

Any reproduction or distribution of this document, in whole or in part, is prohibited unless permission is granted by the authors

EE-559

Deep Learning

What's on today?

- **Graphs**: on nodes, edges and structure
- **Simple graph**: on aggregation and parameter sharing
- **Tasks on graphs**: how to perform regression and classification
- **Graph convolutional networks**: on deep learning with graphs
- **Graph attention**: on weighted, learned, neighbor feature aggregation
- **Training**: how to deal with the structure
- **Line graphs**: on the complementary graph
- **Graph types**: on the diversity of graph representations
- **Exercises**: message passing and graph classification

Graphs

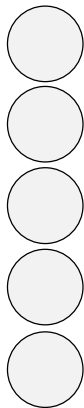
$$y = f(x)$$

$$y = f(x; \Theta)$$

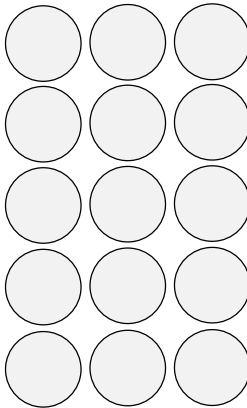
Input



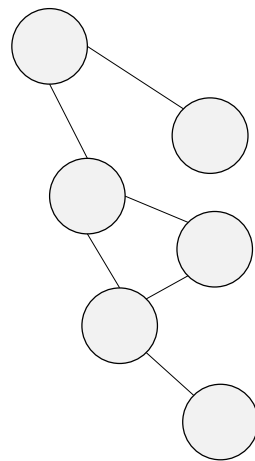
$$x \in \mathbb{R}$$



$$\mathbf{x} \in \mathbb{R}^W$$



$$\mathbf{X} \in \mathbb{R}^{W \times N}$$



$$\mathbf{X} \in \mathbb{R}^{W \times N}$$

$$\mathbf{A} \in \mathbb{R}^{N \times N}$$

slido



What type of real-world problems can be modeled effectively using a graph representation?

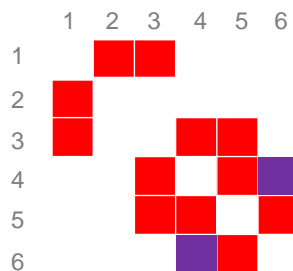
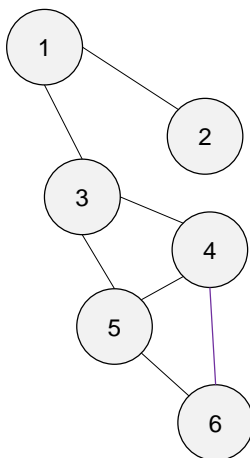
① Start presenting to display the poll results on this slide.

$$y = f(\mathbf{X}, \mathbf{A}; \boldsymbol{\theta})$$

$$y = f(\mathbf{X}, \mathbf{A}; \boldsymbol{\theta})$$

adjacency
matrix

Adjacency matrix



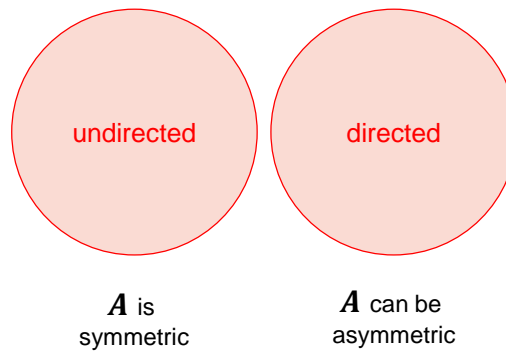
$$\mathbf{A} \in \mathbb{R}^{N \times N}$$

$$N = 6$$

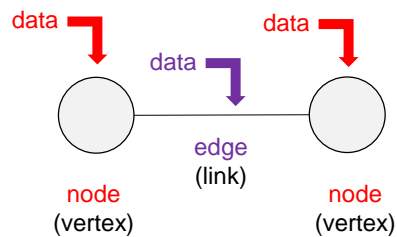
Concepts:

Node indexing, walks of length one

Edges



Node and edge embeddings



Examples: Social networks, citation networks, train map (stations, lines), protein interactions in a cell, molecule (component, bound)

Structure and embeddings

structure of the graph

$$\mathbf{A} \in \mathbb{R}^{N \times N}$$

$$a_{nm} \in \{0,1\}$$

node embeddings

$$\mathbf{X} \in \mathbb{R}^{W_x \times N}$$

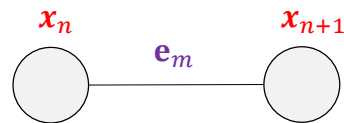
$$\mathbf{x}_n \in \mathbb{R}^{W_x}$$

edge embeddings

$$\mathbf{E} \in \mathbb{R}^{W_e \times E}$$

$$\mathbf{e}_m \in \mathbb{R}^{W_e}$$

N number of nodes
 E number of edges
 W_x size of the node embeddings
 W_e size of the edge embeddings

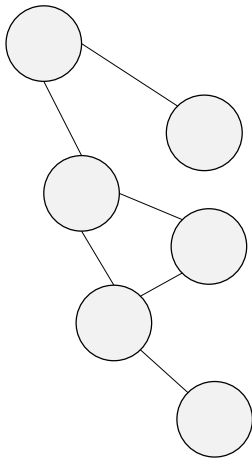


Concept:

The adjacency matrix, \mathbf{A} , is symmetric for undirected graphs

Simple graph

Simple graph



$A \in \mathbb{R}^{N \times N}$ defines neighboring edges
(**structure of the graph**)

\mathcal{N}_n set of all neighbors of node n

$X \in \mathbb{R}^{W_x \times N}$ node embeddings
 $\mathbf{x}_n \in \mathbb{R}^{W_x}$

$E \in \mathbb{R}^{W_e \times E}$ edge embeddings
 $\mathbf{e}_m \in \mathbb{R}^{W_e}$

Concepts:

At most one edge between any two nodes, undirected edges, no self-edges

Walks from a node

$$\mathbf{A} \in \mathbb{R}^{N \times N}$$

$$a_{nm} \in \{0,1\}$$

$$\mathbf{x}_n \in \{0,1\}^N$$

encoding each node
as a one-hot vector

$\mathbf{A}\mathbf{x}_n$ number of walks of length 1
from node n to each other node

$\mathbf{A}^2\mathbf{x}_n$ number of walks of length 2
from node n to each other node

⋮

$\mathbf{A}^L\mathbf{x}_n$ number of walks of length L
from node n to each other node

Concept:

Neighbourhood of each node

Permutation of node indices

P permutation matrix

*one entry in each row and column
is 1, the others are 0*

$$X' = XP$$

permutes the columns

$$A' = P^T A P$$

permutes the rows

Concept:

Node indexing in a graph is *arbitrary* (any processing applied to the graph should be indifferent to permutations)

Importance of parameter sharing

invariance

dependence only
on the structure, A ,
*not on the labelling
of the nodes*

equivariance

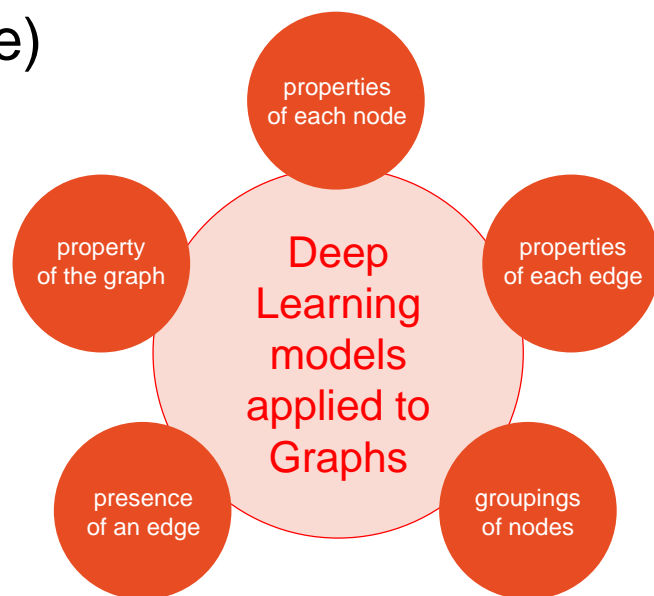
output permuted
consistently
with permutation
of A

scaling

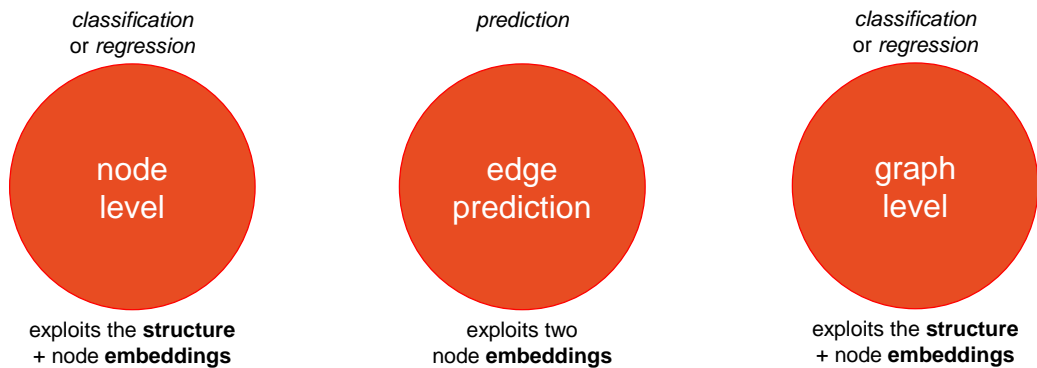
handling the
growing
of the
graph size

Tasks on graphs

Tasks (inference)



Tasks on a graph



Tasks on a graph

X (input) node embeddings (data)
 A adjacency matrix (structure)
 H^k hidden representation (layer k)
 h_n^k (modified) node n embedding

$$P(y_{nm} = 1 | X, A) = \text{sigmoid}((h_m^K)^T h_n^K)$$

edge prediction

$$P(y_n = 1 | X, A) = \text{sigmoid}(\Theta_1^K + \Theta_1^K h_n^K)$$

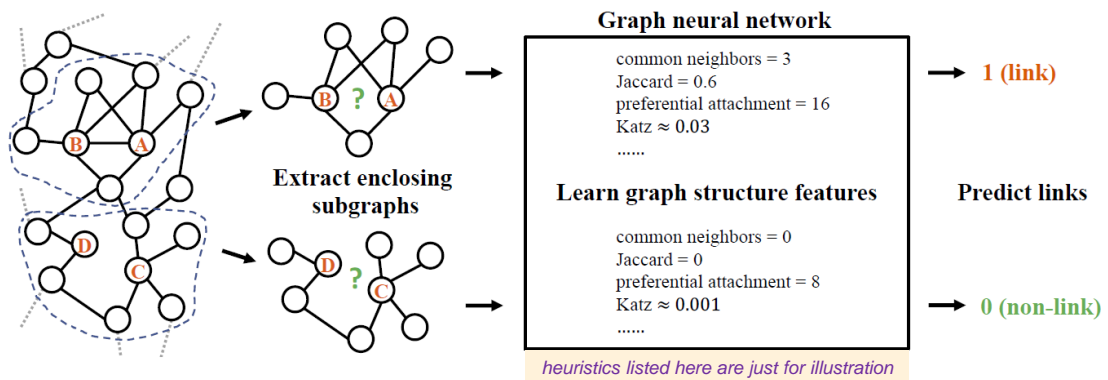
node level

$$P(y = 1 | X, A) = \text{sigmoid}(\Theta_1^K + \Theta_1^K H^K \underset{\text{mean pooling}}{1/N})$$

graph level

Concepts: *Edge prediction*: binary classification task, *node level*: independently for each node, *graph level*: combining output node embeddings

Edge prediction example



[arXiv:1802.09691](https://arxiv.org/abs/1802.09691)

Concepts: Graph structure features inside the observed node and edge structures, relative degree of influence of a node (*Katz centrality*), intersection over the union of the sets of neighbors of a node (*Jaccard index*)

slido



Is there anything about the mini-project that you would like clarification on?

① Start presenting to display the poll results on this slide.

Graph convolutional networks

Graph convolutional networks

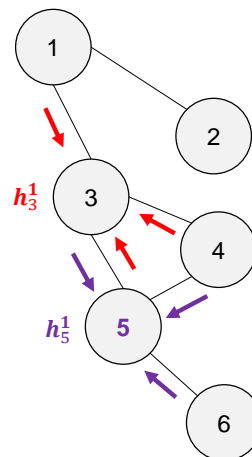
$$\mathbf{H}^1 = f(\mathbf{X}, \mathbf{A}, \boldsymbol{\theta}^0)$$

$$\mathbf{H}^2 = f(\mathbf{H}^1, \mathbf{A}, \boldsymbol{\theta}^1)$$

⋮

$$\mathbf{H}^K = f(\mathbf{H}^{K-1}, \mathbf{A}, \boldsymbol{\theta}^{K-1})$$

\mathbf{X} input node embeddings
 \mathbf{A} adjacency matrix (structure)
 \mathbf{H}^k modified node embeddings (layer k)
 $\boldsymbol{\theta}^k$ parameters that map from layer k to $k + 1$



Concepts: Relational inductive bias, *message passing*, spatial-based convolutional graph neural network (GCN), updating the *local representation* at each node by gathering information from its neighbors by passing messages

GCN layer

$$\mathcal{G}_n^k = \sum_{m \in \mathcal{N}_n} \mathbf{h}_m^k$$

set of indices of the
neighborhood of node n

$$\mathbf{h}_n^{k+1} = a[\Theta_0^k + \Theta_1^k \mathbf{h}_n^k + \Theta_1^k \mathcal{G}_n^k]$$

$$\mathbf{H}^{k+1} = a[\Theta_0^k \mathbf{1}_N^T + \Theta_1^k \mathbf{H}^k + \Theta_1^k \mathbf{H}^k \mathbf{A}]$$

aggregated neighborhood

$$= a[\Theta_0^k \mathbf{1}_N^T + \Theta_1^k \mathbf{H}^k (\mathbf{A} + \mathbf{I})]$$

combine “messages” from adjacent nodes
(sum them with the transformed current node)

Concept: Aggregating information from neighboring nodes (*sum of node embeddings*),
local function of the embedding of the previous layer, combine *messages* from adjacent nodes

Oversmoothing problem

repeated graph convolutions make
node embeddings indistinguishable
(oversmoothing)



$k \uparrow \Rightarrow$ performance gradually decreases



normalization layer

k : number of layers (network depth)

[arXiv:1909.12223](https://arxiv.org/abs/1909.12223)

Graph prediction example

Is a molecule toxic?

$$\mathbf{H}^k = a[\Theta_0^{k-1} \mathbf{1}_N^T + \Theta_1^{k-1} \mathbf{H}^{k-1} (\mathbf{A} + \mathbf{I})]$$

$$f(\mathbf{X}, \mathbf{A}, \Theta) = \text{sigmoid}(\Theta_1^K + \Theta_1^K \mathbf{H}^K \mathbf{1}/N)$$

mean
pooling

$\mathbf{A} \in \mathbb{R}^{N \times N}$
molecular structure

$\mathbf{X} \in \mathbb{R}^{118 \times N}$
matrix of one-hot vectors
indicating which element
of the **periodic table** is present

Variants for aggregation

$$\mathbf{H}^{k+1} = a[\Theta_0^k \mathbf{1}_N^T + \Theta_1^k \mathbf{H}^k (\mathbf{A} + \mathbf{I})]$$

$$\mathbf{H}^{k+1} = a[\Theta_0^k \mathbf{1}_N^T + \Theta_1^k \mathbf{H}^k (\mathbf{A} + (1 + \epsilon_k) \mathbf{I})]$$

learned
scalar

diagonal enhancement

$$\mathbf{H}^{k+1} = \left[a[\Theta_0^k \mathbf{1}_N^T + \Theta_1^k \mathbf{H}^k \mathbf{A}] \right]$$

\mathbf{H}^k

residual connection
concatenation with the node

$$\mathcal{G}_n^k = \max_{m \in \mathcal{N}_n} \mathbf{h}_m^k$$

max pooling aggregation
element-wise maximum

Mean aggregation

$$\mathcal{G}_n^k = \sum_{m \in \mathcal{N}_n} \mathbf{h}_m^k \quad \Rightarrow \quad \mathcal{G}_n^k = \frac{1}{|\mathcal{N}_n|} \sum_{m \in \mathcal{N}_n} \mathbf{h}_m^k$$

mean aggregation

when *embedding information*
is more important than
structural information

$$\mathbf{D} \in \mathbb{R}^{N \times N}$$

diagonal matrix

each non-zero element
is the number of neighbors
of the corresponding node



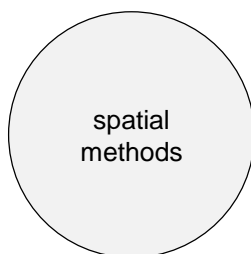
$$(\mathbf{D})^{-1}$$

inverse matrix

each non-zero element
is the denominator to
compute the average

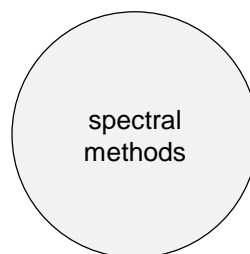
$$\mathbf{H}^{k+1} = a[\Theta_0^k \mathbf{1}_N^T + \Theta_1^k \mathbf{H}^k (\mathbf{A}(\mathbf{D})^{-1} + \mathbf{I})]$$

Spatial and spectral methods



spatial
methods

use the
graph structure



spectral
methods

use the
Fourier domain

Graphs and transformers



handle
variable-length
input

transformer
w/o positional
embedding
~
GNN on a fully
connected
graph

use of
attention

w/o: without

~: equivalent to

GNN: Graph Neural Network

Graph attention

Graph attention layers

$$\widehat{\mathbf{H}}^k = \Theta_0^k \mathbf{1}_N^T + \Theta_1^k \mathbf{H}^k \quad \text{linear transformation on node embedding}$$

$$s_{mn} = a \left[(\underbrace{\boldsymbol{\phi}^k}_{\text{vector of learned parameters}})^T \begin{bmatrix} \widehat{\mathbf{h}}_m^k \\ \widehat{\mathbf{h}}_n^k \end{bmatrix} \right] \quad \mathbf{S} \in \mathbb{R}^{N \times N} \quad \text{similarity of every node to every other node}$$

$$\mathbf{H}^{k+1} = a \left[\widehat{\mathbf{H}}^k \text{softmax}[\mathbf{S}, \mathbf{A} + \mathbf{I}] \right] \quad \text{attention weights applied to transformed embedding}$$

softmax for each column of \mathbf{S}
(set to zero non-neighboring nodes)

Concepts: Aggregation by attention, concatenation of node embeddings, weights depend on the *data at the nodes* (previous cases depended instead on neighbors equally or on graph topology)

Aggregation: summary

$$\mathbf{H}^{k+1} = a[\Theta_0^k \mathbf{1}_N^T + \Theta_1^k \mathbf{H}^k (\mathbf{A} + \mathbf{I})]$$

$$\mathbf{H}^{k+1} = a[\Theta_0^k \mathbf{1}_N^T + \Theta_1^k \mathbf{H}^k (\mathbf{A} + (1 + \epsilon_k) \mathbf{I})] \quad \begin{array}{l} \text{diagonal enhancement} \\ \text{learned scalar } \epsilon_k \end{array}$$

$$\mathbf{H}^{k+1} = \left[a[\Theta_0^k \mathbf{1}_N^T + \Theta_1^k \mathbf{H}^k \mathbf{A}] \right] \quad \begin{array}{l} \text{residual connection} \\ \text{concatenation with the node} \end{array}$$

$$\mathcal{G}_n^k = \max_{m \in \mathcal{N}_n} \mathbf{h}_m^k \quad \begin{array}{l} \text{max pooling aggregation} \\ \text{element-wise maximum of the } \mathbf{h}_m^k \end{array}$$

$$\mathbf{H}^{k+1} = a[\Theta_0^k \mathbf{1}_N^T + \Theta_1^k \mathbf{H}^k (\mathbf{A}(\mathbf{D})^{-1} + \mathbf{I})] \quad \begin{array}{l} \text{mean aggregation} \\ \text{less importance to structural information} \end{array}$$

$$\mathbf{H}^{k+1} = a[\widehat{\mathbf{H}}^k \text{softmax}[\mathbf{S}, \mathbf{A} + \mathbf{I}]] \quad \begin{array}{l} \text{aggregation by attention} \\ \text{attention weights on transformed embedding} \end{array}$$

Training

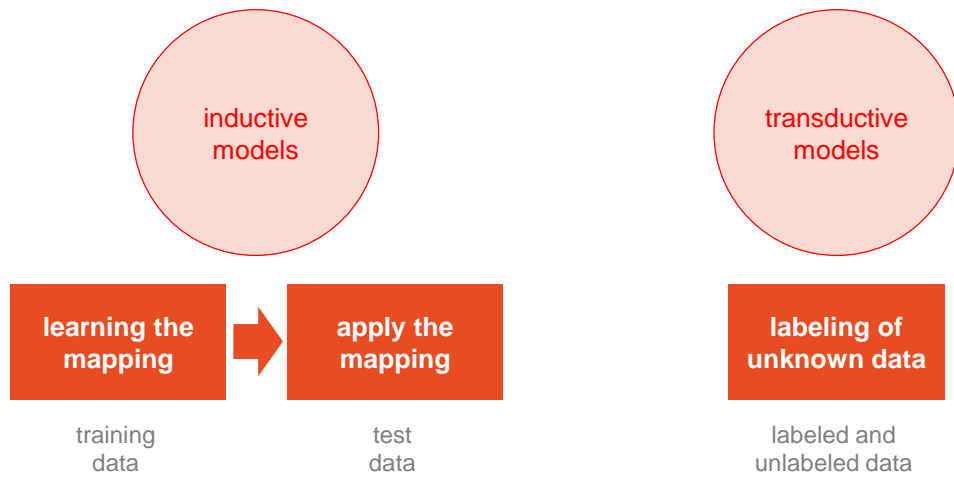
$$\{x_i, y_i\}_{i=1}^M$$

$$\{\mathbf{X}_i, \mathbf{A}_i, y_i\}_{i=1}^M$$

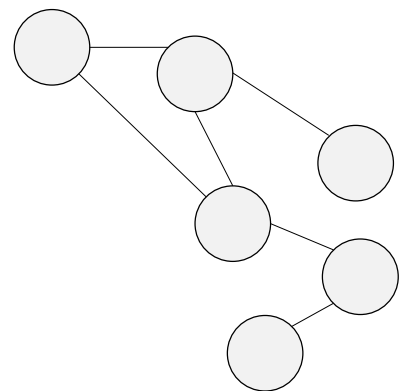
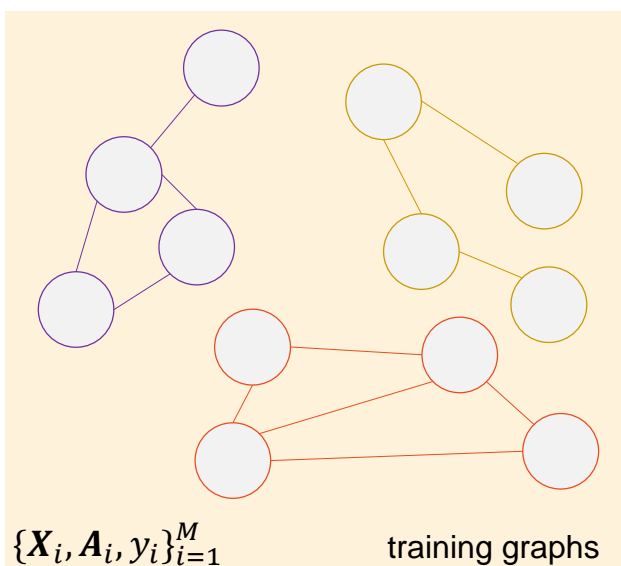
$$y = f(\mathbf{X}, \mathbf{A}; \boldsymbol{\theta})$$

learned with SGD
& binary cross-entropy loss

Inductive and transductive models

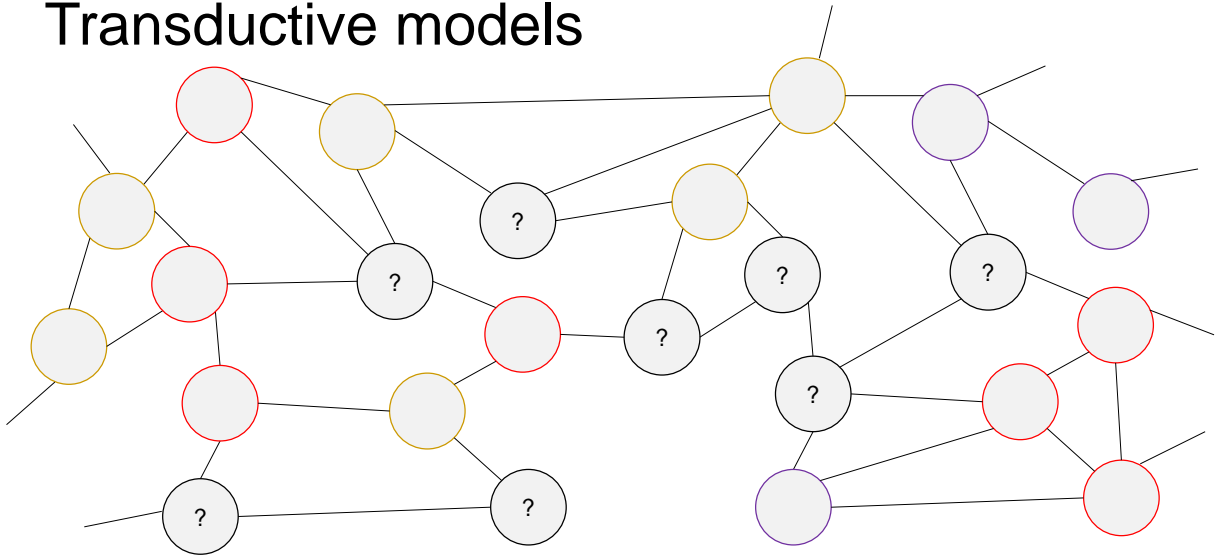


Inductive models



testing graph

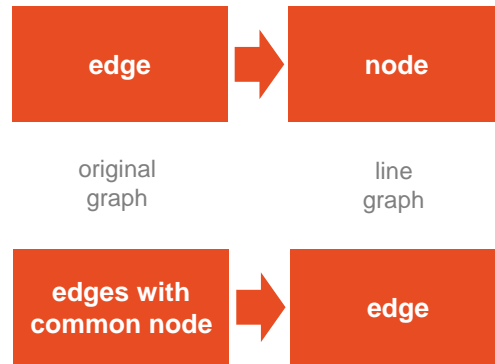
Transductive models

**Concept:**

Train to predict the known labels, then examine the predictions at the unknown nodes

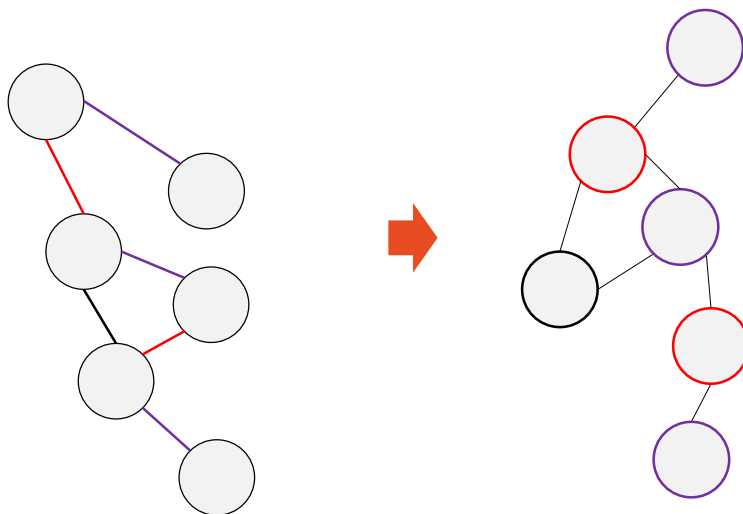
Line graphs

Line or edge graph

**Concept:**

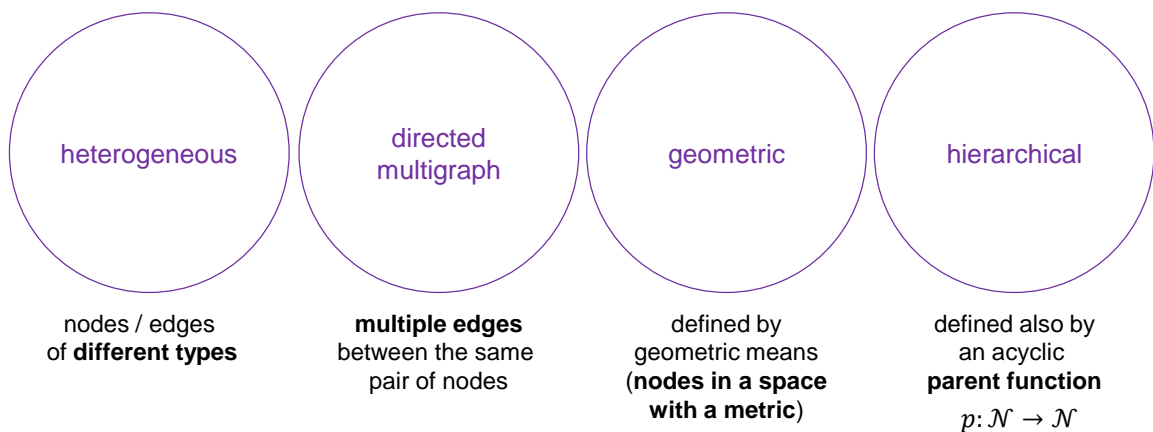
Edge graph, complementary graph to process *edge embeddings*

Example



Graph types

Graph types



Concepts:

Simple graphs, spatial elements associated to geometric objects, parent function defines the hierarchy

Exercises

Today's exercises

Practice. You will become familiar with:

- **PyTorch geometric:** load datasets, view information about the graphs
- **Graph Neural Networks:** message passing and graph classification

What did we learn today?

- Graphs
- Simple graph
- Tasks on graphs
- Graph convolutional networks
- Graph attention
- Training
- Line graphs
- Graph types
- Exercises

EE-559

Deep Learning

andrea.cavallaro@epfl.ch