

# Graph neural networks lecture

TIF 360

Mats Granath

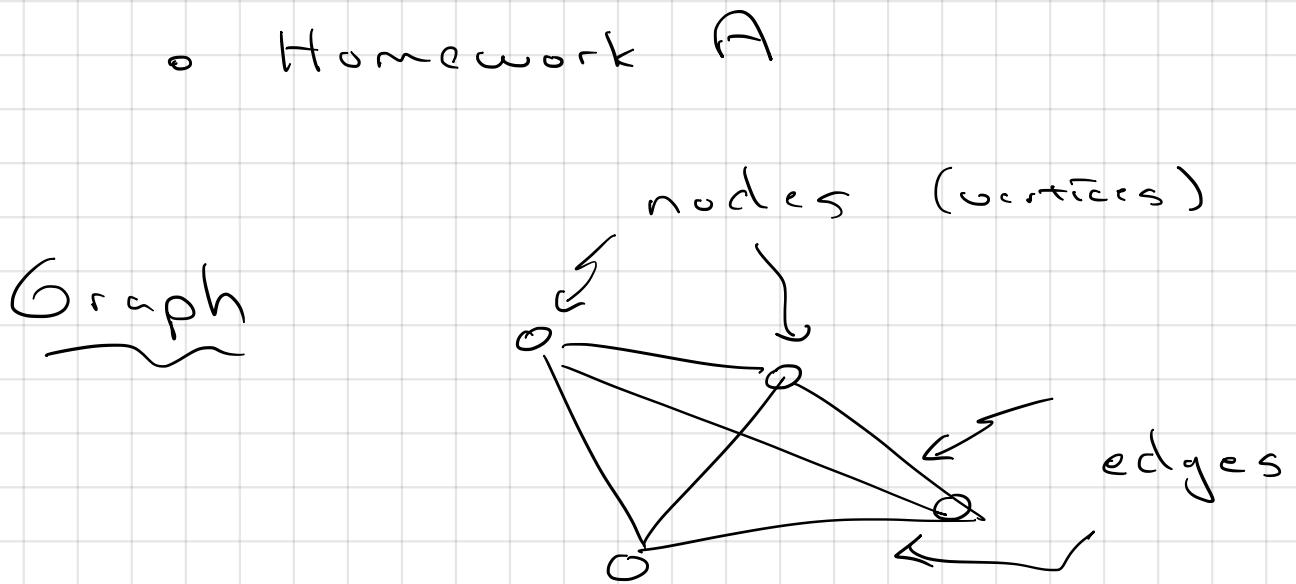
28/3 2023

IV

Graph neural networks =  
= neural networks acting on graphs

## Outline

- Why interesting to apply NN to graph objects?
- How is a graph defined?
- What type of machine learning task may we want to do?
- Basic GNN layer:  
Graph convolution
- Other layers + Graph Pooling
- Homework A

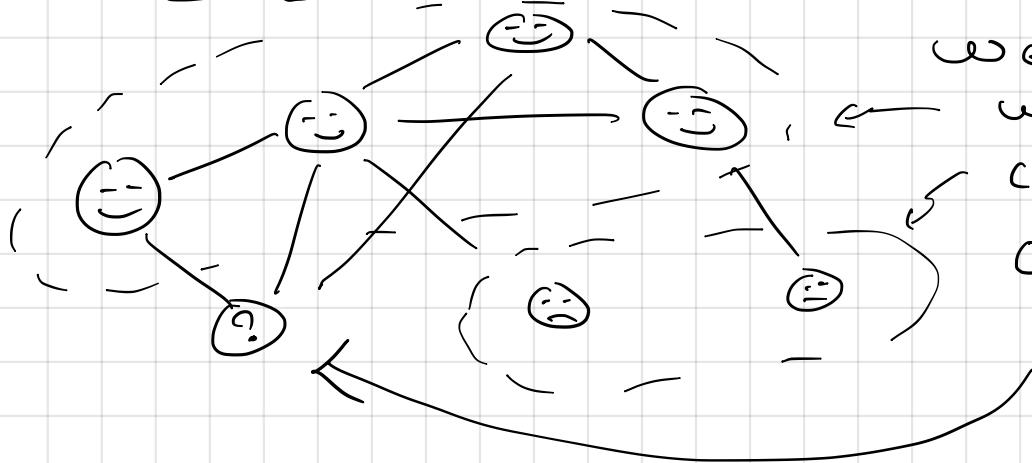


More specific shortly

# Why NN on graphs / what tasks?

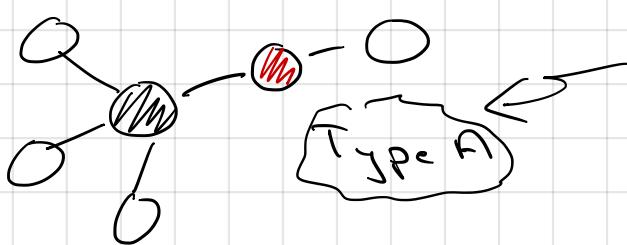
Many interesting forms of data that can be nicely represented as graphs:

## Social networks:



We may want to do clustering.  
Or node prediction.

## Molecules



We may want to do classification

or regression

e.g. how chemically active is the molecule

# Different methodology for different tasks

## 1) Graph classification / regression

↑  
Supervised learning  
↑

Data: large number of labeled graphs

↑  
The network needs to handle  
different size graphs with  
varying connectivity  
(node degree)

## 2) Node-prediction

↑

Semi-supervised learning:

some labeled  
nodes,

also edge  
info taken  
into account

↑  
typically single graph,  
node degree may vary

e.g. neighboring  
nodes may be  
more likely  
to be similar

## 3) Clustering: can be unsupervised

The GNN may do  
some automatic  
clustering

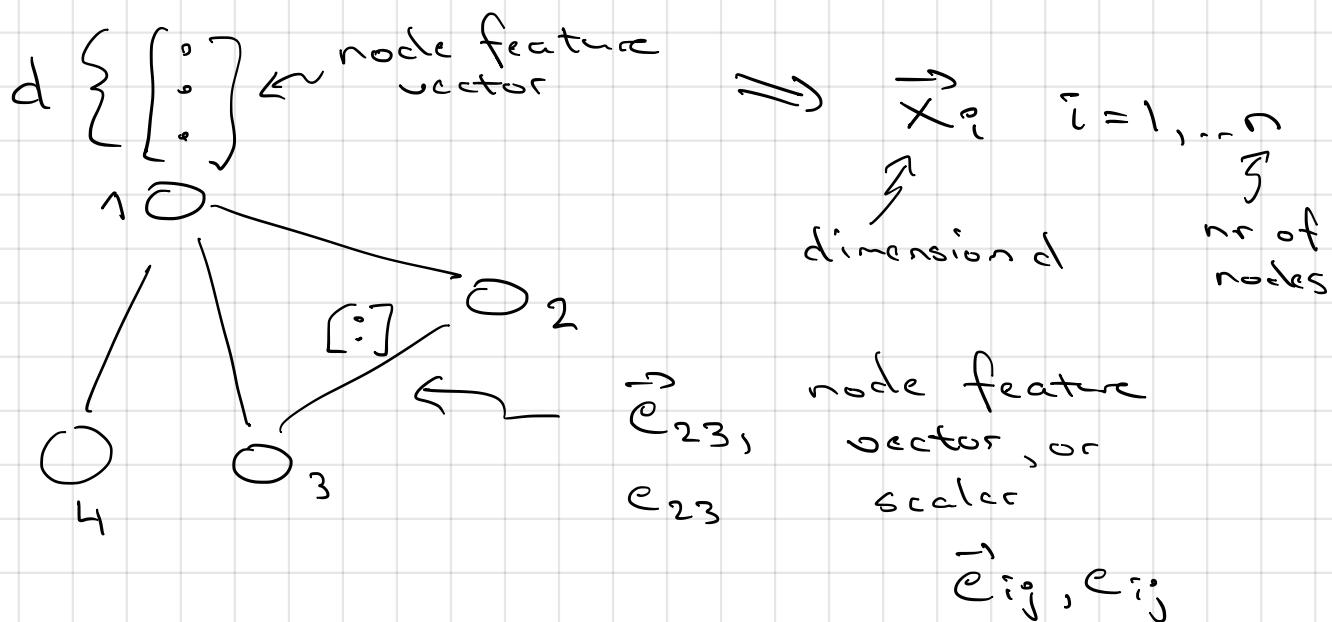
graph → new graph

where info  
from neighboring  
nodes is absorbed  
also, lower dim.  
embedding

Graphs: not only nodes and edges  
(in this context)

"decorated graph"

information on nodes and edges



The structure is given by an  
Adjacency matrix  $A_{ij}$

$A_{ij} = 1$  if edge  $i \rightarrow j$

$A_{ij} = 0$  if no edge

Undirected graph  $A_{ij} = A_{ji}$

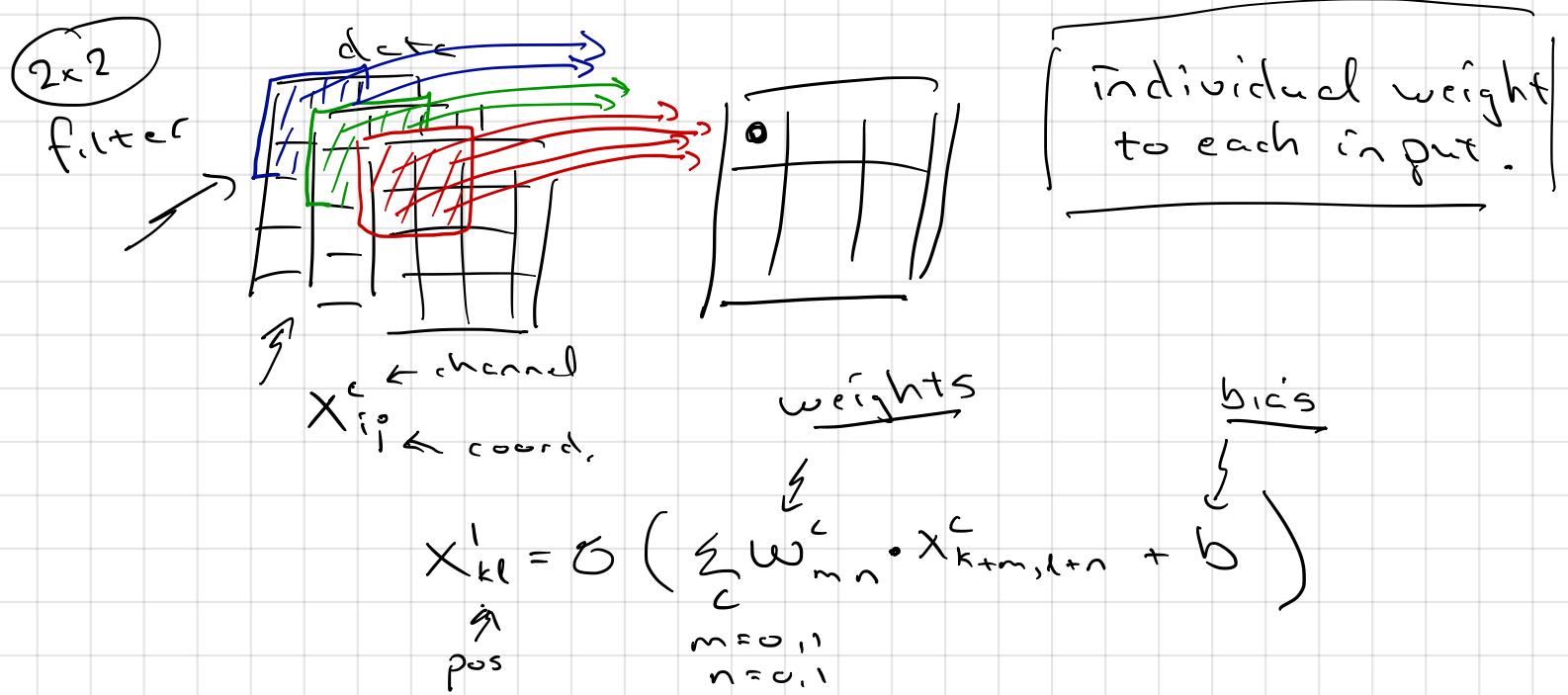
Ex above

$$A = \begin{bmatrix} 0 & 1 & 1 & 1 \\ 1 & 0 & 1 & 0 \\ 1 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 \end{bmatrix}$$

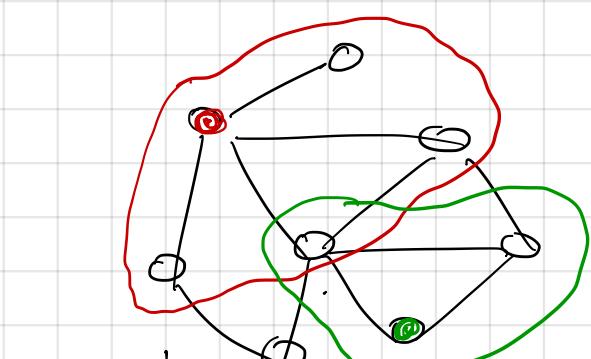
# Convolutions on graphs

Collect information from neighbouring nodes seems like a natural construction.

Ordinary convolution layer on grid :



On a graph ?



Since the number of neighbors very, not simple to have different weights to diff nodes

Discuss  
ingredients

$\vec{x}_i$  feature vector on node  $i$

$A_{ij}$  adjacency matrix

$\sigma$ : non-linear fn  
(i.e ReLU  $\dots$  other)

- Basic construction GCD(G) (Kipf & Welling)

- Isomorphic mapping:  $\tilde{A}_{ij} = A_{ij}$

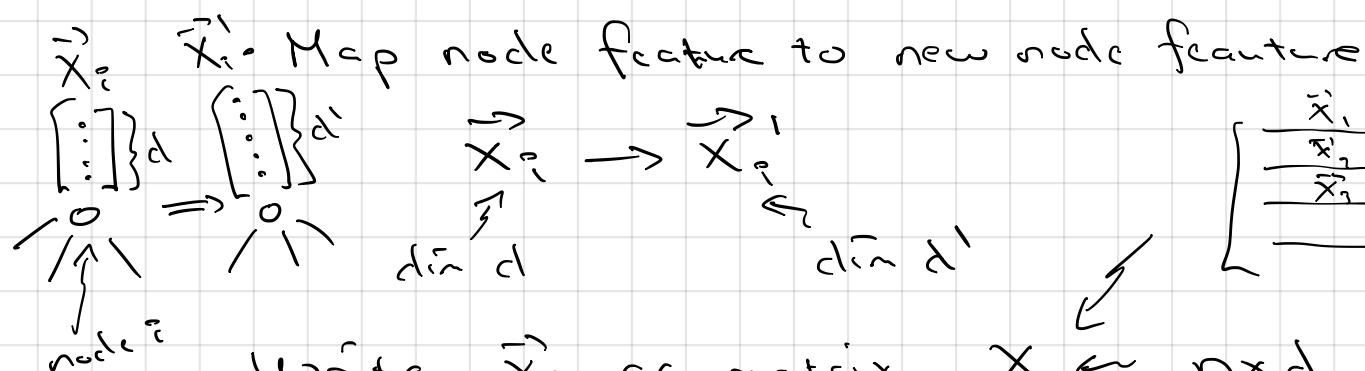
same number of nodes,  
same connectivity

- Include self-loops in  $A$ :

$$\tilde{A} = A + \mathbb{I} \quad (\tilde{A}_{ii} = A_{ii} + \delta_{ii})$$

- Degree matrix  $\tilde{D}_{ii} = \sum_{j=1}^n \tilde{A}_{ij} = \deg_i$   
number of  
neighbors  
incl. self

$$\tilde{D}_{ii} = 0$$



④ 
$$X' = G \left( \underbrace{\tilde{D}^{-\frac{1}{2}} \tilde{A} \tilde{D}^{-\frac{1}{2}}}_{\substack{\text{adds up} \\ \text{all} \\ \text{neighbor} \\ \text{vectors} \\ \text{normalized} \\ \text{by degree}}} \times \omega \right)$$

$\omega$ : weight matrix  $d \times d'$

$\tilde{D}^{-\frac{1}{2}}$  transforms each node vector  $d \rightarrow d'$

⑤ 
$$X'_i = G \left( \sum_j \frac{\tilde{A}_{ij}}{\sqrt{\deg_i} \sqrt{\deg_j}} \omega^T \tilde{x}_j \right)$$

just like a non-activated dense layer on each node

collects information "one-hop"  
per layer. Isotropically, i.e. all  
neighbors have equal influence.

- Simple variation Graph Conv

$$\vec{x}_i' = \sigma(\omega_1 \vec{x}_i + \omega_2 \sum_j e_{ij} \vec{x}_j)$$

edge weight  
 $e_{ij} = 0$  if  $A_{ij} = 0$

takes edge weights into account  
fixed

- We can also make the absorption of information from neighboring nodes adaptive (trainable).

### Graph attention layer GAT Conv

$$\vec{x}_i' = \sigma \left( \sum_j x_{ij} w^T \vec{x}_j \right)$$

fixed weight matrix

$$x_{ij} = \text{softmax}(\varepsilon_{ij}) = \frac{e^{\varepsilon_{ij}}}{\sum_k e^{\varepsilon_{ik}}}$$

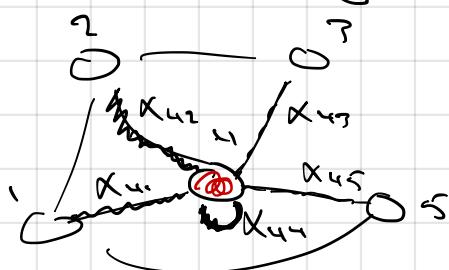
$$\sum_j x_{ij} = 1$$

$$\epsilon_{ij} = \sigma' \left( \vec{a} \cdot \underbrace{\left[ \vec{w_x}_i \| \vec{w_x}_j \right]}_{\text{concatenation}} \right) \tilde{A}_{ij}$$

possibly  
 otherwise  
 activation  
 $\vec{a}$   
 attention  
 vector

to make  
 sure  
 $A_{ij} = 0$   
 if nodes  
 are not  
 adjacent

Attention coefficients,  
trainable edge weights



more input from  
some neighbouring  
nodes

- One can also have several attention heads with attention vcs  $\vec{a}^k$  & weights  $\vec{w}^k$

$$\vec{x}_i^k = \frac{1}{K} \vec{x}_i^K = \begin{bmatrix} \vdots \\ \vec{x}_i^1 \\ \vdots \\ \vec{x}_i^K \\ \vdots \end{bmatrix} \quad \begin{array}{l} \vec{x}_i^1 \text{ head1} \\ \vec{x}_i^2 \text{ head2} \\ \vec{x}_i^3 \text{ head3} \end{array}$$

These types of structures have had great success for Natural Language Processing

- NB! Convolutional graph layers collect information similarly to each node  $\Rightarrow$  after each layer info tends to get smeared out. Too many layers can be bad. (Depending on task.)

Now we have learned how to collect information in the graph.

New graph with node features  $\vec{X}_i^1$  that reflect surrounding of node  $i$ .

What's next? Depends on task

Discuss

- To do node classification: end with a softmax activated convolution layer with  $d$  = number of classes  

nr of node classes {  $\begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$  }  
 $\xrightarrow{\text{conv}} \begin{bmatrix} ? \end{bmatrix}$   $\xrightarrow{\text{softmax}}$   $d \begin{bmatrix} 0.3 \\ 0.2 \\ 0.5 \end{bmatrix}$   
 node is predicted in class  $i$   
 $d = \text{number of classes}$

train over loss w.r.t. labeled nodes.
- To do graph classification: we need to reduce info to one vector, a.k.a. graph embedding.

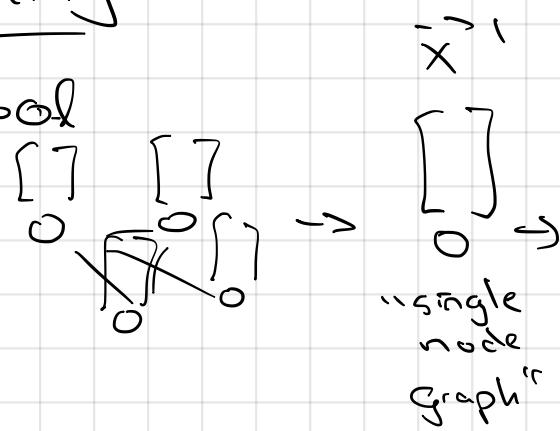
For this use Pooling layer

- Simplest Global pooling

E.g. global mean pool

$$\vec{X}^1 = \frac{1}{n} \sum_{i=1}^n \vec{X}_i^1$$

dense softmax  $\begin{bmatrix} \text{class} \\ \text{predic.} \end{bmatrix}$



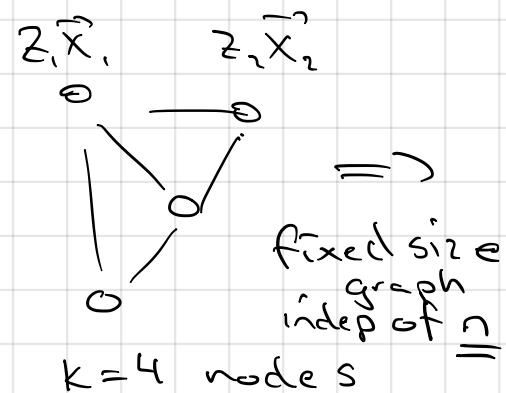
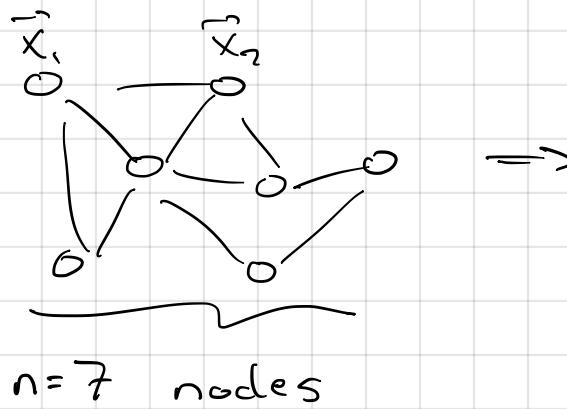
- More sophisticated: max-k pooling

select k most important nodes  
using "self-attention"

$$z_i = \sigma(\vec{a} \cdot \vec{x}_i)$$

↑  
attention vector

Pick k-nodes with highest value  $z$



## Homework A

- Prob. 1 Do semisupervised node-classification  
2 P warm up

- Prob 2 Do graph-classification  
8 P data from a problem related to quantum computing