Graph Neural Networks

Presented by: Enas Tarawneh

From "Shallow" to "Deep"

So far we have focused on "shallow" encoders, i.e. embedding lookups:



From "Shallow" to "Deep"

- Limitations of shallow encoding:
 - O(|V|) parameters are needed: there no parameter sharing and every node has its own unique embedding vector.
 - Inherently "transductive": It is impossible to generate embeddings for nodes that were not seen during training.
 - Do not incorporate node features: Many graphs have features that we can and should leverage.

From "Shallow" to "Deep"

 We will now discuss "deeper" methods based on graph neural networks.

$$ENC(v) = complex function that depends on graph structure.$$

 In general, all of these more complex encoders can be combined with the similarity functions from the previous section.

Outline for this Section

- We will now discuss "deeper" methods based on graph neural networks.
 - 1. The Basics
 - 2. Graph Convolutional Networks (GCNs)
 - 3. GraphSAGE

The Basics: Graph Neural Networks

Based on material from:

- Hamilton et al. 2017. <u>Representation Learning on Graphs: Methods</u> and <u>Applications</u>. *IEEE Data Engineering Bulletin on Graph Systems*.
- Scarselli et al. 2005. <u>The Graph Neural Network Model</u>. *IEEE Transactions on Neural Networks*.

Setup

- Assume we have a graph G:
 - V is the vertex set.
 - A is the adjacency matrix (assume binary).
 - $X \in \mathbb{R}^{m \times |V|}$ is a matrix of 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)

 Key idea: Generate node embeddings based on local neighborhoods.



 Intuition: Nodes aggregate information from their neighbors using neural networks





- Nodes have embeddings at each layer.
- Model can be arbitrary depth.
- "layer-0" embedding of node u is its input feature, i.e. x_u.



Neighborhood "Convolutions"

 Neighborhood aggregation can be viewed as a center-surround filter.



 Mathematically related to spectral graph convolutions (see <u>Bronstein et al., 2017</u>)

 Key distinctions are in how different approaches aggregate information across the layers.



 Basic approach: Average neighbor information and apply a neural network.



The Math

 Basic approach: Average neighbor messages and apply a neural network.



How do we train the model to generate "highquality" embeddings?





- After K-layers of neighborhood aggregation, we get output embeddings for each node.
- We can feed these embeddings into any loss function and run stochastic gradient descent to train the aggregation parameters.

- Train in an unsupervised manner using only the graph structure.
- Unsupervised loss function can be anything from the last section, e.g., based on
 - Random walks (node2vec, DeepWalk)
 - Graph factorization
 - i.e., train the model so that "similar" nodes have similar embeddings.

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



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



Overview of Model Design



Overview of Model Design



Overview of Model

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

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

Many application settings constantly encounter previously unseen nodes. e.g., Reddit, YouTube, GoogleScholar,

Need to generate new embeddings "on the fly"

Quick Recap

- Recap: Generate node embeddings by aggregating neighborhood information.
 - Allows for parameter sharing in the encoder.
 - Allows for inductive learning.
- We saw a basic variant of this idea... now we will cover some state of the art variants from the literature.

 Key distinctions are in how different approaches aggregate messages

Graph Convolutional Networks

Based on material from:

 Kipf et al., 2017. <u>Semisupervised Classification with Graph Convolutional</u> <u>Networks</u>. *ICLR*.

Graph Convolutional Networks

 Kipf et al.'s Graph Convolutional
 Networks (GCNs) are a slight variation on the neighborhood aggregation idea:

$$\mathbf{h}_{v}^{k} = \sigma \left(\mathbf{W}_{k} \sum_{u \in N(v) \cup v} \frac{\mathbf{h}_{u}^{k-1}}{\sqrt{|N(u)||N(v)|}} \right)$$

Graph Convolutional Networks

Graph Convolutional Networks

- Empirically, they found this configuration to give the best results.
 - More parameter sharing.
 - Down-weights high degree neighbors.

Outline for this Section

- 1. The Basics \checkmark
- 2. Graph Convolutional Networks
- 3. GraphSAGE

GraphSAGE

Based on material from:

 Hamilton et al., 2017. <u>Inductive Representation Learning on Large Graphs</u>. NIPS.

GraphSAGE Idea

 So far we have aggregated the neighbor messages by taking their (weighted) average, can we do better?

GraphSAGE Idea

GraphSAGE Differences

Simple neighborhood aggregation:

$$\mathbf{h}_{v}^{k} = \sigma \left(\mathbf{W}_{k} \sum_{u \in N(v)} \frac{\mathbf{h}_{u}^{k-1}}{|N(v)|} + \mathbf{B}_{k} \mathbf{h}_{v}^{k-1} \right)$$

GraphSAGE: $\mathbf{h}_{v}^{k} = \sigma\left(\left[\mathbf{W}_{k} \cdot \operatorname{AGG}\left(\left\{\mathbf{h}_{u}^{k-1}, \forall u \in N(v)\right\}\right), \mathbf{B}_{k}\mathbf{h}_{v}^{k-1}\right]\right)$ generalized aggregation

GraphSAGE Variants

Mean:

Pool

$$AGG = \sum_{u \in N(v)} \frac{\mathbf{h}_u^{k-1}}{|N(v)|}$$

Transform neighbor vectors and apply symmetric vector function.

$$AGG = \bigcap (\{\mathbf{Qh}_u^{k-1}, \forall u \in N(v)\})$$

- LSTM:
 - Apply LSTM to random permutation of neighbors. $AGG = LSTM ([\mathbf{h}_u^{k-1}, \forall u \in \pi(N(v))])$

Summary so far

Graph convolutional networks

- Average neighborhood information and stack neural networks.
- GraphSAGE
 - Generalized neighborhood aggregation.

This Talk

- 1) Node embeddings
 - Map nodes to low-dimensional embeddings.
- 2) Graph neural networks
 - Deep learning architectures for graphstructured data

Application

Decagon: A Graph Convolutional Approach to Polypharmacy Side Effects

Based on material from:

• Zitnik et al. 2018. <u>Modeling polypharmacy side effects with graph</u> <u>convolutional networks</u>. *Bioinformatics & ISMB*.

Polypharmacy Side Effects

Goal: Predict side effects of taking multiple drugs.

Polypharmacy Side Effects

- Polypharmacy is common to treat complex diseases and co-existing conditions
- High risk of side effects due to interactions
- 15% of the U.S. population affected
- Annual costs exceed \$177 billion
- Difficult to identify manually:
 - Rare, occur only in a subset of patients
 - Not observed in clinical testing

Modeling Polypharmacy

- Systematic experimental screening of drug interactions is challenging
- Idea: Computationally screen/predict polypharmacy side effects
 - Use molecular, pharmacological and patient population data
 - Guide strategies for combination treatments in patients

Data: Heterogeneous Graphs

Task Description

- Predict labeled edges between drugs
 - i.e., predict the likelihood that an edge
 (c, r₂, s) exists
- Meaning: Drug combination (c, s) leads to polypharmacy side effect r₂

Neural Architecture: Encoder

- Input: graph, additional node features
- Output: node embeddings

- r₁ Gastrointestinal bleed effect
- r₂ Bradycardia effect

- Drug target interaction
- Physical protein binding

Making Edge Predictions

- Input: Query drug pairs and their embeddings
- Output: predicted edges

O Drug O Gene

- r_1 Gastrointestinal bleed effect
- r₂ Bradycardia effect

Drug target interaction
 Physical protein binding

Feature vector

Experimental Setup

Data:

- Molecular: protein-protein interactions and drug target relationships
- Patient data: Side effects of individual drugs, polypharmacy side effects of drug combinations
- Setup:
 - Construct a heterogeneous graph of all the data
 - Train: Fit a model to predict known associations of drug pairs and polypharmacy side effects
 - Test: Given a query drug pair, predict candidate polypharmacy side effects

Prediction Performance

	AUROC	AUPRC	AP@50
Decagon (3-layer)	0.834	0.776	0.731
Decagon (2-layer)	0.809	0.762	0.713
RESCAL	0.693	0.613	0.476
Node2vec	0.725	0.708	0.643
Drug features	0.736	0.722	0.679

- Up to 54% improvement over baselines
- First opportunity to computationally flag polypharmacy side effects for follow-up analyses