

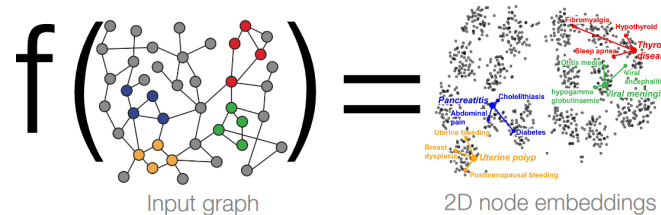
Overview

1. Basics of deep learning
2. Deep learning for graphs
3. Graph Convolutional Networks
4. GNNs subsume CNNs and Transformers

Recap

Node embedding

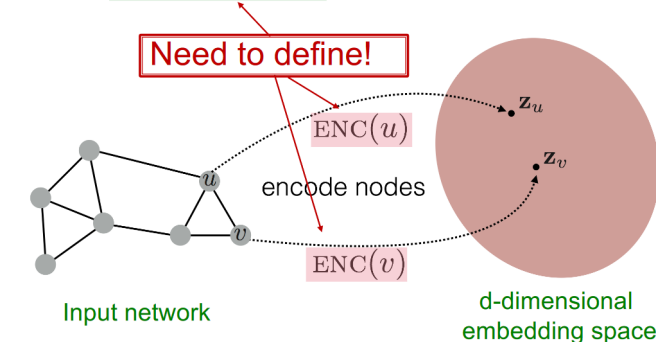
Intuition. Map nodes to d -dimensional embeddings such that similar nodes in the graph are embedded close together.



Question. How to learn mapping function f ?

Goal.

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



Two Key Components. Encoder & Similarity function (decoder)

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

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

d -dimensional embedding

node in the input graph

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

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

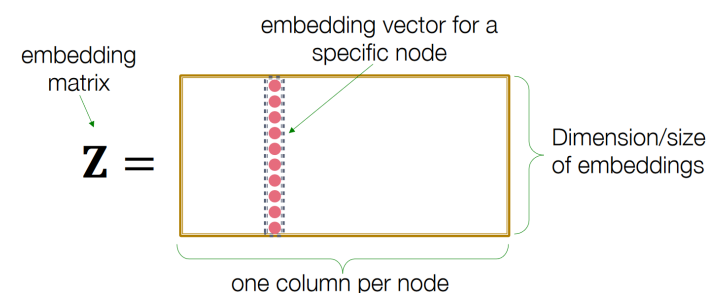
Similarity of u and v in the original network

dot product between node embeddings

Decoder

“Shallow” Encoding. Encoder is just an embedding-lookup.

- shallow since just memorising embedding of every node



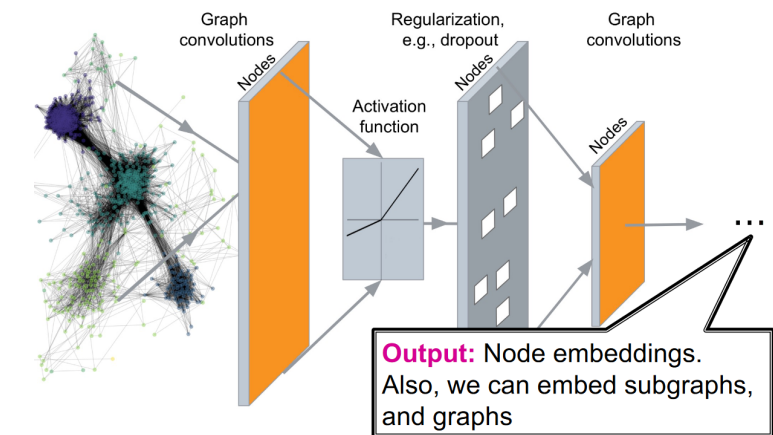
Shallow Encoders.

- $\mathcal{O}(|V|)$ parameters are needed (each column is the embedding for a node)
 - No sharing of parameters between nodes
- Inherently “transductive”
 - Cannot generate embeddings for nodes that are **not seen** during training
- Do not incorporate **node features**

GNN, Deep Graph Encoders

- $\text{ENC}(v)$ = multiple layers of non-linear transformations based on graph structure
 - can be **combined** with node similarity functions

End to end. labels on the right to the graph structure on the left

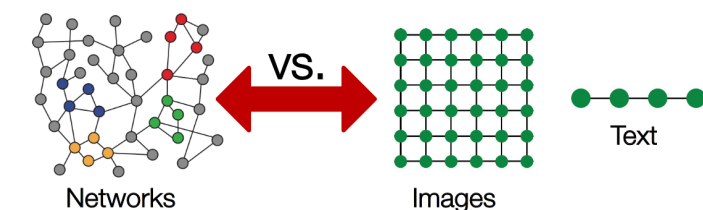


Tasks

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

Graph vs Image

- **Arbitrary size**
- complex **topological structure** (i.e., no spatial locality like grids)
- No fixed node ordering or reference point
- Often dynamic and have **multimodal features**



Basics of Deep Learning

Basics of Deep Learning

See slides.

- Iteration: 1 step of SGD on a minibatch (1 step 1 ~ 1 mini-match)
- Epoch: one full pass over the dataset (# iterations is equal to **ratio** of dataset size and batch size)
- SGD is **unbiased** estimator of full gradient

Objective Function

- **Objective:** $\min_{\Theta} \mathcal{L}(\mathbf{y}, f(\mathbf{x}))$
- In deep learning, function f can be very complex
- **Example:**
 - To start simple, consider linear function $f(\mathbf{x}) = \mathbf{W} \cdot \mathbf{x}$, $\Theta = \{\mathbf{W}\}$
 - Then, if f returns a scalar, then \mathbf{W} is a learnable **vector**

$$\nabla_{\mathbf{W}} f = \left(\frac{\partial f}{\partial w_1}, \frac{\partial f}{\partial w_2}, \frac{\partial f}{\partial w_3} \dots \right)$$

- But, if f returns a vector, then \mathbf{W} is the **weight matrix**

$$\nabla_{\mathbf{W}} f = \begin{bmatrix} \frac{\partial f_1}{\partial w_{11}} & \frac{\partial f_1}{\partial w_{12}} \\ \frac{\partial f_2}{\partial w_{21}} & \frac{\partial f_2}{\partial w_{22}} \end{bmatrix}$$

Jacobian
matrix of f

More complex – use chain rule

- **How about a more complex function:**

$$f(\mathbf{x}) = \mathbf{W}_2(\mathbf{W}_1 \mathbf{x}), \Theta = \{\mathbf{W}_1, \mathbf{W}_2\}$$

- Recall **chain rule:**

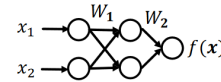
$$\frac{df}{dx} = \frac{dg}{dh} \cdot \frac{dh}{dx} \text{ or } f'(x) = g'(h(x))h'(x)$$

- **Example:** $\nabla_{\mathbf{x}} f = \frac{\partial f}{\partial (W_1 \mathbf{x})} \cdot \frac{\partial (W_1 \mathbf{x})}{\partial \mathbf{x}}$

- **Back-propagation:** Use of **chain rule** to propagate gradients of intermediate steps, and finally obtain gradient of \mathcal{L} w.r.t. Θ .

Example: Simple 2-layer linear network

$$f(\mathbf{x}) = g(h(\mathbf{x})) = \mathbf{W}_2(\mathbf{W}_1 \mathbf{x})$$



$$\mathcal{L} = \sum_{(\mathbf{x}, \mathbf{y}) \in \mathcal{B}} \left\| (\mathbf{y}, -f(\mathbf{x})) \right\|_2$$

- The loss \mathcal{L} sums the L2 loss in a minibatch \mathcal{B} .

- **Hidden layer:**

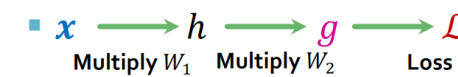
- Intermediate representation of input \mathbf{x}
- Here we use $h(\mathbf{x}) = \mathbf{W}_1 \mathbf{x}$ to denote the hidden layer
- $f(\mathbf{x}) = \mathbf{W}_2 h(\mathbf{x})$

Backprop

- The w.r.t. is backward i.e. the opposite direction to the forward pass

- **Forward propagation:**

Compute loss starting from input

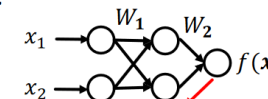


Remember:

$$f(\mathbf{x}) = \mathbf{W}_2(\mathbf{W}_1 \mathbf{x})$$

$$h(\mathbf{x}) = \mathbf{W}_1 \mathbf{x}$$

$$g(z) = \mathbf{W}_2 z$$



- **Back-propagation to compute gradient of**

$$\Theta = \{\mathbf{W}_1, \mathbf{W}_2\}$$

Start from loss, compute the gradient

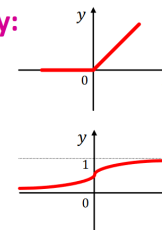
$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}_2} = \frac{\partial \mathcal{L}}{\partial f} \cdot \frac{\partial f}{\partial \mathbf{W}_2}, \quad \frac{\partial \mathcal{L}}{\partial \mathbf{W}_1} = \frac{\partial \mathcal{L}}{\partial f} \cdot \frac{\partial f}{\partial \mathbf{W}_2} \cdot \frac{\partial \mathbf{W}_2}{\partial \mathbf{W}_1}$$

Compute backwards

Compute backwards

Non-linearity

- Improve the expressiveness of the model
- Note that in $f(\mathbf{x}) = \mathbf{W}_2(\mathbf{W}_1 \mathbf{x})$, $\mathbf{W}_2 \mathbf{W}_1$ is another matrix (vector, if we do binary classification)
 - Hence $f(\mathbf{x})$ is still linear w.r.t. \mathbf{x} no matter how many weight matrices we compose
- **We introduce non-linearity:**
 - **Rectified linear unit (ReLU)**
 $\text{ReLU}(x) = \max(x, 0)$
 - **Sigmoid**
 $\sigma(x) = \frac{1}{1 + e^{-x}}$

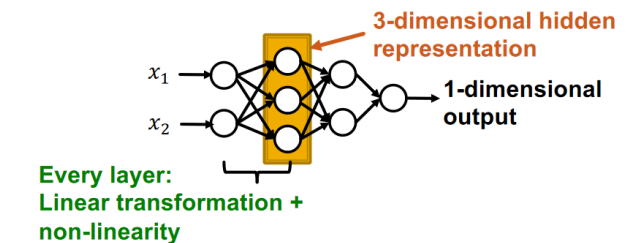


MLP – combine linear transformation & non-linearity

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

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

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



Summary

- **Objective function:**

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

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