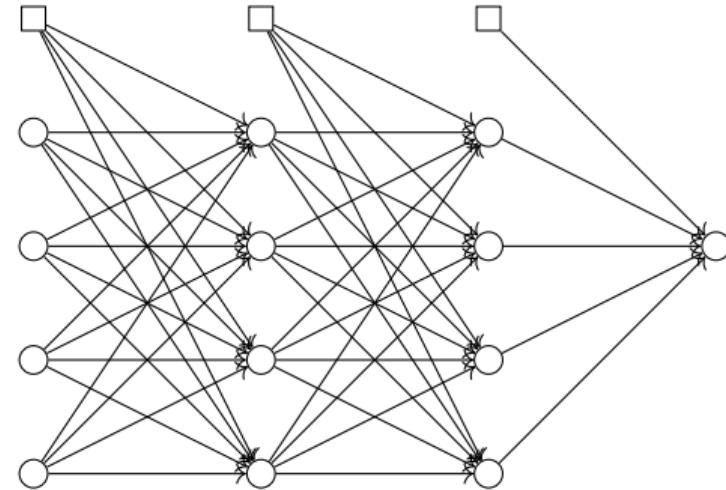


Multi-Layer Perceptron (MLP)

Now, let's consider one input example

(\mathbf{x}, y)

- ▶ Linear regression: $f_1(\mathbf{x}) = \mathbf{w}_1^\top \mathbf{x} + b_1$
- ▶ What if we have multiple outputs?
 - $f_i(\mathbf{x}) = \mathbf{w}_i^\top \mathbf{x} + b_i$
 - $\mathbf{f}(\mathbf{x}) = \mathbf{W}\mathbf{x} + \mathbf{b}$ where
 $\mathbf{W} = [\mathbf{w}_1 \quad \mathbf{w}_2 \quad \dots]^\top$
- ▶ What if we stack multiple layers?
 - $\mathbf{f}^{(k)}(\mathbf{x}^{(k-1)}) = \mathbf{W}^{(k)}\mathbf{x}^{(k-1)} + \mathbf{b}^{(k)}$
 - But this is still a linear function!
- ▶ **Remedy:** add a nonlinear activation function to each hidden unit.
 - E.g., $f_i(\mathbf{x}) = \text{ReLU}(\mathbf{w}_i^\top \mathbf{x} + b_i)$



Input Sample

Hidden Layer

Hidden Layer

Output

Every edge carries a scalar weight

Multi-Layer Perceptron (MLP)

- ▶ MLP is also known as *deep feedforward network* or *feedforward neural network* or *fully-connected neural network*.
- ▶ The functions learned have the form $f(\mathbf{x}) = L^{(K)}(\sigma(\dots(L^{(2)}(\sigma(L^{(1)}(\mathbf{x}))))))$.
 - A composition of affine functions $L^{(k)}(\mathbf{x}^{(k-1)}) = \mathbf{W}^{(k)}\mathbf{x}^{(k-1)} + \mathbf{b}^{(k)}$ with nonlinear activation functions σ (e.g., ReLU)
 - ▶ The activation function acts on each component of a vector or matrix.
 - The hidden layers add depth to the network.
 - ▶ It is that depth which has allowed the composite function f to be so successful in deep learnings
- ▶ The **universal approximation theorem** says that with one hidden layer, MLP can approximate any well-behaved function $f(\cdot)$ within an arbitrarily small error ε .¹

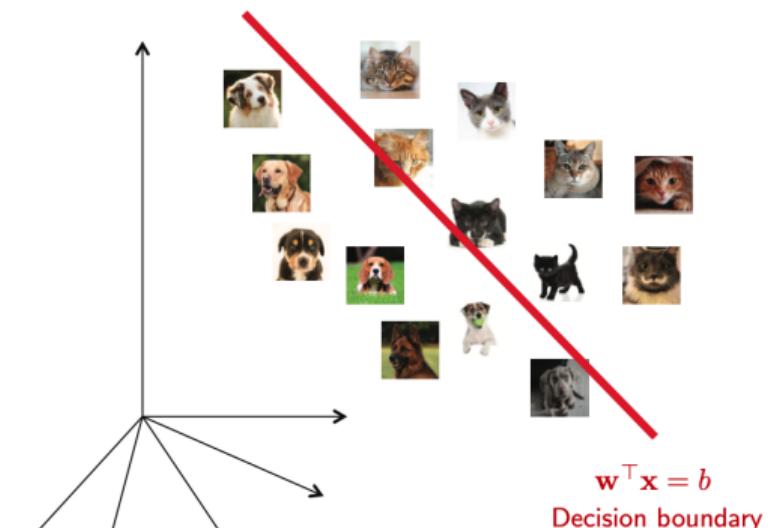
¹https://en.wikipedia.org/wiki/Universal_approximation_theorem

Classification (for Qualitative Output)

► Problem setting:

- Features: x = (pixels of an image at each position), y = cat or dog.
 - $x \in \mathbb{R}^d$ (e.g., for 1024×1024 images $d \approx 10^6$).
- Learn a classification function $f : \mathbb{R}^d \rightarrow \{1, \dots, C\}$ in C classes
- Performance is measured by maximum likelihood

► Traditional methods: logistic regression, linear discriminant analysis, separating hyperplane.



Classification in MLP (for Qualitative Output)

- ▶ There is one output unit for each class

- $z = \mathbf{W}^{(K)} \mathbf{x}^{(K-1)} + \mathbf{b}^{(K)}$

- ▶ Apply the **softmax** function to the output layer

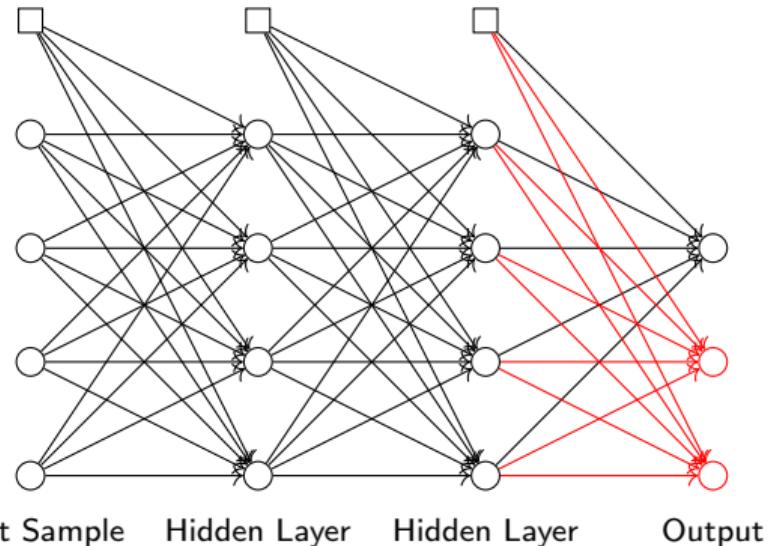
- $\text{softmax}(\mathbf{z})_i = \frac{\exp(z_i)}{\sum_j \exp(z_j)}$

- This is regarded as the probability $P(y = i | \mathbf{x}, \theta)$ of seeing class i for input \mathbf{x} and parameters $\theta = \{\mathbf{b}^{(1)}, \mathbf{W}^{(1)}, \dots, \mathbf{b}^{(K)}, \mathbf{W}^{(K)}\}$

- ▶ **Maximize** the likelihood

$$\prod_{i=1}^N P(y = y_i | \mathbf{x}_i, \theta)$$

- Equivalent to **minimize** the negative log-likelihood **loss** function $\ell(\theta) = -\frac{1}{N} \sum_{i=1}^N \log P(y = y_i | \mathbf{x}_i, \theta)$

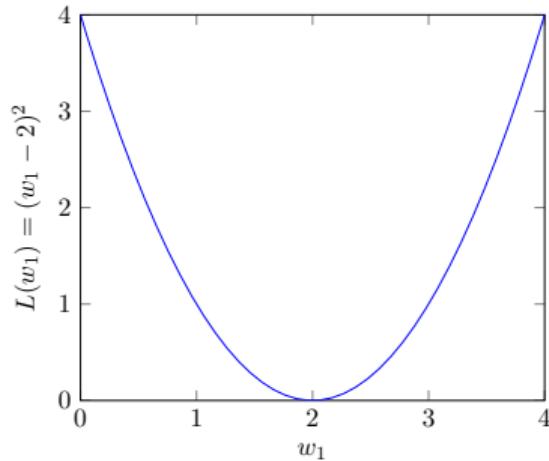


How to Find the Best Parameters (i.e., How to Train the Model)?

- ▶ For linear regression, we aim to minimize $\text{MSE}(\mathbf{w}) = \frac{1}{N} \sum_{i=1}^N (y_i - \mathbf{w}^\top \mathbf{x}_i)^2$.
 - Luckily, we got a closed-form solution $\mathbf{w} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$.
- ▶ In general, we aim to minimize the loss function $\ell(\boldsymbol{\theta}) = -\frac{1}{N} \sum_{i=1}^N \log P(y_i | \mathbf{x}_i, \boldsymbol{\theta})$
 - Unlikely to have closed-form solution

▶ Gradient-based optimization

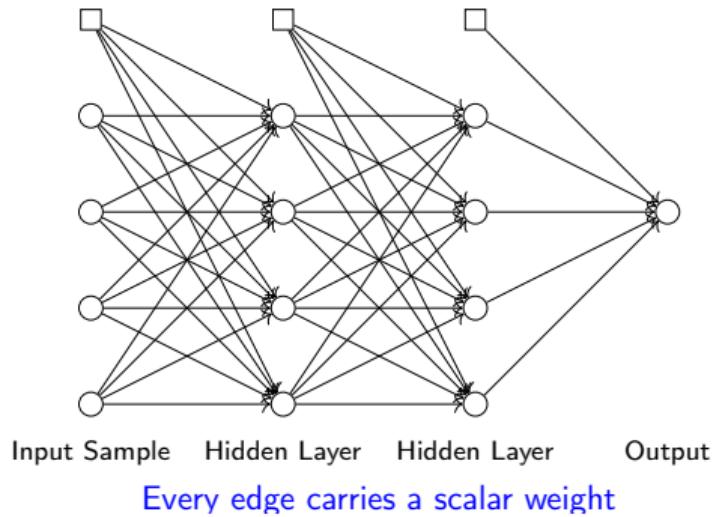
- $\nabla_{\boldsymbol{\theta}} \ell(\boldsymbol{\theta}) = \begin{bmatrix} \frac{\partial \ell(\boldsymbol{\theta})}{\partial \theta_1} & \frac{\partial \ell(\boldsymbol{\theta})}{\partial \theta_2} & \dots \end{bmatrix}^\top$
points the uphill direction
- $\ell(\boldsymbol{\theta})$ decreases the fastest along the vector direction $-\nabla_{\boldsymbol{\theta}} \ell(\boldsymbol{\theta})$
- $\boldsymbol{\theta}^{(t+1)} \leftarrow \boldsymbol{\theta}^{(t)} - \rho \nabla_{\boldsymbol{\theta}} \ell(\boldsymbol{\theta}^{(t)})$
 - ▶ ρ is the learning rate.



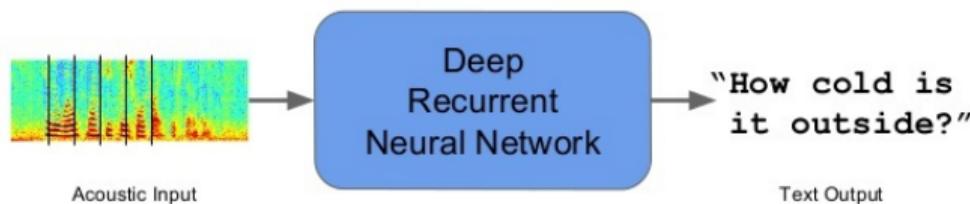
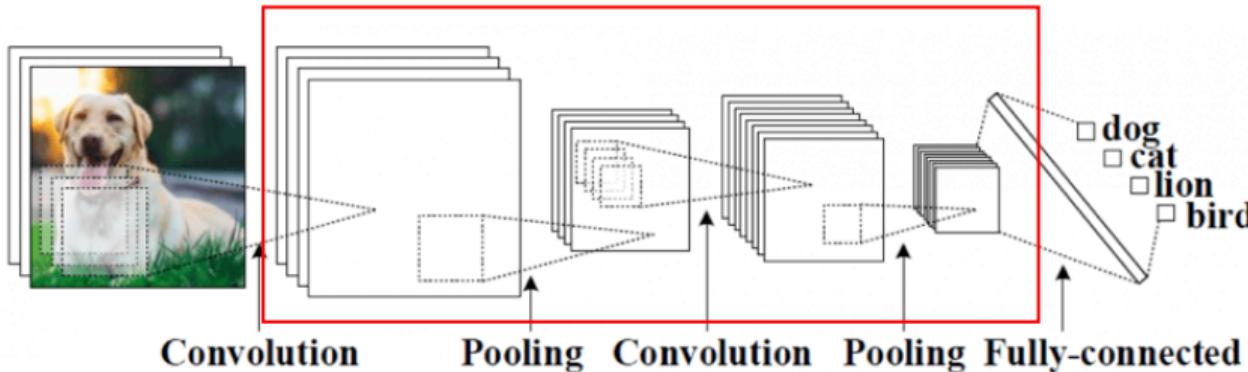
- ▶ Stochastic Gradient Descent (SGD) for computational efficiency
 - Instead of using all training data to compute $\nabla_{\boldsymbol{\theta}} \ell(\boldsymbol{\theta}^{(t)})$, use a small random subset (called mini-batch)

Learning Procedure

- ▶ One-step learning procedure for a mini-batch of the training data:
 - Perform **forward propagation** through the deep neural network for the subset of training data.
 - **Compute the loss** at the output layer for the subset of training data.
 - **Backpropagate** to compute the gradient and to **update** the weight parameters $\mathbf{W}^{(k)}$ in each layer.



Deep Neural Network Architectures for Specific Applications



Outline

Background on Machine/Deep Learning (optional)

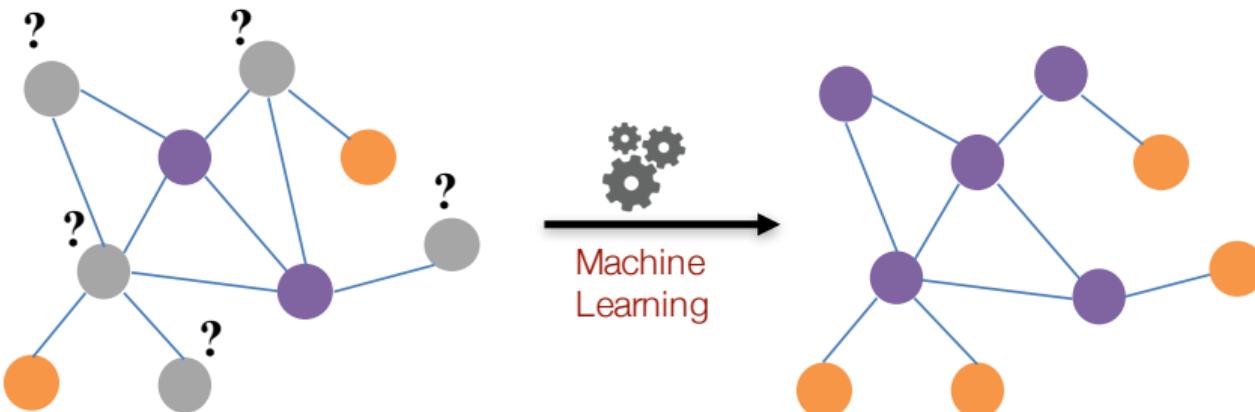
Machine Learning on Graphs

Graph Neural Networks (GNNs)

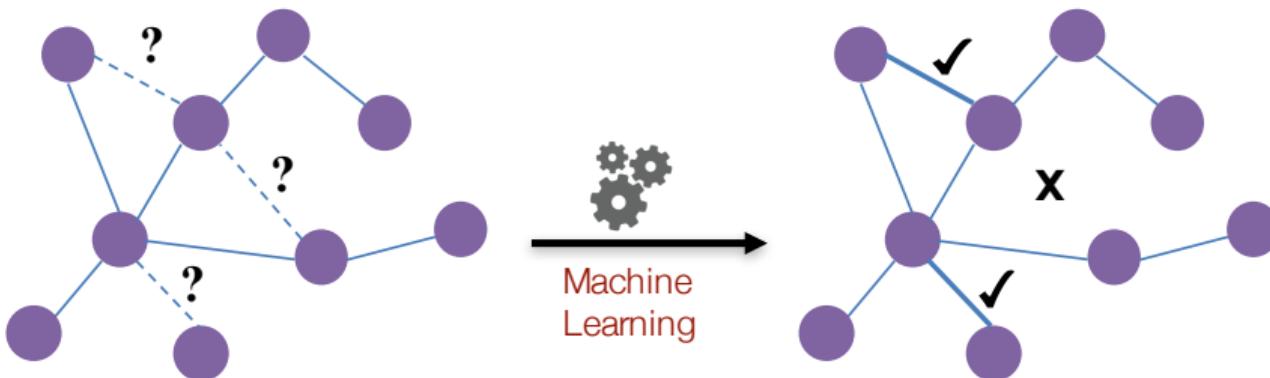
Typical Machine Learning Tasks on Graphs

- ▶ Node classification
 - Predict a type of a given node (e.g., topic of documents in a citation network)
- ▶ Link prediction
 - Predict whether two nodes are (or will be) linked
- ▶ Community detection
 - Identify densely linked clusters of nodes
- ▶ Network similarity
 - How similar are two (sub)networks
- ▶ Graph classification
 - Predict a type of a given graph

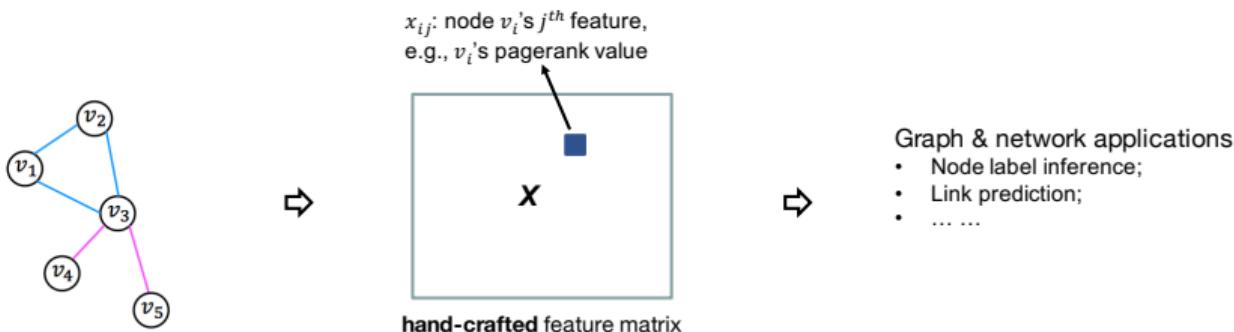
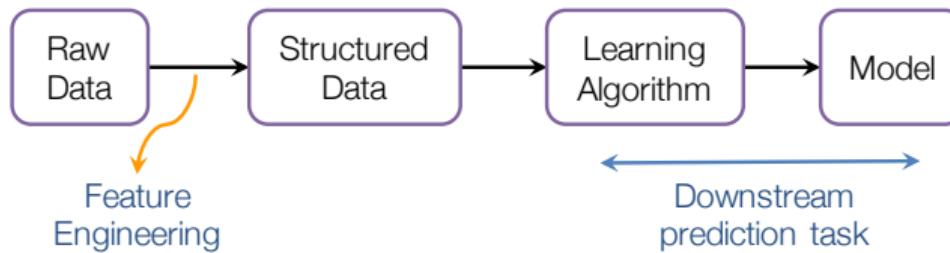
Node Classification



Link Prediction

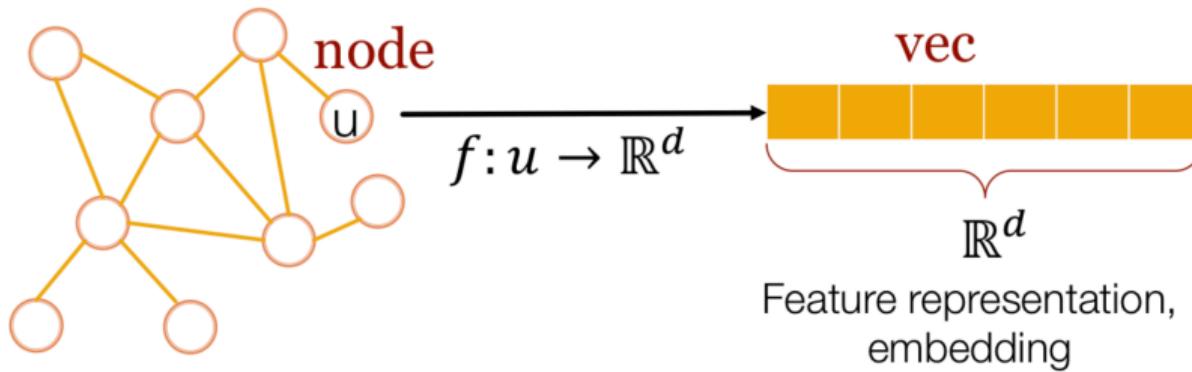


Traditional ML Life Cycle

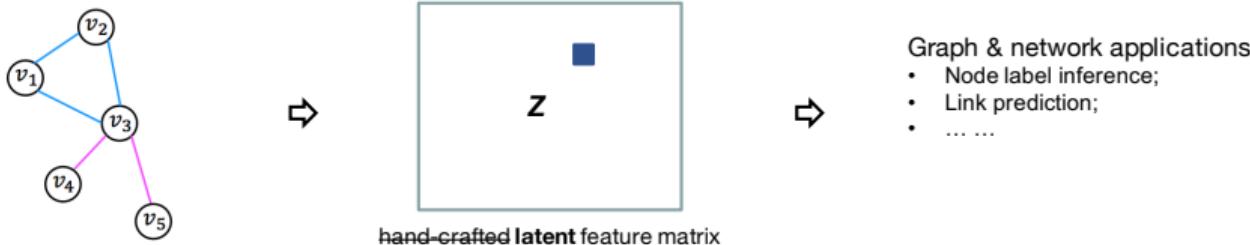
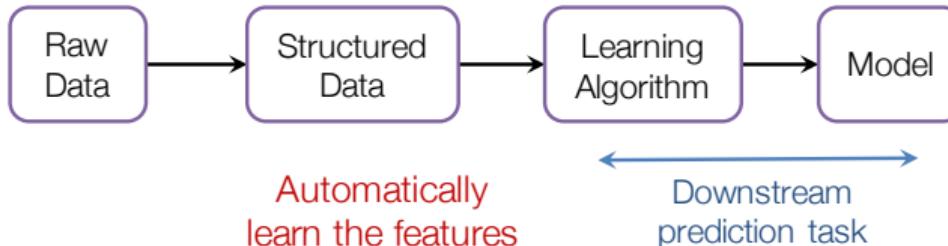


Feature Engineering/Learning

- ▶ Task dependent/independent feature engineering/learning for machine learning on graphs



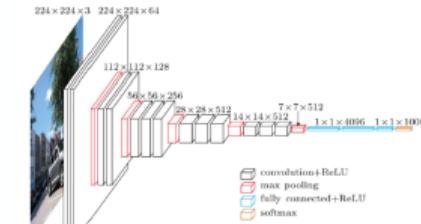
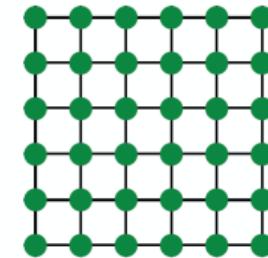
Deep Learning-Driven ML Life Cycle



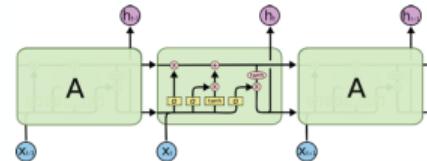
Leskovec, J., Hamilton, W. L., Ying, R., and Sosic, R.: Representation learning on networks. WWW 2018 Tutorial.

Challenges

- ▶ Modern deep learning toolbox is designed for simple sequences or grids
 - Convolutional Neural Networks (CNNs) for fixed-size images/grids ...

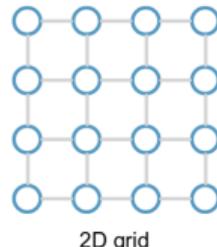
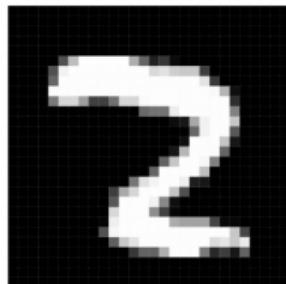


- Recurrent Neural Networks (RNNs) or word2vec for text/sequences ...

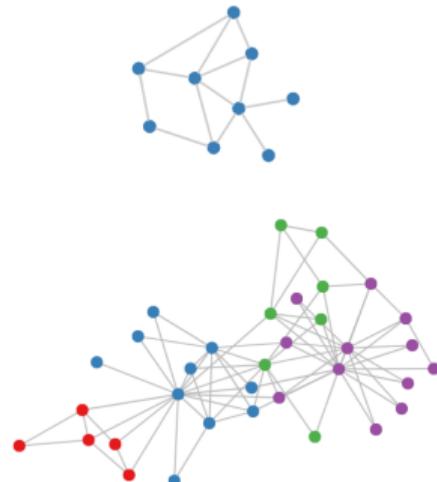


Challenges

- ▶ Graphs are known for the flexibility of the data structure
 - Arbitrary size, complex topographical structure (i.e., no spatial locality like grids), different number of neighbors, no fixed node ordering or reference point



Vs.



Outline

Background on Machine/Deep Learning (optional)

Machine Learning on Graphs

Graph Neural Networks (GNNs)

Setup for Node Classification

Assume we have an **undirected** graph G :

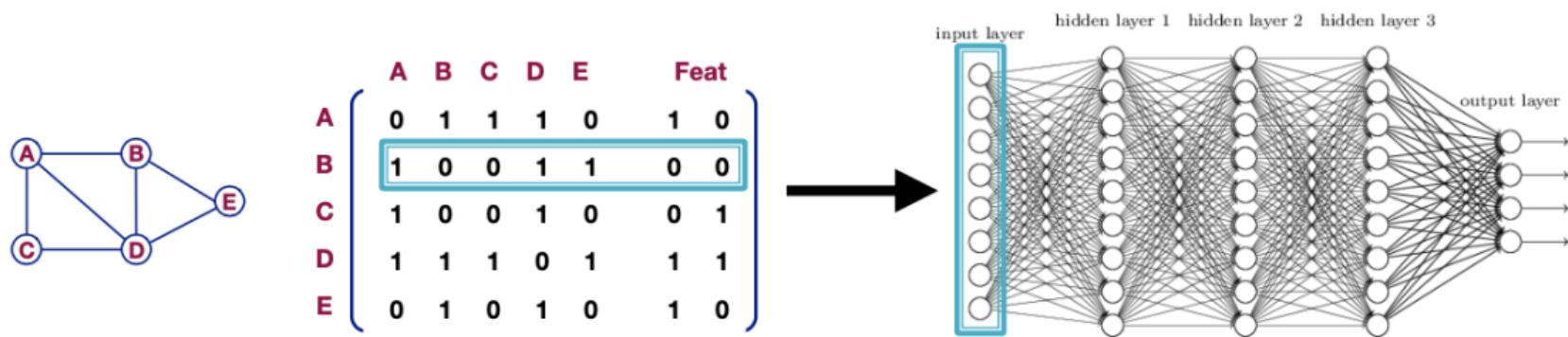
- ▶ V is the **node set** ($|V| = N$)
- ▶ $A \in \mathbb{R}^{N \times N}$ is the **adjacency matrix** (assume binary)
- ▶ $X \in \mathbb{R}^{N \times q}$ is a **matrix of input node features**
 - Categorical attributes, text, image data
 - ▶ E.g., profile information in a social network
 - Node degrees, clustering coefficients, etc.
 - Indicator vectors (i.e., one-hot encoding of each node)
- ▶ For each labeled node $v \in V_{\text{train}}$, its label is denoted by $y_v \in \mathbb{R}^C$, a one-hot encoding vector of the label
- ▶ The task is to predict the labels for nodes in $V \setminus V_{\text{train}}$

A First Naïve Approach

- ▶ Ignore the graph structure information (i.e., ignore \mathbf{A}), and use only \mathbf{X} and y_v for training and prediction
- ▶ This may perform well on some datasets where \mathbf{X} contains a lot of information
- ▶ In general, if we incorporate graph structure information, we can do better

A Second Naïve Approach

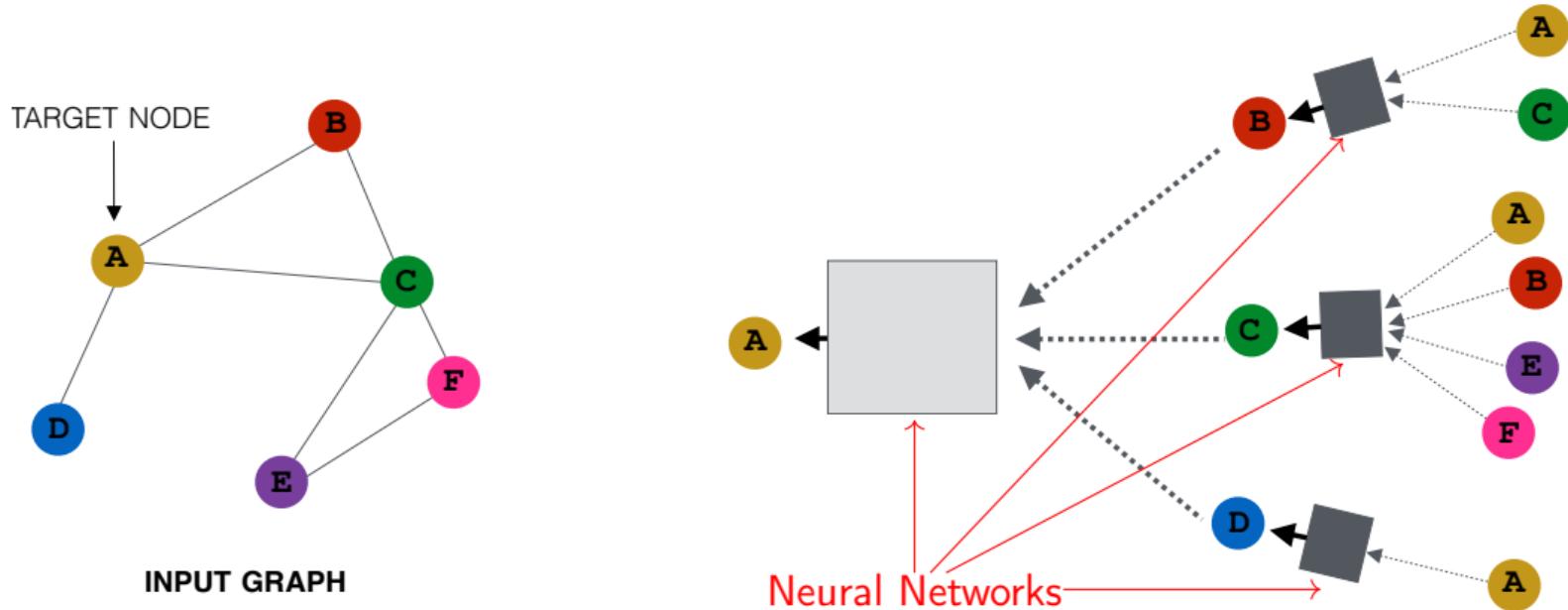
- ▶ Join adjacency matrix and feature matrix, and feed them into a deep neural network



- ▶ Issues with this idea:
 - $O(N)$ parameters
 - Not applicable to graphs of different sizes
 - Not invariant to node ordering

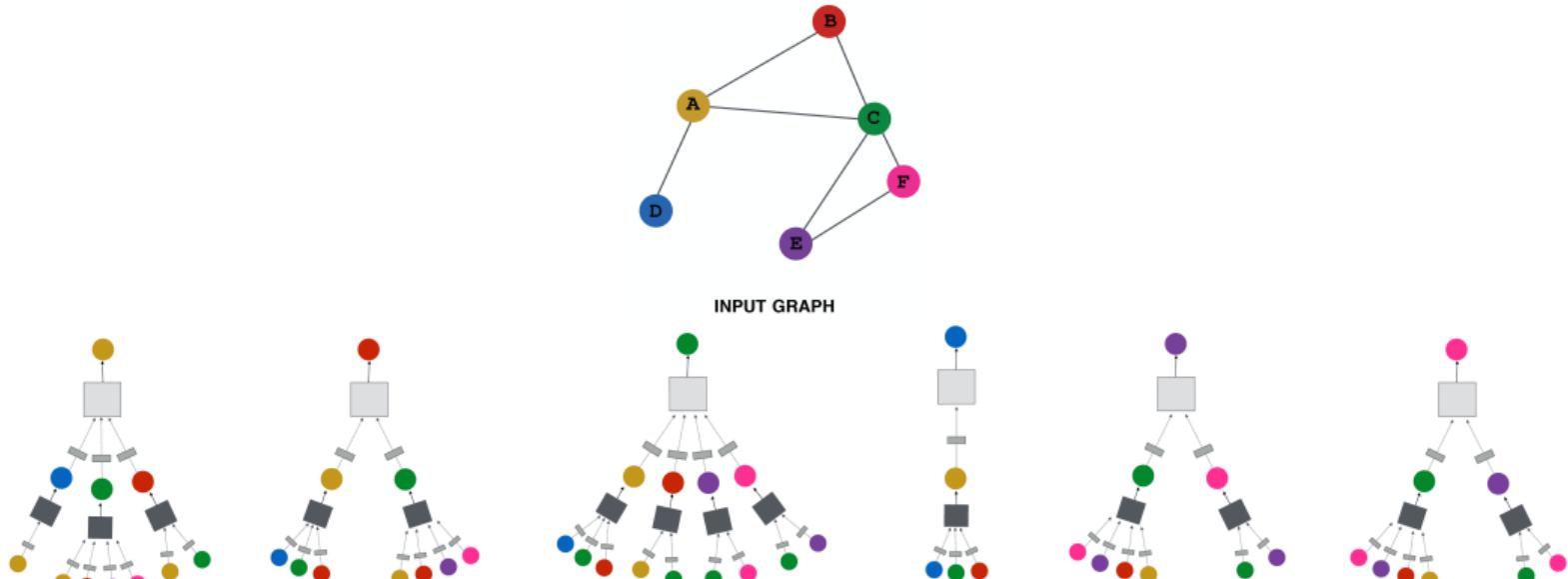
Neighborhood Aggregation

- ▶ Graph neural network approach: combine graph structure and node features via **neighborhood aggregation**, which propagates features based on the graph topology.
- ▶ **Intuition:** nodes aggregate information from their neighbors using neural networks



Neighborhood Aggregation

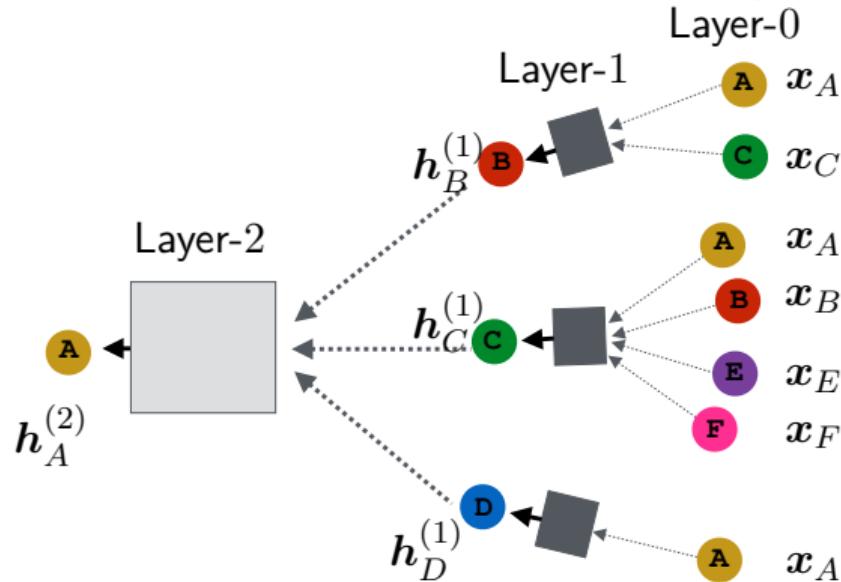
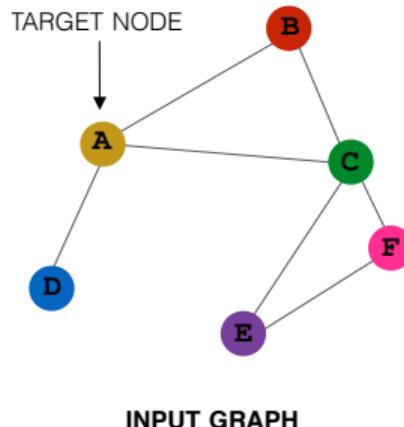
- ▶ **Intuition:** network neighborhood defines a computation graph



Every node defines a computation graph!

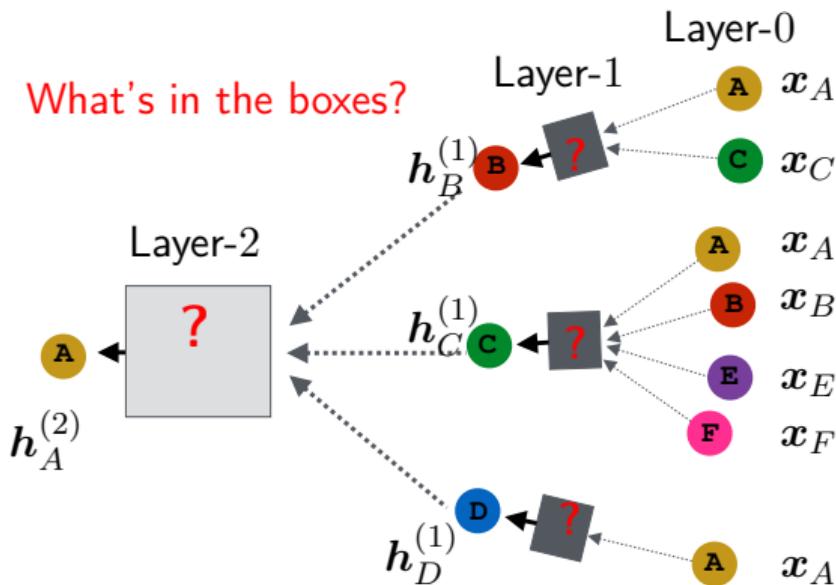
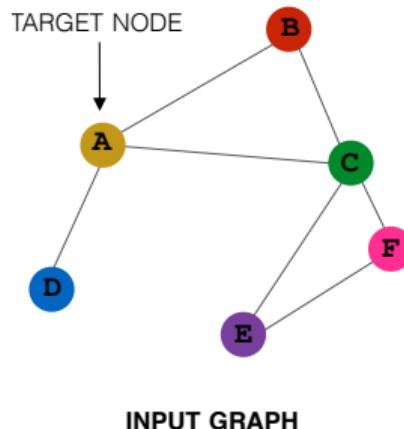
Neighborhood Aggregation

- ▶ Model can be of arbitrary depth
 - Nodes have embeddings at each layer
 - Layer-0 embedding of node u is its input feature \mathbf{x}_u
 - Layer- k embedding gets information from nodes that are k hops away



Neighborhood Aggregation

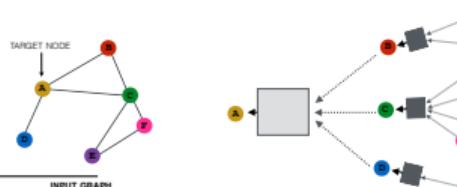
- ▶ Key distinctions of different graph neural networks models are in how they aggregate information across the layers



Graph Convolutional Networks (GCN)²

Weighted average of neighbor messages and then applying a neural network

- ▶ $\mathbf{h}_u^{(0)} = \mathbf{x}_u, \forall u \in V$
- ▶ For $k = 1, \dots, K$
 - $\mathbf{h}_u^{(k)} = \sigma \left(\mathbf{W}^{(k)} \left(\sum_{v \in N(u) \cup \{u\}} \frac{\mathbf{h}_v^{(k-1)}}{\sqrt{(d(v)+1)(d(u)+1)}} \right) \right), \quad \forall u \in V$
 - ▶ $N(u)$ is the neighbors of u , $d(u) = |N(u)|$ is the degree of u
 - ▶ $\mathbf{W}^{(k)}$ is the trainable weight parameters
 - ▶ $\sigma(\cdot)$ is element-wise activation function (e.g., $\text{ReLU}(x) = \max(0, x)$)
- ▶ $\mathbf{z}_u = \mathbf{h}_u^{(K)}$ is the learned representation of u , $\forall u \in V$.
 - These representations can be fed into any loss function to train the weight parameters. E.g., $\{(\mathbf{z}_u, y_u)\}$ can be fed into logistic regression for classification.



²Thomas N. Kipf, Max Welling: Semi-Supervised Classification with Graph Convolutional Networks. ICLR (Poster) 2017

Graph Convolutional Networks (GCN)

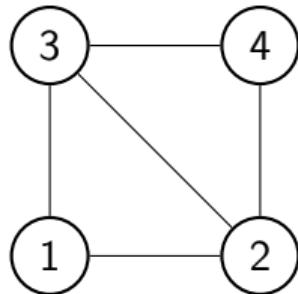
In matrix form

- ▶ $\mathbf{H}^{(0)} = \mathbf{X} \in \mathbb{R}^{N \times q}$
 - $\mathbf{H}^{(0)} = \begin{bmatrix} \mathbf{h}_1^{(0)} & \dots & \mathbf{h}_N^{(0)} \end{bmatrix}^\top$
- ▶ For $k = 1, \dots, K$
 - $\mathbf{H}^{(k)} = \sigma\left(\tilde{\mathbf{D}}^{-1/2}\tilde{\mathbf{A}}\tilde{\mathbf{D}}^{-1/2}\mathbf{H}^{(k-1)}\mathbf{W}^{(k)}\right)$
 - ▶ $\tilde{\mathbf{A}} = \mathbf{A} + \mathbf{I}$
 - ▶ $\tilde{D}_{ii} = \sum_j \tilde{A}_{ij}$
- ▶ $\mathbf{Z} = \mathbf{H}^{(K)}$

This is called graph convolutional network, because $\tilde{\mathbf{D}}^{-1/2}\tilde{\mathbf{A}}\tilde{\mathbf{D}}^{-1/2}\mathbf{H}^{(k-1)}$ is a (first-order) approximation of localized spectral filter, i.e., graph convolution.

GCN Propagation - A simple example

- ▶ Consider the graph below

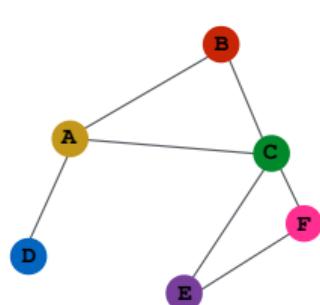


- ▶ We have

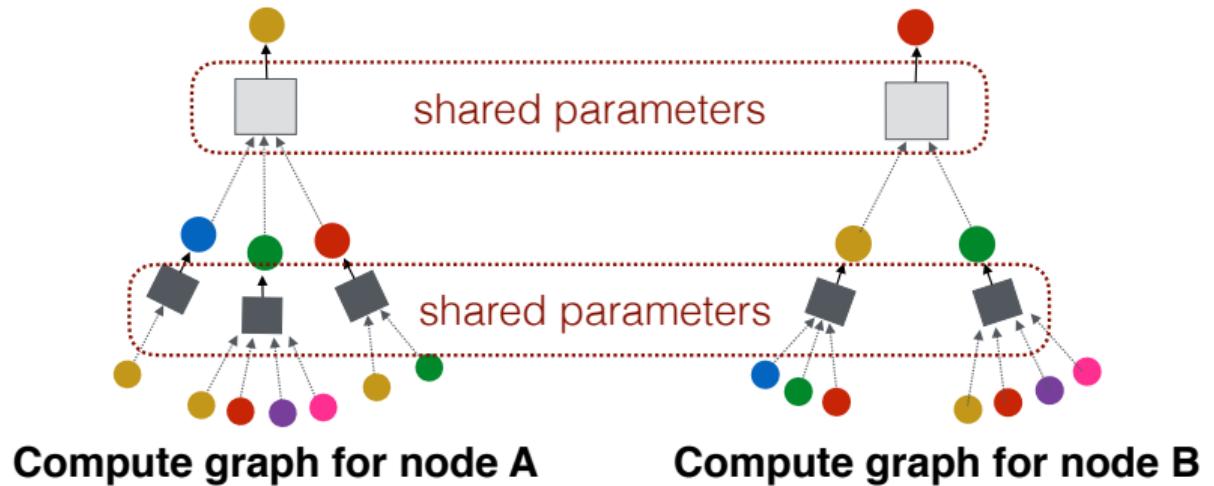
$$\tilde{\mathbf{A}} = \begin{bmatrix} 1 & 1 & 1 & 0 \\ 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \\ 0 & 1 & 1 & 1 \end{bmatrix} \quad \tilde{\mathbf{D}} = \begin{bmatrix} 3 & 0 & 0 & 0 \\ 0 & 4 & 0 & 0 \\ 0 & 0 & 4 & 0 \\ 0 & 0 & 0 & 3 \end{bmatrix} \quad \tilde{\mathbf{D}}^{-1} = \begin{bmatrix} 0.33 & 0 & 0 & 0 \\ 0 & 0.25 & 0 & 0 \\ 0 & 0 & 0.25 & 0 \\ 0 & 0 & 0 & 0.33 \end{bmatrix}$$

Inductive Capability

- ▶ The same weight parameters are shared for the same layer
- ▶ The number of weight parameters is sublinear in N
- ▶ This makes training on large graphs possible.

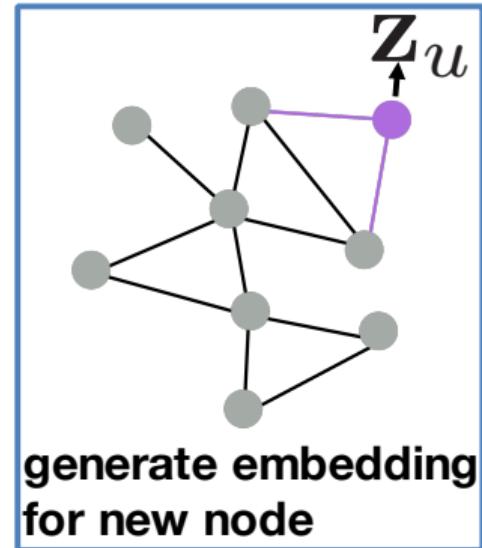
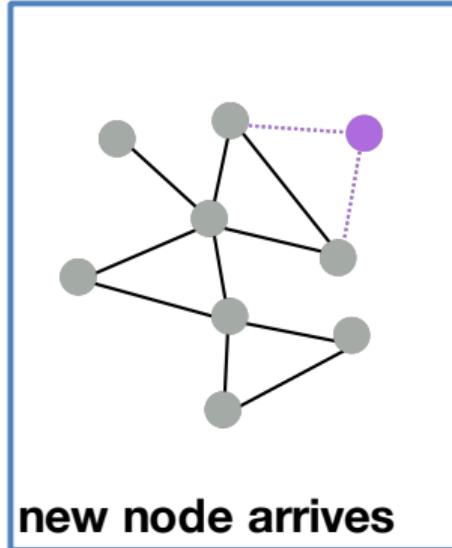
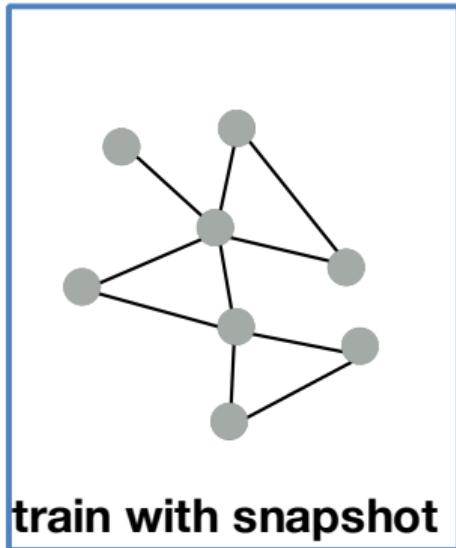


INPUT GRAPH

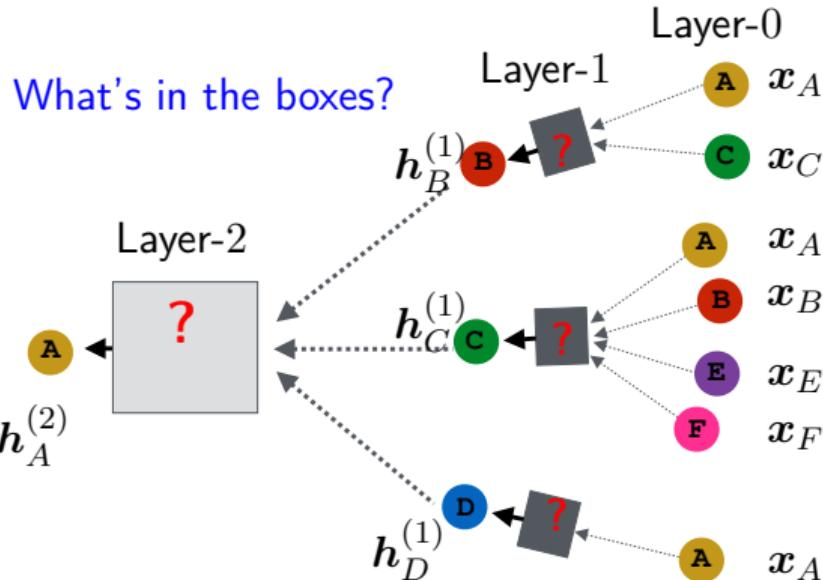
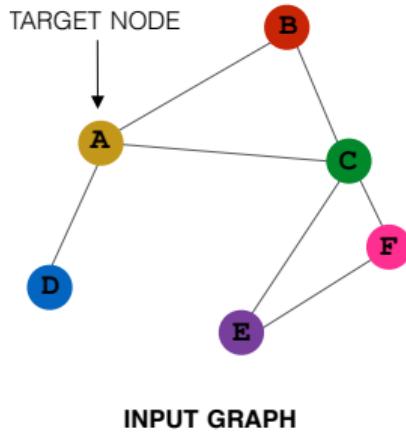


Inductive Capability

- ▶ The same weight parameters are shared for the same layer
- ▶ We can generalize to unseen nodes.



Other Models of Graph Neural Networks (GNNs)



- ▶ In general, each box conducts two operations

- $a_u^{(k)} = \text{AGGREGATE}^{(k)} \left(\{ h_v^{(k-1)} \mid v \in N(u) \} \right)$
- $h_u^{(k)} = \text{COMBINE}^{(k)} \left(h_u^{(k-1)}, a_u^{(k)} \right)$

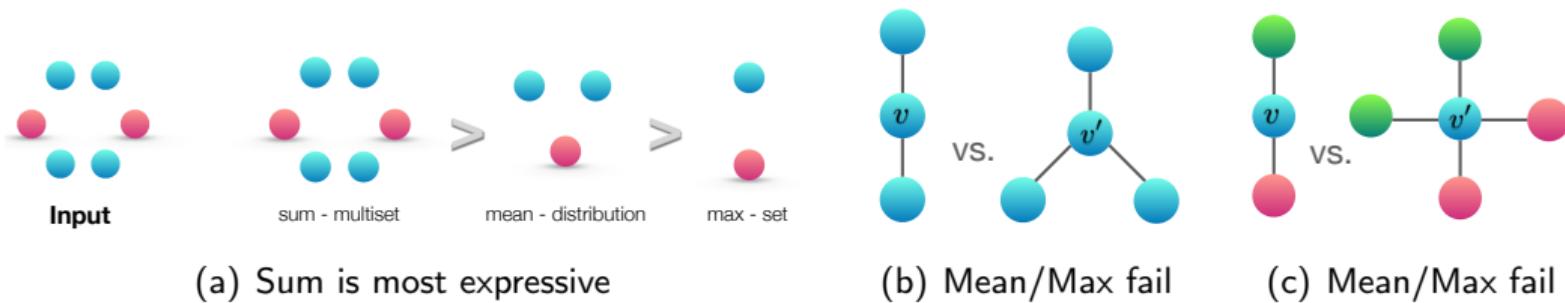
GraphSAGE ³

- ▶ Three variants of $\mathbf{a}_u^{(k)} = \text{AGGREGATE}^{(k)} \left(\left\{ \mathbf{h}_v^{(k-1)} \mid v \in N(u) \right\} \right)$
 - **Mean**: take a weighted average of neighbors
 - ▶ $\mathbf{a}_u^{(k)} = \sum_{v \in N(u)} \frac{\mathbf{h}_v^{(k-1)}}{|N(u)|}$
 - **Pool**: Transform neighbor vectors and apply symmetric vector function
 - ▶ $\mathbf{a}_u^{(k)} = \gamma \left(\left\{ \text{ReLU}(\mathbf{Q}^{(k)} \mathbf{h}_v^{(k-1)}) \mid v \in N(u) \right\} \right)$
 - ▶ $\gamma(\cdot)$ is element-wise mean/max, $\mathbf{Q}^{(k)}$ is weight parameter
 - **LSTM**: Apply LSTM to reshuffled neighbors
 - ▶ $\mathbf{a}_u^{(k)} = \text{LSTM} \left(\left[\mathbf{h}_v^{(k-1)} \mid v \in \pi(N(u)) \right] \right)$
- ▶ $\mathbf{h}_u^{(k)} = \text{COMBINE}^{(k)} \left(\mathbf{h}_u^{(k-1)}, \mathbf{a}_u^{(k)} \right)$: concatenation followed by a linear mapping
 - $\mathbf{h}_u^{(k)} = \text{ReLU} \left(\mathbf{W}^{(k)} \begin{bmatrix} \mathbf{h}_u^{(k-1)} \\ \mathbf{a}_u^{(k)} \end{bmatrix} \right)$

³William L. Hamilton, Zhitao Ying, Jure Leskovec: Inductive Representation Learning on Large Graphs. NIPS 2017: 1024-1034

Graph Isomorphism Network (GIN) ⁴

► $\mathbf{a}_u^{(k)} = \text{AGGREGATE}^{(k)} \left(\{\mathbf{h}_v^{(k-1)} \mid v \in N(u)\} \right) = \sum_{v \in N(u)} \mathbf{h}_v^{(k-1)}$



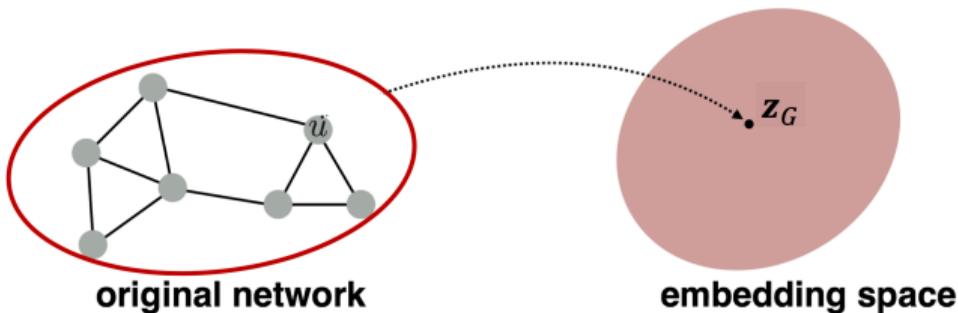
► $\mathbf{h}_u^{(k)} = \text{COMBINE}^{(k)}(\mathbf{h}_u^{(k-1)}, \mathbf{a}_u^{(k)}) = \text{MLP}^{(k)}((1 + \varepsilon^{(k)})\mathbf{h}_u^{(k-1)} + \mathbf{a}_u^{(k)})$

- $\text{MLP}^{(k)}$ is a 2-layer perceptron in the paper
- $\varepsilon^{(k)} = 0$ performs well in practice

⁴Keyulu Xu, Weihua Hu, Jure Leskovec, Stefanie Jegelka: How Powerful are Graph Neural Networks? ICLR 2019

Graph Embeddings

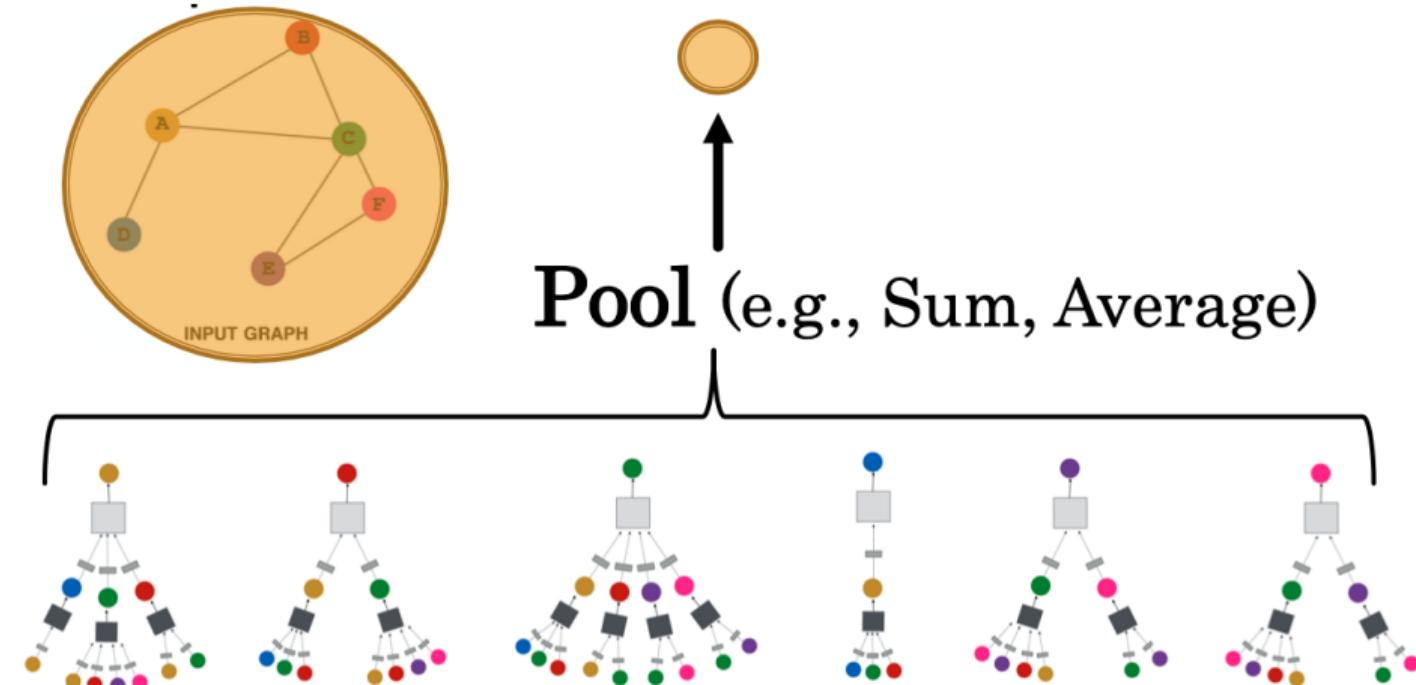
- ▶ So far we discussed generating embeddings for individual nodes
- ▶ How about mapping entire graphs to embeddings?



- ▶ Example applications: graph classification
 - Classifying toxic vs. non-toxic molecules
 - Identifying anomalous graphs

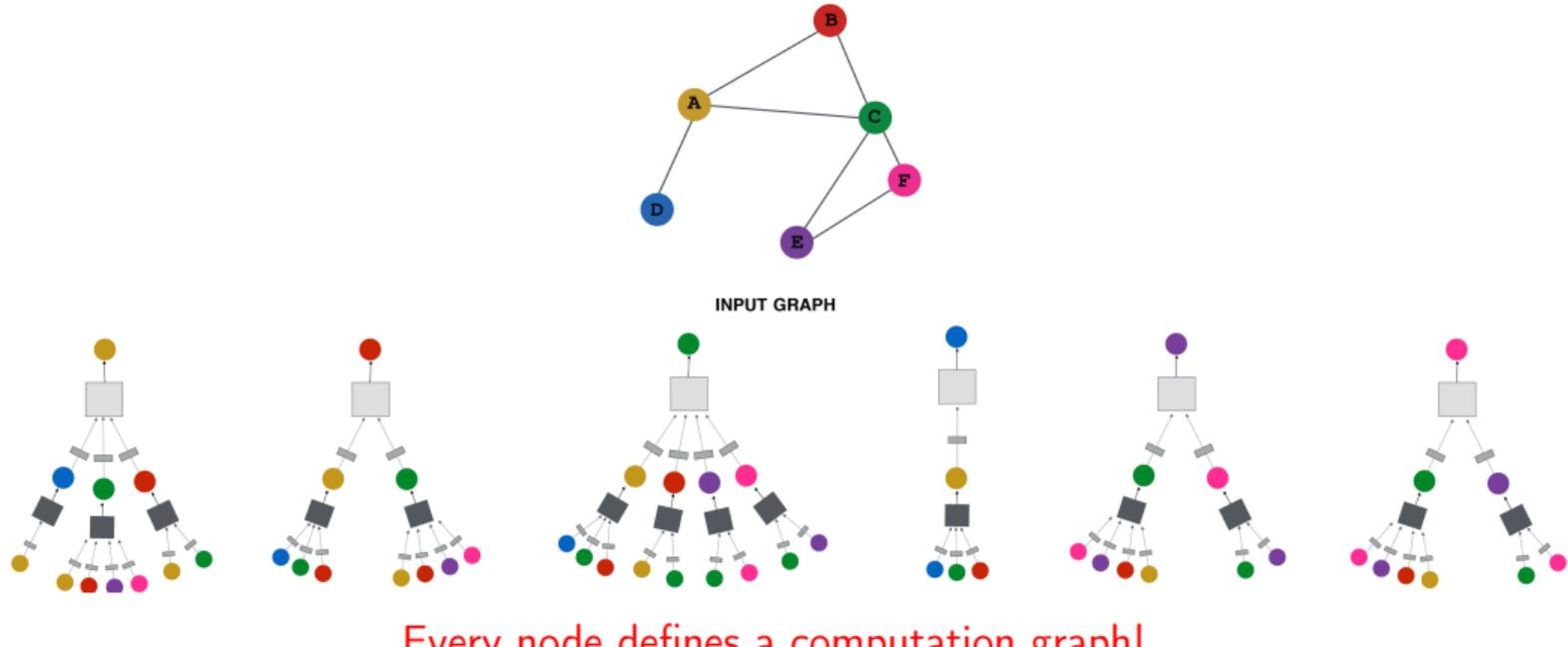
Graph Embeddings

- Obtain graph embedding by pooling node embeddings.



Limitation of GNNs

- ▶ Two nodes with the same computation graph will have the same embedding



Conclusion

- ▶ Deep models are able to learn/approximate complex functions $f(\mathbf{x})$ with few parameters
- ▶ Non-Euclidean nature of graphs makes the problem challenging
- ▶ Machine learning on graphs allows:
 - node classification & link prediction & graph classification ...
- ▶ Graph neural networks
 - provide end-to-end machine learning
 - are inductive

Reference

- ▶ Graph Representation Learning by William L. Hamilton:
https://www.cs.mcgill.ca/~wlh/grl_book/