

Table of Contents

Headers	Page #
Graph Representation using:	
- Matrices	3-4
- Dictionaries	5-6
Neighborhood Aggregation	7-10
GNN Applications	11
Types of GNNs	12-13
Neighborhood Attention and Generalized Update Methods	14
Graph Pooling	15
Fine-Tuning with Graph Augmentation	16
Types of Nonmolecular Graphs	17-18
Common Datasets for Graph Deep Learning	19
GNN Explainability Models	
- Basics	20
- GNNExplainer	21-22
- PGExplainer	23-24
- ProxyExplainer	25-26
Gumbel	27
Diffusion Modeling	
- Denoising Diffusion Probabilistic Modeling	28
- Discrete Denoise Diffusion Modeling	29
- DiffGrad	30
- RePaint Diffusion	31

Graph Representation - Matrices

- Graphs typically represented with...
- Something to define the structure: adjacency matrix / edge index
- For weighted graph the weight is often in A if it is primary attribute

$2 \times M$ - each column $[i]$ is connection from node i to node j

- Edge index representation used instead of adjacency matrix because space wise it is much more efficient

• First row is source node index of edges; second is target node index of edges; each column is an edge
 ↓
 each node is automatically assigned unique index when you create a graph

Node feature matrix X - $(N \times F)$

$$\text{node} \left[\begin{array}{c} f_1 \dots f_F \\ \vdots \\ f_N \end{array} \right] \quad \text{node features}$$

- For categorical features, either label encoding or one-hot encoding is used (usually one-hot); numerical values directly put in
- Label encoding falsely implies ordinal relationship (e.g. red=2, blue=3); should be nominal (no natural ordering)
- One-hot encoding will expand categorical features into multiple columns

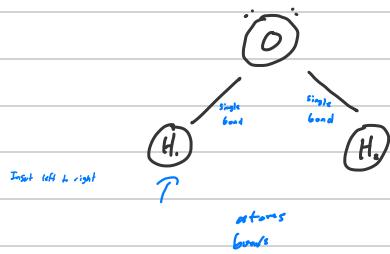
Edge feature matrix E - $(M \times F_e)$; typically used to encode features instead of weighted graph to account for multiple features

$$\text{edge} \left[\begin{array}{c} f_1 \dots f_{F_e} \\ \vdots \\ f_M \end{array} \right] \quad \text{edges features or edge}$$

Label mapping - $\{ \text{label: index, ...} \}$

Graph-level features (e.g. name) - stored in dictionary

Representing H_2O - Matrix Representation



Node Features

• Non-bonding electrons: node 1 ad node 2 = 0, node 2 = 4

• Label encoding only necessary if the labels provide useful information; sometimes also store labels/indices in separate dictionary to visualize more easily
 • One-hot encode node/atom type: hydrogen = $[1 \ 0]$, oxygen = $[0 \ 1]$

$$X = \text{node} \left[\begin{array}{c} \text{atom} \quad \text{non-bond} \\ \hline \text{H} \\ \text{O} \end{array} \right] = \left[\begin{array}{c} \text{H} \quad \begin{matrix} 1 & 0 & 0 \\ 0 & 1 & 4 \end{matrix} \\ \text{O} \quad \begin{matrix} 0 & 1 & 0 \\ 0 & 0 & 0 \end{matrix} \end{array} \right] \quad \text{new atom type} \Rightarrow \text{remove feature matrix and copy old features}$$

Structure/Edges

$$A = \left[\begin{array}{ccc} \text{H-H} & \text{H-O} & \text{O-H} \\ \text{H-H} & \text{H-O} & \text{O-H} \\ \text{H-H} & \text{H-O} & \text{O-H} \end{array} \right] = \left[\begin{array}{ccc} 0 & 1 & 0 \\ 1 & 0 & 1 \\ 0 & 1 & 0 \end{array} \right] \quad \text{Edge index} = \left[\begin{array}{c} 0 \ 1 \ 1 \ 2 \ 1 \end{array} \right]$$

Edge Features

$$E = \left[\begin{array}{c} \text{bond} \\ \text{node1-node2} \\ \text{node2-node3} \end{array} \right] = \left[\begin{array}{c} 1 \\ 1 \end{array} \right]$$

Mapping dictionary

node_indices = { 'H1': 0, 'O': 1, 'H2': 2 } node_indices['H1'] = 2

Global Features

graph = { 'name': 'water' }

Graph Representation - Dictionaries

• NetworkX represents these through dictionaries instead; for output it is list of tuples $[(\text{node}_1, \text{node}_2, \text{attribute_dict}), \dots]$ for easier interpretation

- node - dictionary that stores nodes and their attributes

{node1: attribute_dict1, ...}

e.g.

{

'A': {'color': 'red'}

'B': {'color': 'blue'}

}

- node['A'] = { 'color': 'red' }

- node['A']['color'] = 'red'

↑ key usually uses ID like index, not label

- adj - adjacency list containing node neighbors and edge attributes

{node1: {node1_neigh: {neighbor: 1-neighbor_attribute_dict}}

e.g.

{

'A': {

A-B

'B': { 'weight': 3, 'color': 'red' }

'neighbors': {}

}

'B': {

}

}

-adj['A']['B']['weight'] = 3

- graph - stores global graph attributes

{attribute-name1: attribute_value1,

e.g.

{

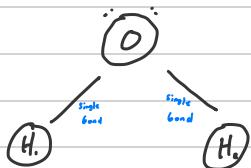
'name': 'water'

'type': 'undirected'

}

-graph['name'] = 'water'

Representing H_2O - Dictionary Representation



Node Features

• One-hot encode atom type : hydrogen = [1 0], oxygen = [0 1]

```
_node = {
    'H1': {
        'atom-type': [1 0], 'non-bonding': 0
    },
    'O': {
        'atom-type': [0 1], 'non-bonding': 4
    },
    'H2': {
        'atom-type': [1 0], 'non-bonding': 0
    }
}
```

Structure/edges

```
_adj = {
    'H1': {
        'O': {
            'bond-order': 1
        }
    },
    'O': {
        'H1': {'bond-order': 1},
        'H2': {'bond-order': 1}
    },
    'H2': {
        'O': {'bond-order': 1}
    }
}
```

Global Features

```
_graph = {'name': 'water'}
```

Message Passing

$$G = (V, E)$$

↓
set of vertices
set of edges

- Each $v \in V$ has features represented in vector
- Each edge (u, v) also has features
- Input to GNN can't just be flattened A because it is not permuting invariant

~~x~~ $\tilde{z}_v = \text{MLP}(A[v]) \oplus \dots \oplus A[u]$

$A \in \mathbb{R}^{N \times N} \Rightarrow A[v] \in \mathbb{R}^{N \times 1}$

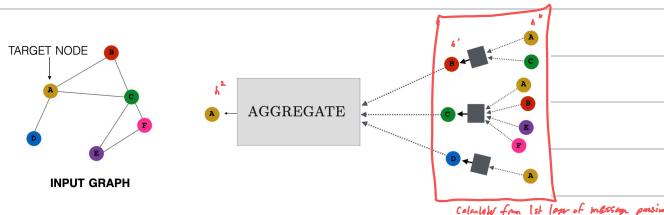
row of A

vector concatenation

- Since graph understanding depends on connection between all nodes, GNN must have some way to transfer information between them (via a CNN)
- Adjacency matrix cannot be used for this as related nodes might not be next to each other
- Message passing is solution to this

- \vec{h}_u^k - k th version of **hidden embedding** for each node $u \in V$ calculated by aggregation function
 - Aggregation takes all u 's neighbors $v \in N(u)$ and final update uses this message and u to form new node embedding
- Message passing can be multiple layers (looks at the neighbors of the neighbors essentially); 2-layer message passing means it is done twice
 - Stacking - having multiple sequential instances of a particular layer (e.g. message pass, fully connected/linear layers, convolution layers)

- The computation graph of GNN for aggregation forms a tree structure



$$\vec{h}_u^{(k+1)} = \text{Update}^{\text{out}}(\vec{h}_u^{(k)}, \text{Aggregate}^{\text{out}}(\{\vec{h}_v^{(k)} \mid v \in N(u)\}))$$

$$= \text{Update}^k(\vec{h}_u^{(k)}, m_{N(u)}^k)$$

includes non-linearity

the message aggregated from neighbours

- $\vec{h}_u^{(k)}$ on the input features $\vec{x}_u \forall u \in V$; aggregate takes set as input \Rightarrow no natural ordering of node's neighbors \Rightarrow aggregation functions must be **permutation invariant**
- Update and aggregate are differentiable functions

- Final embedding output for that final pass is node feature vector after K iterations of message passing
- $\vec{z}_u = \vec{h}_u^{(K)}, \forall u \in V$ more layers requires more convolution layers

- K iterations of message passing means each node contains info within/from a K -hop neighborhood
- \downarrow
K-layer message pass

More specifically...

$$\vec{h}_n^{(k)} = \sigma\left(W_{\text{self}}^{(k)} \vec{h}_n^{(k-1)} + W_{\text{neighbor}}^{(k)} \sum_{v \in N(n)} \vec{h}_v^{(k-1)} + \vec{b}^{(k)}\right) \quad k \geq 1$$

like vs $\vec{h}_n^{(k-1)}$ separate
different things and can make
certain things more clear

W and \vec{b} randomly initialized

- W matrices and \vec{b} contain trainable parameters
- For k layer message pass, you get info from k hops away, and W_{self} , W_{neighbor} , and \vec{b} are used to update the node features
 - You can either have separate W s and \vec{b} for each layer of message passing (i.e. first layer has its $W_{\text{self}}, W_{\text{neighbor}}, \vec{b}_{\text{self}}$) or share it across all layers
 - Separate parameters can mitigate the over-smoothing problem
 - or aggregating in a way where all nodes have the same
- For homogeneous graphs on smaller datasets, may be better to have some transformation at each step of message passing (share parameters)
 - similar architecture, structure
- Similar to MLP (linear operations followed by nonlinearity)
- Many types of message passing...

Sum Aggregation

$$\vec{m}_{N(n)} = \sum_{v \in N(n)} \vec{h}_v$$

→ Highly sensitive to node degrees and can be unstable/difficult to optimize

Mean/Average Aggregation

$$\vec{m}_{N(n)} = \frac{1}{|N(n)|} \sum_{v \in N(n)} \vec{h}_v$$

→ prevents issue of node degree sensitivity

$$\vec{m}_{N(n)} = \frac{1}{|N(n)|} \sum_{v \in N(n)} \vec{h}_v$$

• Solves ↗

Symmetric Normalization

$$\vec{m}_{N(n)} = \sum_{v \in N(n)} \frac{\vec{h}_v}{\sqrt{|N(n)| |N(v)|}}$$

• Motivated from spectral graph theory; normalization used for specific applications

• Normalization can make it hard to distinguish between nodes of different degrees

• Sometimes normalization can lead to a loss of information

$$\text{Update } (\vec{h}_n, \vec{m}_{N(n)}) = \sigma\left(W_{\text{self}} \vec{h}_n + W_{\text{neighbor}} \vec{m}_{N(n)} + \vec{b}\right)$$

W parameters influence how much to weight the current node in its updated representation

• Separate parameters for node's own features and neighbors allow model to learn different transformations for each

• $\sigma(W\vec{h} + \vec{b})$ is just 1 layer MLP (1 input / output)

Above is node level representation for symmetric normalization, but can also be represented graph level with matrices

Symmetrically normalized adjacency matrix = $\tilde{A} = D^{-1/2} A D^{-1/2}$

$d_{ii} = \sum_j a_{ij}$ = # of neighbors for node i

$$D^{-1/2} = \begin{bmatrix} \frac{1}{\sqrt{d_{11}}} & 0 \\ 0 & \frac{1}{\sqrt{d_{nn}}} \end{bmatrix} \quad A = \begin{bmatrix} -A_{11} & - \\ \vdots & \ddots \\ -A_{nn} & - \end{bmatrix}$$

$$D^{-1/2} A D^{-1/2} = \begin{bmatrix} \frac{1}{\sqrt{d_{11}}} & 0 \\ 0 & \frac{1}{\sqrt{d_{nn}}} \end{bmatrix} \begin{bmatrix} -A_{11} & - \\ \vdots & \ddots \\ -A_{nn} & - \end{bmatrix} \begin{bmatrix} \frac{1}{\sqrt{d_{11}}} & 0 \\ 0 & \frac{1}{\sqrt{d_{nn}}} \end{bmatrix}$$

$$= \begin{bmatrix} -\frac{1}{\sqrt{d_{11}}} A_{11} & - \\ \vdots & \ddots \\ -\frac{1}{\sqrt{d_{nn}}} A_{nn} & - \end{bmatrix} \begin{bmatrix} \frac{1}{\sqrt{d_{11}}} & 0 \\ 0 & \frac{1}{\sqrt{d_{nn}}} \end{bmatrix}$$

$$= \tilde{A}$$

Additional Set Aggregators

• Universal approximator - class of functions that can precisely approximate any continuous function given enough capacity/parameters

• Set function - takes a set as input

• Any aggregation function with following form is universal set function approximator

$$\vec{m}_{\phi}(a) = \text{MLP}_\phi \left(\sum_{v \in a} \text{MLP}_\phi(\vec{h}_v) \right)$$

*↑
and adds parameter to set
close to set function
+ this has been proven*

◦ This one isn't used if more so theoretical for a "best case" scenario

◦ Researchers aim to make aggregation functions with the power of this but more efficient

• Reduction function - takes a set/sequence of values and "reduces"/combines them into a smaller set of values or a single value

• Alternative reductions to sum, like element-wise max/min can be used

• Transductive learning - model learns from a and makes predictions on fixed set of data points; only works for data seen during training

◦ Designed for fixed dataset where you don't want to generalize (e.g., GNNExplainer)

• Inductive learning - model learns a function to generate embeddings based on node features and local graph structure

◦ Generalizes to unseen data (e.g., PGExplainer)

• Graph Sample and Aggregate (GraphSAGE)

◦ Samples a fixed-size set of neighbors instead of the full neighborhood; uses message passing scalably to large graphs

◦ Inductive approach

◦ Can generate embeddings for nodes it hasn't seen during training

GNN Applications

- Node classification (i.e. predicting if user is bot in a social network)
 - Each node $v \in V$ has associated label y_v
 - Given $G = (V, E)$, goal is to learn representation \vec{h}_v s.t. $f(\vec{h}_v) = y_v$
- Graph classification
 - Given $\{G_1, \dots, G_n\} \subseteq G$ and labels $\{y_1, \dots, y_n\} \subseteq Y$, learn graph embedding to predict graph label where $f(\vec{h}_G) = y_G$
- Relation/link prediction (i.e. content recommendation)
- Community detection

GNN Types

Basic GNN

- Overall GNN Pass for Classification

1. K-layer message passing and updating node representations repeat step 1 for each gnn layer
2. Get final node representations
3. Turn into graph embedding
4. Feed into classification model like MLP
5. Get classification \rightarrow calculate loss \rightarrow back prop \rightarrow gd \rightarrow update weights \rightarrow repeat until GNN is trained

- For picking subgraph after step 2 compute probability of node being part of subgraph then select the nodes

- Ex. Gumbel-Softmax (differentiable which allows for backprop)

- Can also train a GNN on a single graph for node level task

- Nodes split into train/val/test

- All your nodes in a vector for classification

- After message passing do step 5

Graph Convolutional Network (GCN)

- Uses self-loop message passing - aggregation function also includes self node removes need for update function
 - Aggregate($\{h_v^{(t-1)}, \forall v \in N(u) \cup u\}$)
 - Removes need for update function
- Familiar to sharing parameters of W_{self} or W_{global} since they are already combined with aggregating the message

- Also uses message normalization

$$\text{Thus, } \vec{h}_v^{(t)} = \sigma(W^{(t)} \sum_{\text{neighbors } u} \frac{\vec{h}_u}{\sqrt{|N(u)|}})$$

degree of v

- One "convolution" produces a new feature vector for a single node in the graph
 - Focuses on local neighborhood
 - Like CNN's there can be done in parallel

- Input dimension - # of features for each node or edge in graph

- The aggregation and update of node features occur in a hidden layer

- Dimension of hidden layer = dimension of $\vec{h}_v^{(t)}$ after aggregation/ update

- Message passing doesn't alter dimension of node \vec{h}_v , but $W\vec{h}_v$ does (# of output neurons)
 - Shape of W = output features \times input features

- dim(hidden layer) $>$ dim(input layer) captures more complex patterns

Graph Isomorphism Network (GIN)

- Most GNN's cannot distinguish between similar graph structures
- Injective aggregation update of WL Test makes it powerful to distinguish graphs
 - As a result GIN models injective functions
- Multiset - a set with possibly repeating elements
- $N(u)$ is a multiset because even if the nodes are distinct, their features can be the same
- GIN aggregation function over multiset must be injective and have great discriminative power
 - Discriminative power - how well different aggregation functions can distinguish different multisets
- Certain aggregation/readout function conditions can make GNN as powerful as WL Test

$$\vec{h}_v^{(k)} = \text{MLP}^{(k)} \left((1 + \epsilon^{(k)}) \cdot \vec{h}_v^{(k-1)} + \sum_{u \in N(v)} \vec{h}_u^{(k-1)} \right)$$

↑ weight importance
of central node

• Learnable parameters

- MLP weights/biases
- ϵ

• Hyperparameters

- # of neighborhood aggregations and node update steps (K) - determines how far info will propagate
- MLP architecture
- Readout/pooling function
- Learnable parameter initialization
- :

Neighborhood Attention and Generalized Update Methods

- Applies a trainable weight / attention to each neighbor when forming the message

$$\vec{m}_{(u)} = \sum_{v \in N(u)} \alpha_{u,v} \vec{h}_v$$

↑ weight of neighbor

13

$\vec{1}^T A \vec{1}$ will sum all the entries of A

- Oversmoothing can be thought of as the influence of each $\vec{h}_u^{(0)} = \vec{x}_0$ on all of the other final node embeddings $\vec{h}_v^{(k)}$, $v \in V$ in the graph

Can quantify the influence by the sum of the Jacobian for $\frac{\partial \vec{h}_v^{(k)}}{\partial \vec{h}_u^{(0)}}$

$$J_k(u, v) = \vec{1}^T \left(\frac{\partial \vec{h}_v^{(k)}}{\partial \vec{h}_u^{(0)}} \right) \vec{1}$$

↑ sum

- Other measures for sensitivity of Jacobian other than sum

- Helps prove nodes more likely to be reached by K steps from node u have a stronger impact on node u representation



Thm: For GNN models using self-update approach our specific aggregation function form, $J_k(u, v) \propto P_{G,K}(u|v)$

- Basically inflow of a input feature on another node v in graph after GNN aggregation depends on how likely v is to be reached in random walk

↑ probability of visiting node v in length k random walk from node u

- This theorem motivates various aggregation strategies for weighting each neighbor

- Graph Attention Networks use attention mechanisms to weigh importance of neighbor

- GAT (attention mechanism) uses linear layer multiplied by embeddings of each node

- More generally, we may use edge and graph level features in the GNN

$$\begin{aligned} h_{(u,v)}^{(k)} &= \text{UPDATE}_{\text{edge}}(h_{(u,v)}^{(k-1)}, h_u^{(k-1)}, h_v^{(k-1)}, h_g^{(k-1)}) \quad \text{constant weights possibly applied to each edge} \\ m_{N(u)}^{(k)} &= \text{AGGREGATE}_{\text{node}}(\{h_{(u,v)}^{(k)} \mid v \in N(u)\}) \\ h_u^{(k)} &= \text{UPDATE}_{\text{node}}(h_u^{(k-1)}, m_{N(u)}, h_g^{(k-1)}) \quad \text{↑ can switch order if depends} \\ h_g^{(k)} &= \text{UPDATE}_{\text{graph}}(h_g^{(k-1)}, \{h_u^{(k)} \mid u \in V\}, \{h_{(u,v)}^{(k)} \mid (u,v) \in E\}). \end{aligned}$$

all nodes all edges

Graph Pooling

- Embedding: transformation of one form of representation to a lower dimensional vector space
 - Dense vectors (few non-zero elements)
- f_p : pooling function
- \vec{z}_n : final node embedding for an epoch
- \vec{z}_G : graph embedding used for predictions at graph level

$$f_p(\{\vec{z}_1, \dots, \vec{z}_{|V|}\}) = \vec{z}_G$$

Mean Embedding

$$\vec{z}_G = \frac{\sum_{v \in V} \vec{z}_v}{f_n(|V|)}$$

↑
normalizing function (could be identity)

- Other embeddings like taking the max of each component across graph
- Sometimes multiple embeddings are concatenated for further information for the classifier

Fine Tuning with Graph Augmentation

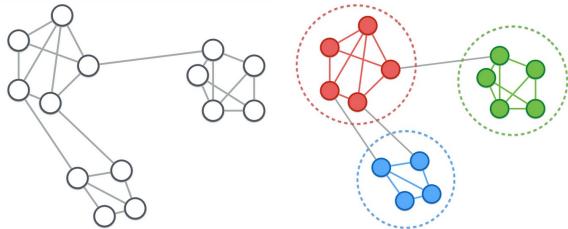
- Fine tuning involves...
 - Taking a pre-trained model and using a new (maybe augmented) dataset to slightly adjust model params to fit new task
 - Initial layers of pre-trained model are frozen (weights not updated) because they capture general features
 - Final layers retrained on new dataset or replaced with layers more suitable to the task
from scratch
- Freezing layer involves turning off `requires_grad` attribute for the tensors containing the parameters of the layer
- Various forms of graph augmentation that can make model more robust
- Node feature masking - using a boolean mask (set with threshold mask rate) and "turning off"/keeping random features for all nodes in the graph
- Edge perturbation/augmentation - uses boolean mask with threshold to remove edges from edge index

Copying Data

- Shallow copy - creates a new object; for nested object (e.g. 2D array) only copies references
- Deep copy - creates new object entirely independent of original (including nested)

Types of Nonmolecular Graphs

- Community graphs - vertices are grouped into clusters with dense connections internally and sparser connections between clusters
 - Modularity - a measure to assess how well-defined communities are within network
 - Some times you conceptually **division** / partition the network/graph in order to have high internal density and low external density of communities without actually changing the structure of the graph; various division algorithms to achieve best modularity
 and get high modularity ✓
 - This division can provide simplified analysis and better interpretability of the graph
- Clustering coefficient - measures tendency of nodes to cluster together locally (e.g. tendency for node's neighbors to also be neighbors)

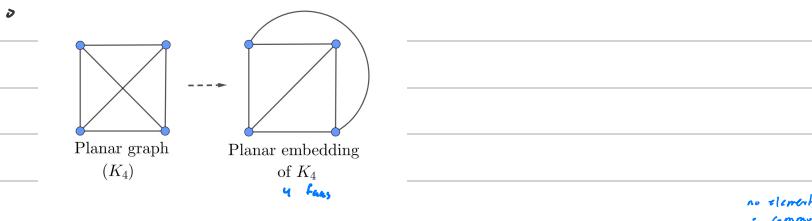


- Planar graphs - can be drawn in \mathbb{R}^2 without any edge intersections

Connected planar graph: planar graph with at least 1 edge between every pair of vertices

Faces: regions of planar graph surrounded by edges (including the **outer region**, the infinitely large region that surrounds the graph)

- K_n** - undirected graph where every pair of distinct vertices is connected by a unique edge



Bipartite Graph: no edges between vertices in the same set when divide into two **disjoint** sets U and V

- $K_{m,n}$ (Complete Bipartite Graph)** - a graph that has two disjoint sets of vertices where each vertex in one set is connected to every vertex in the other set (all edges between the sets exist)

Euler Formula: $|V| - |E| + |F| = 2 \Leftrightarrow G = (V, E)$ is connected planar graph

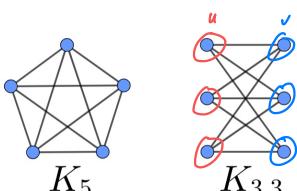
↑
Cardinality
of set of vertices ↑
faces

$$A-B \rightarrow A-C-B$$

$$A-B \leftarrow A-C-B$$

- Two graphs are **homeomorphic** if they can be transformed into the other through **edge subdivision** or **edge contraction** (the reverse)

Kuratowski's Theorem: A graph is planar \Leftrightarrow a graph does not contain a subgraph that is homeomorphic to K_5 or $K_{3,3}$



- Planar graphs used in areas like transportation to minimize overlap for efficient networks/layouts

Spectral Graph Theory - analyzes properties of graphs through eigenvalues/eigenvectors of matrices representing the graph

- A **spectrum** of a matrix is the set of its eigenvalues
- **Degree matrix (D)** - diagonal matrix where d_{ii} is degree of vertex i
- **Laplacian matrix (L)** - $L = D - A$
 \uparrow
adjacency matrix

- L is symmetric since A and D are symmetric
- 2nd largest eigenvalue of A refers to connectivity

Datasets

MUTAG dataset

- Collection of graphs representing molecules
- Nodes/atoms have categorical feature detailing type of atom encoded as one-hot vectors; some with edges for type of covalent bond
- Mutagenic - having the ability to cause permanent change in an organism's genes
 - Each of the 188 graphs has binary label indicating whether compound is mutagenic
 - Around 17 atoms per molecule on average

Baerbasi-Albert (BA) Model

- Generates graphs that exhibit scale-free property
- Node degree - number of edges connected to a node
 - When adding node to graph, node with higher degree is likely to make another (preferential attachment)
- Keep doing this until desired graph size is reached
- Motif - recurring, significant subgraph patterns within larger graph
 - Eg (cycle - nodes connected in closed loop, house - looks like a house)

Certain motifs can often explain various behavior/attributes for an entire graph (especially for molecules)

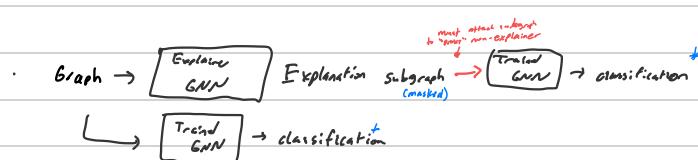
2 different motifs

BA2 Motifs

- Used to evaluate performance of GNN/explainer in detecting specific patterns within larger networks
- Contains graphs created using BA that contain motifs (can label parts of graph when motif is)
- Input BA into GNN \rightarrow learns to identify motifs \rightarrow classifies whole graph based on presence of motif
- Deterministic - premade dataset
 - 1000 BA graphs in dataset; 500 have 5-node acyclic motif and 500 have house motif
 - Goal is to classify graph into 1 of 2 classes based on motif it contains

Explainability

- Explainability can...
 - Increase trust in the model!
 - Improve the model's transparency
 - Allow people to better understand/interpret a GNN model and correct mistake patterns they make
- Subgraph shard matches labels of original graphs (e.g., mutagenicity); due to make sure it properly encapsulates whole graph
- Hard to identify "explainable" subgraphs in a graph because you can't feed them back into the network due to different size/distribution
 - Known as **Out-of-Distribution** problem
 - More specifically, distribution of data points fall outside distribution of training data (different connectivity, density, structure)
- Difference in structural properties prevents us by evaluating prediction capability of subgraph by feeding sub and full network into GNN (should be same label/output
→ Assume GNN already trained)
 - do or properly explained



- Explainer starts with random masks & edges/node features (parameters that will be adjusted)
 - Edge matrix applied to adjacency matrix, feature mask applied to feature matrix

Loss function

- Prediction loss: comparison between \hat{y} and y
 - Things with masks to highlight minimum edges/nodes for subgraph
 - Train until mask parameters converge and produce good explanations

GNNExplainer

- First general graph model explainer
- Given trained GNN classifier and its predictions, the explainer outputs subgraph structure and small subset of important node features that are important in GNN classifier prediction
- Incorporates relational information as well as node features
 - Works for an individual graph or set of graphs
- Mask - set of learned weights
- "Select" the important edges (for the subgraph) and node features by applying a mask
 - Each edge and node feature type gets a weight between 0 and 1 weighted by importance
 - Adjusted as the explainer is trained
 - Masks applied to A and X through element-wise multiplication
- Edge mask M is importance of all possible edges while $A \odot M$ represents important edges that exist in the graph
- Mask approach uses continuous relaxation - making a discrete problem into one with continuous values
- Mean field variational approximation - approximates complex distribution of all possible subgraphs by assuming the presence of an edge is not dependent on the other edges
- Mutual information - measure of how much information from one variable provides about another
- Optimization process
 - Randomly initialize probabilities of mask M (with all edges being important)
 - Apply mask to graph
 - Calculate conditional entropy (the loss)
• use COM (mixed graph) in place of G₀; it's not needed because container is out of distribution can't do
 - Calculate $\frac{\partial L}{\partial M}$ with backward propagation
 - Use gradient descent to update weights/probabilities of M
 - Repeat
- One does set threshold on masked graph to extract subgraph (Gumbel-Softmax not needed if container already trained)
- Maximizing mutual information = minimizing conditional entropy (as shown below)

2.9. Node classification

- $\mathcal{Y} = \{1, \dots, C\}$: all possible node classes

GNN aggregates info based off message passing (computation graph)

- Goal is to see which features and part of computation graph best "explains" GNN prediction

G_c - Node v computation graph (nodes from K iteration message pass)

$$\vec{z} \xrightarrow{G_c} \vec{z}' \xrightarrow{\text{out}} \hat{y}$$

- Goal is to find $G_s \subseteq G_c$ and node features $X_s \subseteq X_c$ that is most important in prediction \hat{y}

H - entropy term

$$H(Y | G=G_s, X=X_s)$$

↑
out
only

- This conditional entropy of Y given G_s and X_s measures how much inputting G_s and X_s (assuming i.i.d.) would reduce uncertainty of Y

$$\max_{G_s} \text{MI}(Y | G_s, X_s) = H(Y) - H(Y | G=G_s, X=X_s)$$

looking for subgraph G_s that
maximizes mutual information function
explain our GNN node class prediction

Assume all node features explain for now

- Node stably lowers probability of prediction \hat{y} : when removed from $G_c \Rightarrow$ import in classification \Rightarrow should be part of subgraph explanation

Expected value (E) - average outcome of experiment if repeated many times ($\in \mathbb{R}_{[0,1]}$)

$$H(Y | G_s, X_s) = -E_{Y|G_s, X_s} [\log(P_{\hat{y}}(Y | G_s, X_s))]$$

$\log_e(x)$ when x is 0-1

- $x > 1 \Rightarrow \log_e(x) > 0$ e.g. $\log_e(2) \rightarrow e^2 = 2 \rightarrow ?$ due to $e < 1 \rightarrow$ multiplying a number by itself x times only increasing product when > 1
- $x = 1 \Rightarrow \log_e(x) = 0$
- $0 < x < 1 \Rightarrow \log_e(x) < 0$ can't get negative number no matter how small number is

- Each possible $\hat{y} \in \mathcal{Y}$ has its own probability so E does weighted average of all these; allows explain to consider correct and incorrect decisions

Good explainer \Rightarrow given G_s and X_s prob. of predicting Y is high $\Rightarrow P_{\hat{y}}$ (between 0-1) is high / close to 1 $\Rightarrow \log(P_{\hat{y}})$ is small negative (near 0)

$\Rightarrow -\log(P_{\hat{y}})$ is small positive \Rightarrow max mutual info is higher

↑
large positive if explain is bad

Ultimate takeaway is we want to minimize conditional entropy across all subgraphs sampled from G ($G_s \sim G$)

- $\min_G H(Y | E_G[G_s], X_s)$
 - This finds entropy of average subgraph rather than every one
 - Added regularization will encourage smaller subgraph explainers

Although good, it does not have a global understanding of the classifier and does not take into account explanations for other graphs in the dataset

- For each new instance, optimization process is required (each graph learns its own mask for GNNExplainer)

- Each edge gets a weight/parameter: H parameters of num.edges

PGExplainer

- Parameterized explainer that uses MLP to parameterize generation process of explanations (only includes G_s)
- Has global understanding of explanations; don't need to retrain the explainer for each graph instance
 - MLP, parameterized by γ , is trained on all graphs in the dataset; bNNExplainer does the optimization independently and has to retrain for each new graph
 - Faster than bNNExplainer for new instances/large datasets

Different ways to interpret DNN (learning)

- **Whitebox / transparent method** - use the model's internal structure/parameters to provide insight into the model's decisions
 - e.g. Gradient to see which input features have greatest effect on decision
 - More model specific; focuses on forward/backward propagation
- **Blackbox method** - use only the input and output of the model to understand its behavior
 - i.e. creates a simpler, approximate model of a more complex model

$$G_o \xrightarrow{\text{GNN classifier}} Y_o$$

- PGExplainer also aims to maximize mutual information between $G_s \subseteq G_o$ and $Y_o \Rightarrow$ minimize the conditional entropy

$$\max_{G_s} MI(Y_o, G_s) = H(Y_o) - H(Y_o | G = G_s)$$

- $2^{\binom{|V|}{2}}$ possibilities for G_s , which is inefficient to optimize this function
 - Needs relaxation for picking G_s efficiently; assume graph edges are independent (Gilbert random graph)

$$P(G) = \prod_{(i,j) \in E} P(e_{i,j})$$

- Edge probabilities start off in bernoulli distribution

$$P(e_{i,j}) = \theta_{i,j}$$

- $\pi(\Theta)$ - distribution over different configurations of G_s ; basically the mask of edge probabilities?
 - The optimization will try to learn $\pi(\Theta)$

- $G_s \sim \pi(\Theta)$ involves randomly generating explainer subgraphs from this distribution; we approximate this

• Ω - Latent variables (all w_{ij})

$$\Omega = \text{MLP}_\Phi(\mathcal{Z})$$

- MLP outputs the probabilities of certain edges being in explainer for G_o ; input edges which contain edge/node features

$$G_s \approx G_o = f(G_o, \Omega, T, \epsilon)$$

- f comes as Gumbel-Softmax; $\Omega = \text{all } w_{ij}$

- Temperature (T) - adjusts sharpness/determinism of edge selection process; applied after MLP output

- $T \rightarrow 0 \Rightarrow P(e_{ij}) \rightarrow 0 \text{ or } 1$

- $T \rightarrow \infty \Rightarrow P(e_{ij}) \rightarrow 0.5$

- High temp \Rightarrow more exploration; usually used during training

- For node classification, explainer subgraph differs per node v , so this should be part of explainer function (MLP)

$$w_{ij} = \text{MLP}_\Psi([z_i^*; z_j^*; z_v])$$

; trying to see which edges of computation graph to \vec{z}_v are most important

- For graph classification graph prediction is not conditional per node so \vec{z}_v not needed

- $w_{ij} = \text{MLP}_\Psi(z_i^* \oplus z_j^*) = w_{ij}$

- Do this for all the edges simultaneously

sigmoid

- Final edge weight is some form of $\sigma\left(\frac{w_{ij}}{T}\right)$

- Loss function = Cross-entropy ($\text{GNN}(G_o)$, $\text{GNN}(G_s)$)

- Could be between nodes in the graph

- Variants regularization techniques to get more compact explainer and include neighboring edges (nearest to same node) in the explainer

• Training PGExplainer

- Per epoch For G_o & G_s

- $\text{GNN}(G_o) = Y_o$ note embeddings from GNN classifier

- $\text{MLP}_\Psi(Z) = W$ latent variables

- Extract subgraph

- Obtain prediction for subgraph (assumes in distribution)

- Compute loss (minimize cross-entropy/maximize similarity between predictions)

- Update Ψ with backprop and grad

- Uses the reparameterization trick in * to allow for backpropagation

- See below for VAE

- For inference after MLP extracts edge probabilities threshold can be used to extract subgraph

Algorithm 2: Training algorithm for explaining graph classification

```

1: Input: A set of input graphs with  $i$ -th graph represented by  $G_o^{(i)}$ , node features  $\mathbf{X}^{(i)}$ , and a label  $Y^{(i)}$ , a trained GNN model:  $\text{GNNE}_{\Phi_0}(\cdot)$  and  $\text{GNNC}_{\Phi_1}(\cdot)$ .
2: for each graph  $G_o^{(i)}$  do
3:    $\mathbf{Z}^{(i)} \leftarrow \text{GNNE}_{\Phi_0}(G_o^{(i)}, \mathbf{X}^{(i)})$ .
4:    $Y_o^{(i)} \leftarrow \text{GNNC}_{\Phi_1}(\mathbf{Z}^{(i)})$ .
5: end for
6: for each epoch do
7:   for each graph  $G_o^{(i)}$  do
8:      $\Omega \leftarrow$  latent variables calculated with Eq. (11)
9:     for  $k \leftarrow 1$  to  $K$  do
10:       $\hat{G}_s^{(i,k)} \leftarrow$  sampled from Eq. (4).
11:       $\hat{Y}_s^{(i,k)} \leftarrow \text{GNNC}_{\Phi_1}(\text{GNNE}_{\Phi_0}(\hat{G}_s^{(i,k)}, \mathbf{X}^{(i)}))$ 
12:    end for
13:   end for
14:   Compute loss with Eq. (9).
15:   Update parameters  $\Psi$  with backpropagation.
16: end for

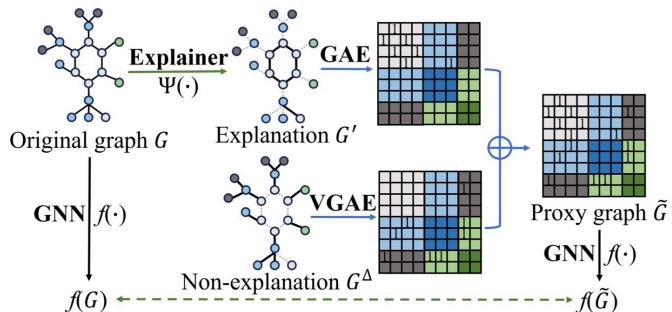
```

Variational Autoencoder (VAE)

- Autoencoder - compresses data into smaller representation
 - Encoder-layers who output have less dimensions than input (encoder maps to the latent space)
 - Decoder - takes the output of encoder to try to reconstruct the input
- Variational autoencoder - rather than $\text{encoder}(\vec{x}) = \vec{\text{latent}}$, $\text{encoder}(\vec{x}) = \frac{\vec{\text{mean}}}{\vec{\text{sd}}} \xrightarrow{\text{Sample from distribution}} \text{decoder}(\vec{z}) = \text{input}$
 - Loss = Reconstruction loss + KL Divergence
- Sampling prevents backprop so use reparameterization trick
 - $\vec{z} = \mu + \sigma \odot \epsilon$
↳ μ is mean, σ is standard deviation, ϵ is random noise
 - Sampling not done within \vec{z}
 - (1) \vec{z} $\xrightarrow{(2)} \mu + \sigma \odot \epsilon$ $\xrightarrow{(3)} \vec{z}$
- Can do backprop on the deterministic nodes and ignore the Sampling
- Disentangled variant autoencoders includes a β that weights how much KL Divergence is present
- Variational Graph Autoencoder - uses encoded node/graph embeddings for graph reconstruction

Proxy Explainer

- In practice difficult to train explainer with classifier G_x since G_x is out of distribution
- Paper suggests to use graph VAEs & generate classifier and non-explainer and combine to form in distribution \tilde{G} which can then be put in the classifier



- Kullback-Leibler (KL) divergence is a measure of the difference between two probability distributions

- Always ≥ 0

- Measures the amount of information lost when Q is used to approximate P

$$KL(P||Q) = H(P, Q) - H(P)$$

- Equivalent to difference in conditional entropy

$$D_{KL}(P||Q) = \sum_i P(i) \log \frac{P(i)}{Q(i)}$$

- Can confuse distributions P and $Q \Rightarrow$ they are similar \Rightarrow assign similar probabilities to the same sequences

Algorithm

- Generate observations \mathcal{O}

- Calculate probability of Q generating the observations

Gumbel

- Function for picking edges for subgraph needs to be continuous so it is differentiable and explainer can be trained
- Additionally, threshold such as 0.51 probability of being in explainer subgraph is treated same as 0.99

Uses Gumbel distribution

- Gumbel noise - allows for exploration for less plausible outcomes, but introduces randomness that mostly preserves edge importance

• Argmax - returns index of max value

• Simplex - space of probability distribution

Gumbel-Max

- Input: probabilities of each edge importance
- Convert to logits: $\theta_i = \log(\frac{p_i}{1-p_i})$ for numerical stability later on
- Calculate and add Gumbel noise to logit
- Choose edge with highest val with argmax (discrete / not differentiable)

$$\vec{z} = \text{one-hot}(\underset{\substack{\text{argmax} \\ \text{prob}}}{\text{argmax}} [\vec{g}_i + \text{log } \pi_i])$$

Gumbel-softmax

- Uses softmax as differentiable approximation instead of argmax
- $\frac{1}{T}$ is used for edge sharpness
- Density of Gumbel-softmax distribution is more ~~smooth~~ with low temp and ~~smooth~~ with high temp
- Forward pass will discrete edges with argmax but backward pass will use approximation for derivative?

Discrete Denoise Diffusion Modeling (DDDM)

• Forward Diffusion Process

- Through Markov chain, add some noise to the sample each step t ; noise defined by transition matrices \mathbb{Q}^t amount of times you apply the noise
- markov chain
current state

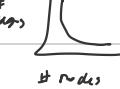
• Reverse Denoising Process

- Neural network trained to predict original data from noisy samples at each step
- Model learns to remove noise

• Has GNN, embed timestep t of diffusion into it to see how corrupt it is

• The model adds discrete graph edits in steps (edge addition/deletion, feature edit)

- Train graph transformer network
- Perturbed graph
- Sparsity - # of edges < max possible edges
- Gaussian is continuous while discrete is more real; Gaussian blurs distinction between connected or not connected
- Connection/Lack of is important for nodes/edges so use discrete
- Distribution in graphs rel to the statistical patterns of degrees, node/edge features (e.g. path length distribution, edge feature distribution, clustering coefficient)

Ex. node degree distribution is  **power law**

"Noise" for graph means adding/removing edges for discrete; much more realistic and has meaning (e.g. bond either does or doesn't exist)



- Loss function is between $*$ and $*$

- Multiple components to define loss between graphs (e.g. node/edge features, degrees, weighted sum of these losses)
- Once transformer is trained, add more noise to your dataset and repeat the process
- Training the DDDM for all graphs in dataset makes DDDM (random graph) = in distribution graph

Denoising Diffusion Probabilistic Models (DDPM)

• T - maximum number of steps in the diffusion process

• Goal of DDPM is to reverse the forward diffusion process

◦ Loss function is between the probability distributions

• Per training example $\vec{x}_0 \in \text{Dataset}$

◦ Randomly choose / sample timestep $t \in [1, T]$ and noise ϵ

◦ Add appropriate amount of denoising noise for that timestep to create \vec{x}_t

◦ Model predicts the noise added

◦ Calculate MSE between predicted noise and actual noise

◦ Backprop and grad

• Once model is accurate, we can increase T and continue to train

• Inference

◦ Perform denoising for each timestep

$$\vec{x}_{t-1} = \text{scaling}(\vec{x}_t - \text{scaling} \epsilon_{\theta}(\vec{x}_t, t) + \sigma_t \vec{\epsilon})$$

◦ $\epsilon_{\theta}(\vec{x}_t, t)$ - noise predicted by model

◦ σ_t - 2.0 for last t

◦ scaling noise
scale
decreasing variance

DiGress (Discrete Graph Denoising Diffusion)

- Transformers are good at handling sequential data (like iteratively adding noise to a graph which is why they are the model used)
- A PDDM model that uses a transformer and explicitly models both nodes and edges
 - Graph transformer has attention mechanism that can simultaneously consider all nodes
- Discrete PDDM is better since it maintains structural properties of graph
- Rather than trying to predict noise added from previous timesteps, it directly tries to predict \vec{x}_0
 - This results in model having more consistent target which results in better performance
- For diffusion model to be efficient...
 - \vec{x}_0 should be calculated directly from \vec{x}_t
 - It should directly control \vec{x}_0
 - The state of the data as $t \rightarrow \infty$ should not depend on \vec{x}_0 (starting point for inference generation is always the same)

$$\vec{x} \rightarrow z^t = \bar{Q}^t \vec{x} = Q^t \dots Q^0 \vec{x}$$

- $G^t \rightarrow G^{t+1} \Rightarrow$ Noise is applied independently to each node and edge feature

Cross-entropy transformer (G^t), G^0)

- Parameters ϕ for transformer ϕ are adjusted for denoising (e.g. attention weights)
- G^t is not a graph but rather probabilities for each node and edge
- Cross entropy loss is between these probabilities and the original graph
- λ - hyperparameter for importance of edges over nodes

DiGress code uses MSE

$$l(p^t, G) = \sum_{i,j} \text{cross-entropy}(x_{ij}, \hat{p}_{ij}^t) + \lambda \sum_{i,j} \text{cross-entropy}(e_{ij}, \hat{p}_{ij}^t)$$

- DiGress is permutation invariant (as \vec{x}_0 since $n!$ matrices can represent the same graph)

Training

$$G_0 \rightarrow G^t \sim X \bar{Q}^t_x \text{ and } E \bar{Q}^t_e \rightarrow z = f(G^t, t) \rightarrow \text{prob-transformer}(G^t, z) \rightarrow \text{calculate loss and optimize}$$

get structural and spectral features

Sample row partitions for nodes and edges

- More features describing the graph can be calculated and added which improves DiGress performance (like eigenvalues / values of each atom)
- Singular value - eigenvalue of matrix $A^T A$ (square)
 - Takes long time to compute for large graphs (3 days on A100)

Inference

- Iteratively denoises from $t=T$ to 1
- $z = f(G^t, t) \rightarrow \text{prob-transformer}(G^t, z) \rightarrow$ find probability distributions of graph at $t=1$ based on predictions for $\vec{x}_0 \rightarrow$ sample new graph of $t=1$

RePaint Diffusion

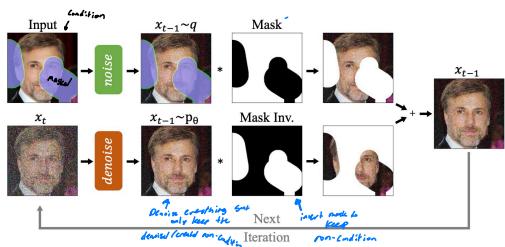
- Condition - does not change when passed through the diffusion model
- The binary mask
 - Masked region(0) - defines which parts of the input need to be generated/inpainted
 - Unmasked region(1) - the condition

- Input is completely noisy x
- At each time step replace unmasked region with the condition

Resampling allows Repaint to make multiple attempts at denoising per timestep

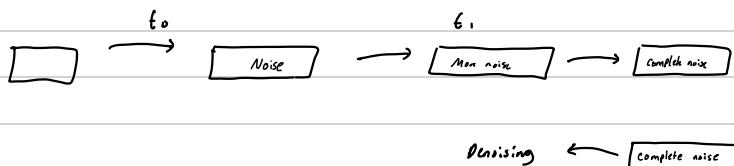
Inference

- Initialize \hat{x}^t with random noise
- For unmasked regions, use condition with sample noise
- For masked regions, use models prediction to set mask step
- Apply the mask
- Resample



- Condition is still noised but not with the DDPM
- Repaint is an algorithm that uses a pre-trained DDPM

At



- Using trained diffusion model, at each step of denoising insert the original condition in the corresponding space