Graph Neural Networks

Applied Machine Learning, KU

Daniel Murnane - May 8th, 2024



8/05/2024

1

Introduction & Goals

- Goals for today:
 - Understand point cloud structure
 - Understand graph structure
 - Look at some real data that is represented as a graph
 - Understand how to generalize from grid (image) to a graph
 - Think about some choices of GNN architectures
 - Zoom in to the Transformer as a kind of GNN
- Have borrowed content from Troels' slides from last year!





How can we represent data



8/05/2024

Graph Neural Networks – Daniel Murnane

Data types

- Let's tidy up our language: There are data types, and data structures
- A data type is an input to our neural network, it falls into one of the following types





- Let's tidy up our language: There are data types, and data structures
- A data type is an input to our neural network, it falls into one of the following types
- A data structure defines how each component of our dataset is related
- There are many data structures...



• Tabular: Each sample has one entry





- Tabular: Each sample has one entry
- Set: Each sample has a variable number of entries





- Tabular: Each sample has one entry
- Set: Each sample has a variable number of entries
- Point Cloud: Each sample has a variable number of entries, with some notion of distance



Class 1



- Tabular: Each sample has one entry
- Set: Each sample has a variable number of entries
- Point Cloud: Each sample has a variable number of entries, with some notion of distance
- Grid: Each sample has a fixed number of entries, binned into a grid, with a natural distance



Image 1



- Tabular: Each sample has one entry
- Set: Each sample has a variable number of entries
- Point Cloud: Each sample has a variable number of entries, with some notion of distance
- Grid: Each sample has a fixed number of entries, binned into a grid, with a natural distance
- Sequence: Each sample has an ordered list of variable number of entries



Class 1



- Tabular: Each sample has one entry
- Set: Each sample has a variable number of entries
- Point Cloud: Each sample has a variable number of entries, with some notion of distance
- Grid: Each sample has a fixed number of entries, binned into a grid, with a natural distance
- Sequence: Each sample has an ordered list of variable number of entries
- Time series: Each sample has an ordered list of variable number of entries, with a neighbor distance given by time



Day 1

• Tabular: Each sample has one entry

UNIVERSITY OF COPENHAGEN

- Set: Each sample has a variable number of entries
- Point Cloud: Each sample has a variable number of entries, with some notion of distance
- Grid: Each sample has a fixed number of entries, binned into a grid, with a natural distance
- Sequence: Each sample has an ordered list of variable number of entries
- Time series: Each sample has an ordered list of variable number of entries, with a neighbor distance given by time
- Graph: Each sample has a variable number of entries, with neighborhoods given by explicit pairwise relationships



Class 1

UNIVERSITY OF COPENHAGEN

- Graph: Each sample has a variable number of entries, with neighborhoods given by explicit pairwise relationships
- A graph is a collection of *nodes* (objects or entries) and *edges* (relationships between each object)
- Nodes can have features, edges can have features
- A graph may also have graph-level or "global" properties, e.g. class_1["topic"] = "applied ML"



Class 1 = Applied ML

• I encourage you to look at all data structures through the eyes of a graph

8/05/2024





Class 1 = Applied ML

Graph Neural Networks – Daniel Murnane

Graph Data



8/05/2024

Classic Problems with Graphs



8/05/2024

Graph Neural Networks – Daniel Murnane

Representing a Graph

- Nodes a list of node vectors n_i^k , where $i = 0, ..., N_n$ up to number of nodes, $k = 0, ..., N_f$ up to number of features
- Edges an adjacency matrix A_{ij} , where $A_{ij} = 1$ when there is an edge between node i and node j, and 0 otherwise
- In practice, A_{ij} is mostly 0s, so let's represent it with a more efficient structure: E_{ij} , which is a list of *pairs* of node indices



Removing the Grid



8/05/2024

Graph Neural Networks – Daniel Murnane

• Remember how the image convolution worked:





• Let's rewrite this big tensor multiplication into smaller pieces...





8/05/2024

Now this looks like each pixel in the window is passed through a simple linear layer of a feedforward NN





- Now this looks like each pixel in the window is passed through a simple linear layer of a feedforward NN
- Once we have passed each pixel through, we aggregate (in this case sum) the updated pixels







- Now this looks like each pixel in the window is passed through a simple linear layer of a feedforward NN
- Once we have passed each pixel through, we aggregate (in this case sum) the updated pixels
- The "center" of our 2x2 neighborhood is updated with the convolution

8/05/2024

outputs

UNIVERSITY OF COPENHAGEN



- The "center" of our 2x2 neighborhood is updated with the convolution outputs
- These centers form the inputs for the *next* convolution, which again happens in an $k \times k$ grid neighborhood





Now, let's throw away the grid structure

- Grids are convenient for images, but the world is not a grid
- Instead of "pixels", let's call each hidden vector now a "node"





Now, let's throw away the grid structure

- Grids are convenient for images, but the world is not a grid
- Instead of "pixels", let's call each hidden vector now a "node"
- Let's also make each node the center of its own neighborhood
- And let's do our convolution over nearby nodes

COPENHAGEN

• E.g. this would be one such neighborhood:



One final tweak!

- Let's also make each node the center of its own neighborhood
- And let's do our convolution over nearby nodes
- E.g. this would be one such neighborhood
- In the CNN, each pixel had its own dedicated W. This was easy, because every neighborhood had a fixed number of pixels



• Now there are arbitrary number of nodes in each neighborhood, let's be even simpler and use the same W_{MN} for every node multiplication **COPENHAGEN** 8/05/2024 Graph Neural Networks – Daniel Murnane

This is a GNN!

- It's a simple one a graph convolution network
- The $h_M \times W_{MN}$ step is a "node update"

8/05/2024

- Passing these node features to the center node is called "message passing"
- Summing all the features is called "aggregation"
- Node update, message passing and aggregation are the building blocks for basically all graph neural networks



0.2

1.9

1.3

GNN Architectures



8/05/2024

Graph Neural Networks – Daniel Murnane

Message passing

- We can create a language for convolutions that will work for almost every other ML model
- It's called "message passing", and it has two parts: calculate the message between pairs of objects, then aggregate the messages coming into each object's neighborhood
- For a graph, a message is the features that are passed along an edge a learnable function that takes in the two nodes on either side of that edge



8/05/2024

- *f_e* is a learnable function (feed-forward neural network)
- n_i^k are the N_k hidden channels of the N_i nodes in the graph



Aggregation & Permutation Invariance

- We have a message function $m_{ij} = f_e(n_i^k, n_j^k)$, so how do we combine these messages around each node?
- We could stack them together and pass them through a FFNN? E.g. $f([m_{01}, m_{02}])$
- Two problems:
- 1. A FFNN has a fixed size, but we might have any number of incoming messages



2. If we switch the order that we receive the messages (which is meaningless), the output of the FFNN will be different!

Aggregation & Permutation Invariance

- What we are looking for is a "permutation invariant" way to combine any number of incoming messages
- There are a few ways to do this, but the simplest is to take the sum (or mean/max/min)
- The choice of aggregation function is a hyperparameter (like the CNN pooling step)







Node update

UNIVERSITY OF

- The final step is to pass this aggregated information through a node FFNN
- The f_e is an edge-wise function/network, it is the same for every edge
- The f_n is a node function/network, it is the same for every node
- Note that even though the function is the same across all nodes or edges, the *outputs* of the function depend on the node or edge features



$$n_0'^k = f_n(f_e(n_0^k, n_1^k) + f_e(n_0^k, n_2^k))$$

Let's put it all together...



Let's put it all together...



Let's put it all together...



Let's assume our "learned" message function is just:

$$f_e(n_i^k, n_j^k) = n_i^k - n_j^k$$

and our "learned" node update function is just:

$$f_n(a_i^k, n_i^k) = \frac{1}{2}(a_i^k + n_i^k)$$



1. Calculate messages

$$m_{ij}^l = f_e\big(n_i^k, n_j^k\big)$$

$$E_{ij} = \begin{bmatrix} 0 & 0 & 0 & 2 & 1 & 2 & 3 & 3 \\ 1 & 2 & 3 & 3 & 0 & 0 & 0 & 2 \end{bmatrix}$$

$$\downarrow$$

$$m_{01}^{l} = [20 - 50, 170 - 160] = [-30, 10]$$

Let's assume our "learned" message function is just:

$$f_e(n_i^k, n_j^k) = n_i^k - n_j^k$$

and our "learned" node update function is just:

$$f_n(a_i^k, n_i^k) = \frac{1}{2}(a_i^k + n_i^k)$$



1. Calculate messages

$$m_{ij}^l = f_e\big(n_i^k, n_j^k\big)$$

Let's assume our "learned" message function is just:

$$f_e(n_i^k, n_j^k) = n_i^k - n_j^k$$

and our "learned" node update function is just:

$$f_n(a_i^k, n_i^k) = \frac{1}{2}(a_i^k + n_i^k)$$



1. Calculate messages

$$m_{ij}^l = f_e\big(n_i^k, n_j^k\big)$$

$$\begin{split} E_{ij} = \begin{bmatrix} 0 & 0 & 0 & 2 & 1 & 2 & 3 & 3 \\ 1 & 2 & 3 & 3 & 0 & 0 & 0 & 2 \end{bmatrix} \\ \downarrow & \downarrow \\ m_{ij}^{l} = [\begin{bmatrix} -30 \\ 10 \end{bmatrix}, \begin{bmatrix} -15 \\ 20 \end{bmatrix}, \begin{bmatrix} -55 \\ -20 \end{bmatrix}, \begin{bmatrix} -40 \\ -40 \end{bmatrix}, \begin{bmatrix} 30 \\ -10 \end{bmatrix}, \begin{bmatrix} 15 \\ 20 \end{bmatrix}, \begin{bmatrix} 55 \\ 20 \end{bmatrix}, \begin{bmatrix} 40 \\ 40 \end{bmatrix}] \end{split}$$

2. Aggregate messages

$$a_i^l = \sum_j (m_{ij}^l)$$

$$a_{0}^{l} = m_{01}^{l} + m_{02}^{l} + m_{03}^{l}$$

= $\begin{bmatrix} -30\\ 10 \end{bmatrix} + \begin{bmatrix} -15\\ 20 \end{bmatrix} + \begin{bmatrix} -55\\ -20 \end{bmatrix}$
= $\begin{bmatrix} -100\\ 10 \end{bmatrix}$

Let's assume our "learned" message function is just:

$$f_e(n_i^k, n_j^k) = n_i^k - n_j^k$$

and our "learned" node update function is just:

 $f_n(a_i^k, n_i^k) = \frac{1}{2}(a_i^k + n_i^k)$



1. Calculate messages

$$m_{ij}^l = f_e \big(n_i^k, n_j^k \big)$$

$$\begin{split} E_{ij} &= \begin{bmatrix} 0 & 0 & 0 & 2 & 1 & 2 & 3 & 3 \\ 1 & 2 & 3 & 3 & 0 & 0 & 0 & 2 \end{bmatrix} \\ & \downarrow \\ m_{ij}^{l} &= [\begin{bmatrix} -30 \\ 10 \end{bmatrix}, \begin{bmatrix} -15 \\ 20 \end{bmatrix}, \begin{bmatrix} -55 \\ -20 \end{bmatrix}, \begin{bmatrix} -40 \\ -40 \end{bmatrix}, \begin{bmatrix} 30 \\ -10 \end{bmatrix}, \begin{bmatrix} 15 \\ 20 \end{bmatrix}, \begin{bmatrix} 55 \\ 20 \end{bmatrix}, \begin{bmatrix} 40 \\ 40 \end{bmatrix}] \end{split}$$

2. Aggregate messages

$$a_i^l = \sum_j (m_{ij}^l)$$

 $a_i^l = \begin{bmatrix} -100\\10 \end{bmatrix}, \begin{bmatrix} 30\\-10 \end{bmatrix}, \begin{bmatrix} -25\\-20 \end{bmatrix}, \begin{bmatrix} 95\\60 \end{bmatrix}]$

Let's assume our "learned" message function is just:

$$f_e(n_i^k, n_j^k) = n_i^k - n_j^k$$

and our "learned" node update function is just:

$$f_n(a_i^k, n_i^k) = \frac{1}{2}(a_i^k + n_i^k)$$



1. Calculate messages

$$m_{ij}^l = f_e\big(n_i^k, n_j^k\big)$$

$$\begin{split} E_{ij} = \begin{bmatrix} 0 & 0 & 0 & 2 & 1 & 2 & 3 & 3 \\ 1 & 2 & 3 & 3 & 0 & 0 & 0 & 2 \end{bmatrix} \\ \downarrow & \downarrow \\ m_{ij}^{l} = {}_{[10]}^{-30}, {}_{[20]}^{-15}, {}_{[20]}^{-55}, {}_{[40]}^{-40}, {}_{[40]}^{-30}, {}_{[20]}^{15}, {}_{[20]}^{55}, {}_{[40]}^{40}]_{]} \end{split}$$

2. Aggregate messages $a_{i}^{l} = \sum_{j} (m_{ij}^{l})$ $a_{i}^{l} = [\begin{bmatrix} -100\\10 \end{bmatrix}, \begin{bmatrix} 30\\-10 \end{bmatrix}, \begin{bmatrix} -25\\-20 \end{bmatrix}, \begin{bmatrix} 95\\60 \end{bmatrix}]$

3. Update $n'_{i}^{k} = f_{n}(a_{i}^{l}, n_{i}^{k})$ nodes

$$n_{0}^{\prime k} = f_{n}(a_{0}^{l}, n_{0}^{k}) = \frac{1}{2} \left(\begin{bmatrix} -100\\10 \end{bmatrix} + \begin{bmatrix} 20\\170 \end{bmatrix} \right) = \begin{bmatrix} -40\\90 \end{bmatrix}$$

Let's assume our "learned" message function is just: $f_e(n_i^k, n_j^k) = n_i^k - n_j^k$

and our "learned" node update function is just:

 $f_n(a_i^k, n_i^k) = \frac{1}{2}(a_i^k + n_i^k)$



1. Calculate $m_{i\,i}^l = f_e(n_i^k, n_i^k)$ messages $E_{ij} = \begin{bmatrix} 0 & 0 & 0 & 2 & 1 & 2 & 3 & 3 \\ 1 & 2 & 3 & 3 & 0 & 0 & 0 & 2 \end{bmatrix}$ $m_{ij}^{l} = [\begin{smallmatrix} -30 \\ 10 \end{smallmatrix}], \begin{bmatrix} -15 \\ 20 \end{smallmatrix}], \begin{bmatrix} -55 \\ -20 \end{smallmatrix}], \begin{bmatrix} -40 \\ -40 \end{smallmatrix}], \begin{bmatrix} 30 \\ -10 \end{smallmatrix}], \begin{bmatrix} 15 \\ 20 \end{smallmatrix}], \begin{bmatrix} 55 \\ 20 \end{smallmatrix}], \begin{bmatrix} 40 \\ 40 \end{smallmatrix}]$ 2. Aggregate $a_i^l = \sum_i (m_{ij}^l)$ messages $a_i^l = \begin{bmatrix} -100 \\ 10 \end{bmatrix}, \begin{bmatrix} 30 \\ -10 \end{bmatrix}, \begin{bmatrix} -25 \\ -20 \end{bmatrix}, \begin{bmatrix} 95 \\ 60 \end{bmatrix}$ $n_i^{\prime k} = f_n(a_i^l, n_i^k)$ 3. Update nodes $n_{i}^{\prime k} = \begin{bmatrix} -40 & 90\\ 40 & 75\\ 5 & 65 \end{bmatrix}$ 125

GNN tasks vs. architectures

- Just like with the CNN, the choice of GNN convolution is usually separate from the final training task
- Any GNN convolution that does message passing, and updates hidden node features, allows us to predict:



The Transformer as a GNN



8/05/2024

Graph Neural Networks – Daniel Murnane

• For each node, attach three sets of hidden features: K, Q, V





8/05/2024

- For each node, attach three sets of hidden features: K, Q, V
- Define a message function on each edge:

$$m_{ij} = f_e(K_i^k, Q_j^k) = softmax\left(\sum_k K_i^k Q_j^k\right)$$

8/05/2024





- For each node, attach three sets of hidden features: K, Q, V
- Define a message function on each edge:

$$m_{ij} = f_e(K_i^k, Q_j^k) = softmax\left(\sum_k K_i^k Q_j^k\right)$$

- Define an aggregation function around each node:
- $a_i^k = \sum_j m_{ij} V_j^k$, this is just a weighted sum





- For each node, attach three sets of hidden features: K, Q, V
- Define a message function on each edge:

$$m_{ij} = f_e(K_i^k, Q_j^k) = softmax\left(\sum_k K_i^k Q_j^k\right)$$

- Define an aggregation function around each node:
- $a_i^k = \sum_j m_{ij} V_j^k$, this is just a weighted sum
- The node update is just FFNNs to get the next K', Q', V'





COPENHAGEN

- For each node, attach three sets of hidden features: K, Q, V
- Define a message function on each edge:

 $m_{ij} = f_e(K_i^k, Q_i^k) = softmax \sum_{i=1}^{k} K_i^k Q_i^k$ The mean aggregation runction around each mode: $a_i^k = \sum_j m_{ij} V_i^k$, this is just a weighted sum

• The node update is just FFNNs to get the next K', Q', V'





https://arxiv.org/pdf/2012.09699.pdf

Case Studies



8/05/2024

Google Maps

- DeepMind worked with Google Maps developers to apply a GNN to estimate travel time
- Improved estimates by up to 50% in large cities
- Maps are *really* hard to optimize on, as they exhibit combinatorial behavior
- GNNs can give fast, approximate solutions to map problems



Edge Classification for Particle Tracking

- In the ATLAS experiment, we need to connect points in the detector to reconstruct particle tracks
- It's a massive connect-the-dots game, with 300,000 dots to connect, every 25 nanoseconds





Edge Classification for Particle Tracking



- In the ATLAS experiment, we need to connect points in the detector to reconstruct particle tracks
- It's a massive connect-the-dots game, with 300,000 dots to connect, every 25 nanoseconds
- We are building graph neural networks to do this faster than traditional algorithms
 ^{8/05/2024}



Edge Classification for Particle Tracking



- In the ATLAS experior properties and a long the detector to reconstruct particle tracks
- It's a massive connect-the-dots game, with 300,000 dots to connect, every 25 nanoseconds
- We are building graph neural networks to do this faster than traditional algorithms

10⁵ 10⁴ 10² 10²1

Ice Cube neutrino prediction



- Neutrino observatory inside the ice of Antarctica
- Goal is to detect neutrinos and measure precisely the direction they came from, their energy, what type of neutrino they are, etc.

Ice Cube neutrino prediction



Ice Cube neutrino prediction

