

# Graph Neural Networks, Homework A

**Course:** Advanced machine learning using neural networks, TIF360/FYM360 (2022)

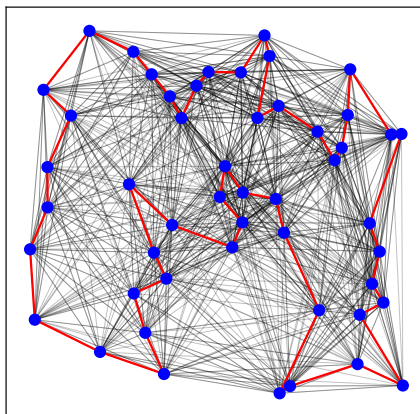
**Contact:** Mats Granath, mats.granath@physics.gu.se

David Fitzek, davidfi@chalmers.se

---

## 1 Background

The archetypical problem for neural network applications is image classification. Here the convolutional network (CNN) is the standard architecture. CNN's can be trained to identify generic features of images, and apply those in a parameter efficient way across the input. Mathematically images are represented as data points on a regular grid, e.g. integer values in a  $3 \times 28 \times 28$  tensor for the CIFAR10 RGB images. However, other types of real world data may be better represented as graphs, i.e. nodes (or vertices) connected by edges (or links). Both the nodes and the edges may have features, typically represented by scalars or vectors. What characterizes a graph as opposed to a grid, is the irregular structure where the number of adjacent neighbors will vary from node to node, and the number of nodes may vary from graph to graph even within the same data set.



Figur 1: Complete graph with 50 nodes where edges have a scalar feature correspond to euclidean distance. The red edges marks the TSP solution on this graph. Image from [1].

An example of graph data (fig. 1) is the euclidean distances between cities, where finding the shortest cycle through all the nodes corresponds to the travelling salesperson (TSP) problem. To express this data on a regular grid would give a very sparse representation, which is likely to make the standard convolutional network inefficient and difficult to train and to generalize between problem instances. The machine learning objective would be to predict the correct edges that proposes a solution to the TSP for a graph, given a set of labeled training graphs. This corresponds to what's known as a link prediction task; for each potential link predict whether it's in the TSP cycle of the graph or not. Other types of machine learning tasks that may be relevant for graphs is node classification (e.g. categorizing people according to subcommunities in social networks), graph regression (e.g. to what extent does the molecule described by a graph have a particular chemical property), or graph classification. Depending on application and type of data,

graph neural networks can be used for both supervised learning (e.g. graph classification) or semi-supervised learning (e.g. node classification). The first part of the homework will focus on the latter, where the graph proximity of labeled and unlabeled nodes can be utilized to improve performance over supervised learning algorithms on the same data. The second part of the homework is a supervised graph classification task.

Two recommended references that provide summaries of graph neural networks are [1] and [2]. A more general background to 'Geometric deep learning' can be found in [3].

## 2 Graph convolutions

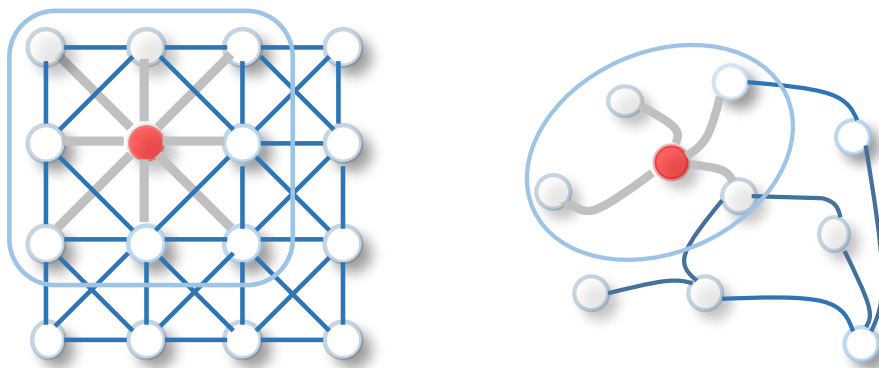


Figure 2: Illustration of a standard convolutional kernel with a fixed size compared to a graph convolution where the kernel size effectively depends on the connectivity of the graph. Image from [2].

A challenge with handling graphs in a neural network pipeline is that the data is irregular: the size of the graph given by the number of nodes and the degree of each node (corresponding to the number of neighbors) will vary both between and within graphs. A graph with  $n$  nodes is represented by an adjacency matrix  $A_{ij}$  with  $i, j \in \{1, \dots, n\}$  with binary (1 or 0) entries indicating edge or no edge. Typically, this is a sparse object represented numerically as a list of edges  $\{\{i, j\}\}$ . (The graph may or may not be directed, corresponding to a symmetric or not adjacency matrix.) The diagonal degree matrix is given by  $D_{ii} = \sum_j A_{ij}$ . Each node has  $d$  features represented by a vector  $\vec{x}_i$  or the corresponding  $n \times d$  matrix  $X$ . There may also be  $m$  edge features given by a vector  $\vec{e}_{ij}$  for the edge connecting node  $i$  to node  $j$ .

The basic structure of a graph convolution is a mapping from a graph with feature matrix  $X$  to an isomorphic graph (i.e. with the same adjacency matrix) with feature matrix  $X'$ . The feature vectors after the graph convolution may have a different dimension  $d'$ , so that  $X'$  is an  $n \times d'$  matrix. This is similar to how a dense layer with  $d'$  neurons would map a  $d$  dimensional input vector to a  $d'$  dimensional output, but implemented in a way that takes into account information from the feature vectors on adjacent nodes of the graph.

One of the earlier and standard forms of a graph convolution was suggested in [4] (GCNConv in PyTorch Geometric) as follows. The node feature matrix  $X^{l+1}$  in layer

$l + 1$ , is given by

$$X^{l+1} = \sigma(\tilde{D}^{-1/2} \tilde{A} \tilde{D}^{-1/2} X^l W^l) \quad (1)$$

where  $\tilde{A} = A + I_N$  (added self edges),  $\tilde{D}$  is the corresponding adjacency matrix (which acts to normalize the feature vector irrespectively of the node connectivity),  $W$  is a trainable  $d_l \times d_{l+1}$  weight matrix, and  $\sigma$  is the non-linearity (e.g. ReLU) that acts element-wise (and may also include a trainable bias per feature). After a trainable linear transformation (with  $W$ ) the new feature vector at each node is thus an average over all the nodes to which it is connected processed through the non-linearity. With each successive layer, feature information travels in from further out in the graph.

Note that the weight matrix  $W$  is equivalent to that of a dense layer  $d_\ell \rightarrow d_{\ell+1}$ , which makes the graph convolution quite distinct from the standard grid convolution which has different trainable weights to each input pixel in the kernel. The reason for this is the variable size of the input (number of neighboring nodes) which would require a kernel of varying size to use the standard CNN type. Note also that this form of the graph convolution does not take edge features into account but only the structure of the graph.

There are of course a large number of more advanced convolutions that expand on this. For example, in the case that the edge features are a scalar, i.e. an edge weight,  $e_{ij}$  which is 0 if there is no edge, there is [GraphConv](#) which has the simple intuitive form

$$\vec{x}_i^{l+1} = \sigma(W_1^l \vec{x}_i^l + \sum_j e_{ij} W_2^l \vec{x}_j^l) \quad (2)$$

where  $W_1$  and  $W_2$  are two independent  $d_{l+1} \times d_l$  weight matrices. Another example of graph convolutions, [GATConv](#) [5] is based on the attention mechanism where the idea is to train a more selective (anisotropic) sampling of the feature vectors of the adjacent nodes. The convolution (feature vector  $d \rightarrow d'$ ) is standard,

$$\vec{x}_i' = \sigma\left(\sum_j \alpha_{i,j} W \vec{x}_j\right), \quad (3)$$

but with sophistication hiding in the attention coefficients  $\alpha_{i,j}$ . The latter are given by  $\alpha_{ij} = \text{softmax}(\epsilon_{ij})$  with

$$\epsilon_{ij} = \text{LeakyReLU}(\vec{a} \cdot [W \vec{x}_i || W \vec{x}_j]) \tilde{A}_{ij} \quad (4)$$

where  $\vec{a}$  is a dimension  $2d'$  trainable weight vector,  $[||\cdot]$  is the concatenation (i.e. putting the two vectors on top of each other), and  $\tilde{A}$  is included to only allow attending on adjacent nodes or the node itself. This can also be extended to having several attention heads  $\epsilon_{ij}^k$ ,  $k = 1, \dots, K$ , each with its own weight matrices and attention vectors.

The graph convolutional layers are typically capped by a dense network and pooling layers that give the required output, depending on the structure of the targets. (See, e.g. appendix of [1].) For the purpose of node classification on a single graph there is actually no need for any post convolution layers. The graph convolution output is of the form of a feature vector per node that may be given appropriate length on which to apply softmax activation to provide a distribution over the labels. (Additional dense layers could of course also be added.) For graph classification, on the other hand, a single label per graph requires that the information from all the nodes must be gathered and processed to give a single output vector for the graph. Various pooling structures may be used for this, with the simplest being a mean pooling layer that takes the mean feature vector of all nodes which can then be processed through a dense network or go directly to softmax.

### 3 Semi-supervised learning

As mentioned previously, a hallmark of machine learning on graphs is the use of semi-supervised learning, as exemplified by the task of node classification. Given a number of labeled nodes, the objective is to classify all nodes. Using only the node features of the labeled nodes individually for training would correspond to supervised learning, but this would not make use of the implicit information about the relations between labeled and unlabeled nodes that is contained in the graph adjacency matrix. Using a graph convolution, the feature vectors of adjacent unlabeled nodes are transferred to the labeled nodes and vice versa and incorporated in the network training, making it semi-supervised through the use of both labeled and unlabeled data.

An intuitive example of semi-supervised learning is shown in figure 3, demonstrating how a classifier can be constructed from two labeled datapoints complemented by a clustering algorithm. (The latter could be in the form of a non-linear support vector machine, that aims to make a minimal complexity separation of the data into two sets.)

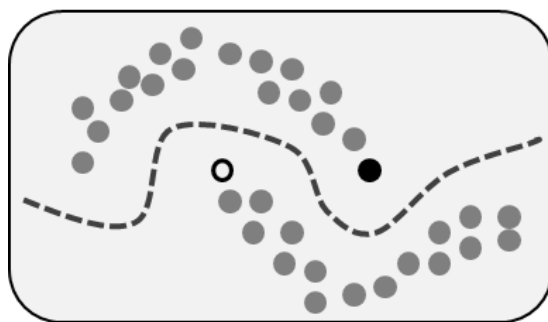


Figure 3: Semi-supervised learning, illustration. Two labeled data points (empty/black) complemented by unlabeled (gray) that are identified by an unsupervised clustering algorithm. A GNN for node classification works along the same principles by using the graph connectivity to identify unlabeled with labeled nodes. (By Techerin, [wikimedia.org](https://www.wikimedia.org/)).

## 4 Assignment

For this assignment you will consider two different datasets, using semi-supervised learning for node classification and supervised learning for graph classification respectively. The recommended library to use for the assignment is [PyTorch Geometric \(PyG\)](https://pytorch-geometric.com/).

### 4.1 Dataset 1

The first set of problems deal with the Cora dataset[6], which is citation graph of scientific papers. The nodes (papers) have features given by sparse bag-of-words vector that correspond to 1433 keywords. There are seven categories of papers based on a classification of the keywords. Out of the 2708 nodes 140 are labeled in the standard training set, 20 in each class. The dataset can be found: <https://github.com/kimiyoung/planetoid> or through PyG.

You are welcome to use the [Colab Notebook](https://colab.research.google.com/) on node classification on the Cora dataset provided in the PyTorch Geometric documentation. Either modify and run directly in

Google Colab or reuse relevant parts of the code in your own Jupyter notebook or other environment.

1) [3p]

- a) From the Planetoid Cora dataset, extract the number of nodes that are in the training set, validation set, and test set. Discuss briefly the implications for supervised versus semi-supervised learning.
- b) Train both a standard dense network using supervised learning and network using the convolution from equation 1 using semi-supervised learning.
- c) In both cases, implement 'early stopping' using the validation set, i.e. the training should evaluate against the validation set to stop training when this fails to improve.
- d) Use t-SNE to visualize the results.

Part (b) and (d) are essentially provided in the Colab notebook referred to above. What you need to implement additionally is the early stopping. You are also welcome, but not required, to try out varying the depth and width of the network. (In the next exercise you are expected to experiment more.)

2) [2p] Try to outperform the results from (1) using one of the more advanced algorithms such as the graph attention network, from eqn. 3-4. Try to optimize performance, varying the number of heads (for GAT), different pooling layers (?), dropout, etc. (Extensive hyperparameter tuning however is not asked for.)

It may not be that easy to improve the results from (1) much (see [Cora Benchmark](#)), the main objective here is to show that you have tried.

## 4.2 Dataset 2

In brief, this dataset contain a collection of 'syndromes' that correspond to (fictive) measurements on a quantum-bit (qubit) memory, a 'surface code', that is constructed by a concatenation of 49 physical qubits. The syndrome gives an indication of where error have occurred and the objective is to use this information to mitigate errors, so called quantum error correction. (In contrast to a classical computer, bit-errors in a quantum computer cannot be measured individually so a more elaborate scheme for error correction based on quantum non-demolition measurements is required.) Each syndrome (graph) comes with a target distribution over four classes. These are called equivalence classes and to correct errors successfully it is important to know the most likely equivalence class of the syndrome. The machine learning objective is to train a graph neural network that given an arbitrary syndrome (graph) outputs the most likely equivalence class. The methodology is supervised learning for graph classification. (It should not be necessary for the assignment task, but if you want to get some more background on quantum error correction, see [7] and references therein.)

The data is provided as a `.pt` file on Canvas that can be loaded using `torch.load`. Note that the target values for the data points are one hot encoded corresponding to the correct class. Each node has a four dimensional feature vector. Note also that the data has an edge attribute corresponding to scalar weights on the edges (these are constructed from inverse distances between nodes on a grid), and that the graph is complete (fully connected).

3) [5p]

- Divide the dataset 80/20 into a training set and a test set. (If you want to, use a partition of the training set for validation.)
- Design a graph neural network with the objective to optimize the accuracy on the test set. Make sure to make use of the edge weights.

## Referenser

- [1] Dwivedi, V. P., Joshi, C. K., Laurent, T., Bengio, Y. & Bresson, X. Benchmarking graph neural networks. *arXiv preprint arXiv:2003.00982* (2020). URL <https://arxiv.org/abs/2003.00982>.
- [2] Wu, Z. *et al.* A comprehensive survey on graph neural networks. *IEEE transactions on neural networks and learning systems* (2020). URL <https://arxiv.org/abs/1901.00596>.
- [3] Bronstein, M. M., Bruna, J., LeCun, Y., Szlam, A. & Vandergheynst, P. Geometric deep learning: going beyond euclidean data. *IEEE Signal Processing Magazine* **34**, 18–42 (2017). URL <https://arxiv.org/abs/1611.08097>.
- [4] Kipf, T. N. & Welling, M. Semi-supervised classification with graph convolutional networks. *arXiv preprint arXiv:1609.02907* (2016). URL <https://arxiv.org/abs/1609.02907>.
- [5] Veličković, P. *et al.* Graph attention networks. *arXiv preprint arXiv:1710.10903* (2017). URL <https://arxiv.org/abs/1710.10903>.
- [6] Yang, Z., Cohen, W. & Salakhudinov, R. Revisiting semi-supervised learning with graph embeddings. In *International conference on machine learning*, 40–48 (PMLR, 2016). URL <https://arxiv.org/abs/1603.08861>.
- [7] Fitzek, D., Eliasson, M., Kockum, A. F. & Granath, M. Deep q-learning decoder for depolarizing noise on the toric code. *Physical Review Research* **2**, 023230 (2020). URL <https://arxiv.org/abs/1912.12919>.