

# Graph Neural Networks

Jan Tönshoff, Timo Gervens, Prof. Dr. Martin Grohe

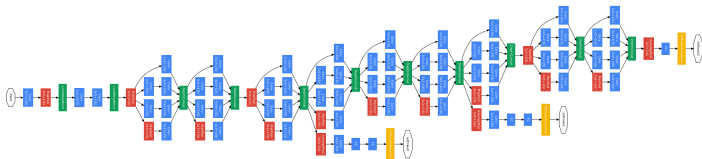
RWTH Aachen

*toenshoff@informatik.rwth-aachen.de*

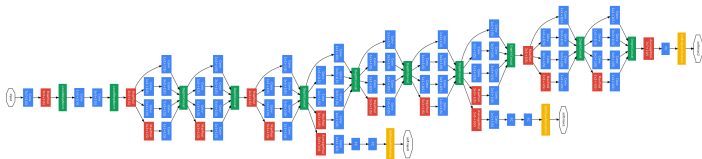
*gervens@informatik.rwth-aachen.de*

November 16, 2022

# Neural Networks



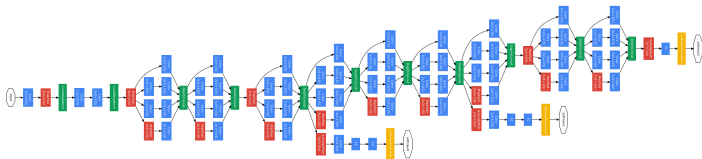
## Neural Networks



Basic Idea:

- Build (large) computational graphs with differentiable tensor operations.

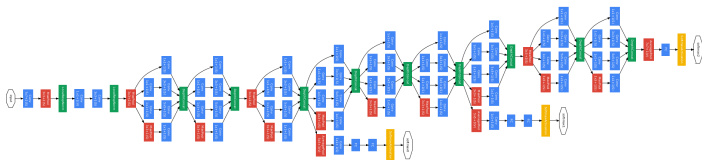
# Neural Networks



## Basic Idea:

- Build (large) computational graphs with differentiable tensor operations.
- Train with Stochastic Gradient Descent on large amounts of data.

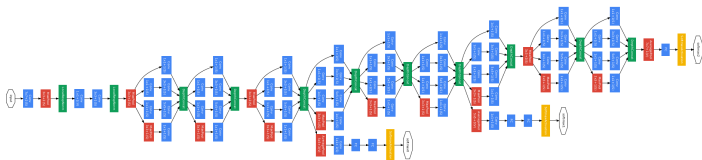
# Neural Networks



## Basic Idea:

- Build (large) computational graphs with differentiable tensor operations.
- Train with Stochastic Gradient Descent on large amounts of data.
- Usually used for end-2-end learning with little feature engineering.

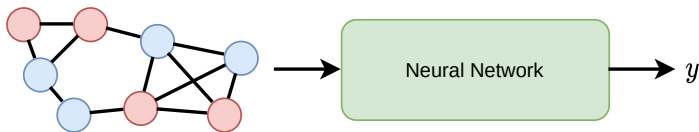
# Neural Networks



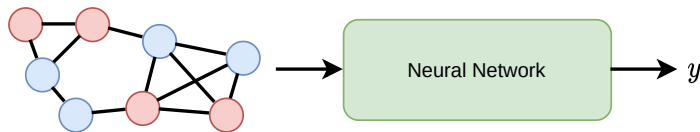
## Basic Idea:

- Build (large) computational graphs with differentiable tensor operations.
- Train with Stochastic Gradient Descent on large amounts of data.
- Usually used for end-2-end learning with little feature engineering.
- State-of-the-art for almost all learning domains, including graphs.

# Graph Neural Networks



# Graph Neural Networks

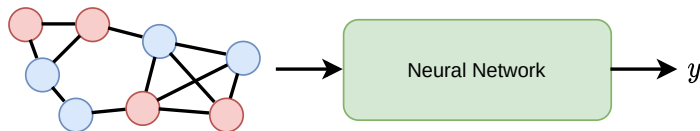


Fundamental Challenges of Graph Learning:

- Process graphs of any size and structure
- Permutation Invariance



# Graph Neural Networks



Fundamental Challenges of Graph Learning:

- Process graphs of any size and structure
- Permutation Invariance

Solution: Message Passing

# Message Passing

Input:

- Graph  $G = (V, E)$
- Node Features  $X_V : V \rightarrow \mathbb{R}^d$
- Edge Features  $X_E : E \rightarrow \mathbb{R}^{d'}$ . (optional)

# Message Passing

Input:

- Graph  $G = (V, E)$
- Node Features  $X_V : V \rightarrow \mathbb{R}^d$
- Edge Features  $X_E : E \rightarrow \mathbb{R}^{d'}$ . (optional)

Idea of Message Passing:

- Initial embedding  $h^0(v) = X_V(v)$ .
- Update latent vertex embedding  $h^\ell(v) \in \mathbb{R}^{d^\ell}$  in each layer  $\ell$ .

# Message Passing

Input:

- Graph  $G = (V, E)$
- Node Features  $X_V : V \rightarrow \mathbb{R}^d$
- Edge Features  $X_E : E \rightarrow \mathbb{R}^{d'}$ . (optional)

Idea of Message Passing:

- Initial embedding  $h^0(v) = X_V(v)$ .
- Update latent vertex embedding  $h^\ell(v) \in \mathbb{R}^{d^\ell}$  in each layer  $\ell$ .
- The update depends on the states of the neighbors.
- Use shared local neural network for the update.

## Message Passing

Elements of a GNN Layer [Wu et al., 2020]:

- Trainable functions  $\mathbf{M}^\ell$ ,  $\mathbf{U}^\ell$
- Aggregation function  $\oplus$  (Sum, Mean, Max, ...)

## Message Passing

Elements of a GNN Layer [Wu et al., 2020]:

- Trainable functions  $\mathbf{M}^\ell$ ,  $\mathbf{U}^\ell$
- Aggregation function  $\bigoplus$  (Sum, Mean, Max, ...)

Function computed by the layer:

$$h^\ell(v) = \mathbf{U}^\ell \left( h^{\ell-1}(v), \bigoplus_{u \in \mathcal{N}(v)} \mathbf{M}^\ell \left( h^{\ell-1}(u), X_E(vu) \right) \right)$$

## Downstream Tasks

### Node-Level Tasks

Apply predictive MLP to final vertex embedding:

$$y(v) = \text{MLP}(h^L(v))$$

## Downstream Tasks

### Node-Level Tasks

Apply predictive MLP to final vertex embedding:

$$y(v) = \text{MLP}(h^L(v))$$

### Graph-Level Tasks

Pool vertex embedding into a graph embedding with *readout* function (Sum, Mean, ...):

$$h(G) = \sum_{v \in V} h^L(v)$$



## Downstream Tasks

### Node-Level Tasks

Apply predictive MLP to final vertex embedding:

$$y(v) = \text{MLP}(h^L(v))$$

### Graph-Level Tasks

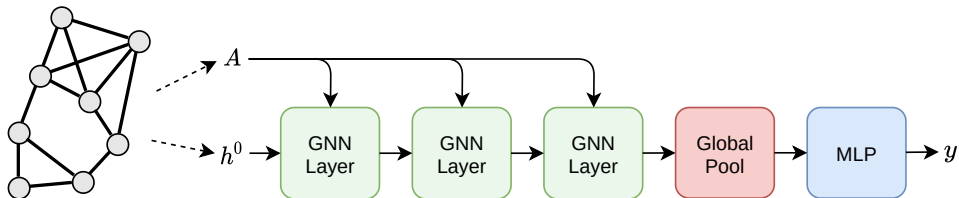
Pool vertex embedding into a graph embedding with *readout* function (Sum, Mean, ...):

$$h(G) = \sum_{v \in V} h^L(v)$$

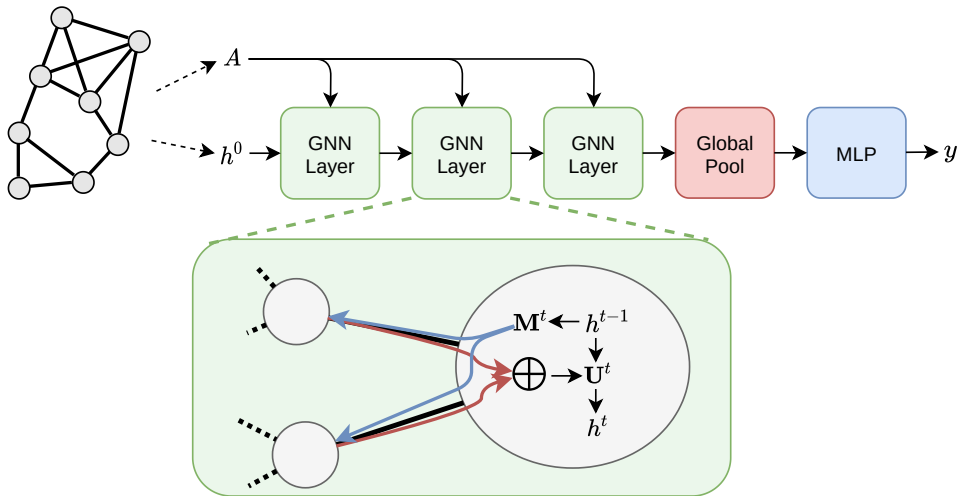
Apply predictive MLP to pooled vector:

$$y(G) = \text{MLP}(h(G))$$

# Graph Neural Networks



# Graph Neural Networks



## Expressiveness of GNNs

Do GNNs seem familiar?

## Expressiveness of GNNs

Do GNNs seem familiar?

We can view latent states as vertex colors.

⇒ **M**, **U** and  $\oplus$  define a HASH function for the colors in the neighborhood of a vertex.

## Expressiveness of GNNs

Do GNNs seem familiar?

We can view latent states as vertex colors.

⇒ **M**, **U** and  $\oplus$  define a HASH function for the colors in the neighborhood of a vertex.

GNNs are a differentiable version of Color Refinement!

## Expressiveness of GNNs

Do GNNs seem familiar?

We can view latent states as vertex colors.

⇒ **M**, **U** and  $\oplus$  define a HASH function for the colors in the neighborhood of a vertex.

GNNs are a differentiable version of Color Refinement!

GNNs are equivalent to the WL-Test. They can not detect substructures that the WL-Kernel could not already detect!

## Expressiveness of GNNs

GNNs are **at most** as expressive as the WL-Kernel.



## Expressiveness of GNNs

GNNs are **at most** as expressive as the WL-Kernel.

Why can we still expect them to (usually) work better?

## Expressiveness of GNNs

GNNs are **at most** as expressive as the WL-Kernel.

Why can we still expect them to (usually) work better?

- WL views all subtrees as strictly different.

## Expressiveness of GNNs

GNNs are **at most** as expressive as the WL-Kernel.

Why can we still expect them to (usually) work better?

- WL views all subtrees as strictly different.
- Slightly different subtrees may have identical effect on the label.

## Expressiveness of GNNs

GNNs are **at most** as expressive as the WL-Kernel.

Why can we still expect them to (usually) work better?

- WL views all subtrees as strictly different.
- Slightly different subtrees may have identical effect on the label.
- A non-injective HASH function may therefore yield a more robust classifier.

## Expressiveness of GNNs

GNNs are **at most** as expressive as the WL-Kernel.

Why can we still expect them to (usually) work better?

- WL views all subtrees as strictly different.
- Slightly different subtrees may have identical effect on the label.
- A non-injective HASH function may therefore yield a more robust classifier.
- GNNs effectively learn a fitting HASH function from data.

## Expressiveness of GNNs

GNNs are **at most** as expressive as the WL-Kernel.

Why can we still expect them to (usually) work better?

- WL views all subtrees as strictly different.
- Slightly different subtrees may have identical effect on the label.
- A non-injective HASH function may therefore yield a more robust classifier.
- GNNs effectively learn a fitting HASH function from data.
- GNNs also seamlessly handle real-valued node and edge features.

## Exercise 2 and 3

In the next two sheets you will implement GNNs:

- Sheet 2: Implement the GCN using simple dense data structures.
- Sheet 3: Implement a general GNN Layer with efficient sparse data structures.

# Graph Convolutional Neural Network

The Graph Convolutional Neural Network (GCN) [Kipf and Welling, 2016]:

- One of the first GNNs that actually worked.
- Simple and effective.
- Still often used as the default GNN for many tasks.
- Limitations: Does not incorporate edge features.



## The GCN Layer

Components:

- Trainable weight matrix  $W^{(\ell)} \in \mathbb{R}^{d^{(\ell-1)} \times d^{(\ell)}}$
- Activation Function  $\sigma$  (usually ReLU)

# The GCN Layer

Components:

- Trainable weight matrix  $W^{(\ell)} \in \mathbb{R}^{d^{(\ell-1)} \times d^{(\ell)}}$
- Activation Function  $\sigma$  (usually ReLU)

Update rule:

$$h^{\ell}(v) = \sigma\left(\sum_{u \in N(v) \cup \{v\}} \frac{1}{\sqrt{d_v d_u}} \cdot h^{\ell-1}(u) \cdot W^{(\ell)}\right)$$

with  $d_i = \deg(v_i) + 1$ .

# The GCN Layer

Components:

- Trainable weight matrix  $W^{(\ell)} \in \mathbb{R}^{d^{(\ell-1)} \times d^{(\ell)}}$
- Activation Function  $\sigma$  (usually ReLU)

Update rule:

$$h^{\ell}(v) = \sigma\left(\sum_{u \in N(v) \cup \{v\}} \frac{1}{\sqrt{d_v d_u}} \cdot h^{\ell-1}(u) \cdot W^{(\ell)}\right)$$

with  $d_i = \deg(v_i) + 1$ .

In the general GNN framework: **M** is the identity function, **U** is a one layer perceptron without bias and  $\oplus$  is a weighted sum.

## Implementing the GCN Layer

Stack vertex embeddings in matrix:  $H^{(\ell)} \in \mathbb{R}^{|V| \times d^{(\ell)}}$

## Implementing the GCN Layer

Stack vertex embeddings in matrix:  $H^{(\ell)} \in \mathbb{R}^{|V| \times d^{(\ell)}}$

Pre-Compute normalized adjacency matrix  $\tilde{A} \in \mathbb{R}^{|V| \times |V|}$ :

$$\tilde{A}_{i,j} = \begin{cases} \frac{1}{\sqrt{d_i d_j}} & \text{if } v_i v_j \in E \text{ or } i = j \\ 0 & \text{otherwise} \end{cases}$$

## Implementing the GCN Layer

Stack vertex embeddings in matrix:  $H^{(\ell)} \in \mathbb{R}^{|V| \times d^{(\ell)}}$

Pre-Compute normalized adjacency matrix  $\tilde{A} \in \mathbb{R}^{|V| \times |V|}$ :

$$\tilde{A}_{i,j} = \begin{cases} \frac{1}{\sqrt{d_i d_j}} & \text{if } v_i v_j \in E \text{ or } i = j \\ 0 & \text{otherwise} \end{cases}$$

The layer can now be expressed with efficient matrix operations:

$$H^{(\ell)} = \sigma(\tilde{A} \cdot H^{(\ell-1)} \cdot W^{(\ell)})$$

# Considerations

Implementing GNNs with dense tensors:

- Dimension 0 is reserved as the *batch* dimension.
- For a batch of size  $b$  we have  $A \in \mathbb{R}^{b \times |V| \times |V|}$  and  $H^0 \in \mathbb{R}^{b \times |V| \times d}$
- Graphs have different sizes: Pad  $A$  and  $X_V$  with zeros to uniform sizes over the whole dataset.

## Exercise 2

Implement the GCN in PyTorch:

- Pre-Process graphs to compute padded, normalized adjacency matrices.
- Implement the GCN Layer as PyTorch Module.
- Construct Node- and Graph-Level GNNs with your GCN Layer.



## Exercise 2

Implement the GCN in PyTorch:

- Pre-Process graphs to compute padded, normalized adjacency matrices.
- Implement the GCN Layer as PyTorch Module.
- Construct Node- and Graph-Level GNNs with your GCN Layer.

Perform 10-fold Cross-Validation on 4 datasets:

- Cora
- Citeseer
- NCI1
- ENZYMES

# PyTorch

Numpy-like math framework with GPU support and backpropagation.

Defacto Standard for Graph Neural Networks.

Many useful open source extensions (`torch-scatter`, `torch-vision`, ...)

# PyTorch

Numpy-like math framework with GPU support and backpropagation.

Defacto Standard for Graph Neural Networks.

Many useful open source extensions (`torch-scatter`, `torch-vision`, ...)

We will build a Multi-Layer Perceptron as an example.

# A Basic PyTorch Module

```
1 import torch
2
3 class Layer(torch.nn.Module):
4
5     def __init__(self, dim_in, dim_out, is_linear=False):
6         super(Layer, self).__init__()
7         self.is_linear = is_linear
8
9         # use Kaiming Init when using ReLU
10        self.W = torch.nn.Parameter(torch.zeros(dim_in, dim_out))
11        torch.nn.init.kaiming_normal_(self.W)
12
13        self.b = torch.nn.Parameter(torch.normal(0.0, 1.0, size=(dim_out,)))
14
15    def forward(self, x):
16        # linear transformation on input
17        y = torch.matmul(x, self.W) + self.b
18
19        # apply activation
20        if not self.is_linear:
21            y = torch.relu(y)
22        return y
23
```

# Combining Modules

```
25 class MLP(torch.nn.Module):
26
27     def __init__(self, input_dim, output_dim, hidden_dim, num_layers):
28         super(MLP, self).__init__()
29         self.num_layers = num_layers
30
31         # add sub-modules as attribute
32         self.input_layer = Layer(input_dim, hidden_dim)
33         self.output_layer = Layer(hidden_dim, output_dim, is_linear=True)
34
35         # Store multiple submodules in "ModuleList"
36         self.hidden_layers = torch.nn.ModuleList(
37             [Layer(hidden_dim, hidden_dim) for _ in range(num_layers-2)]
38         )
39
40     def forward(self, x):
41         # apply layers
42         y = self.input_layer(x)
43         for i in range(self.num_layers-2):
44             y = self.hidden_layers[i](y)
45         y = self.output_layer(y)
46         return y
47
```

## Before Training

```
51 import numpy as np
52 from torch.utils.data import TensorDataset, DataLoader
53
54 # load numpy data and cast to torch tensors
55 X, Y = np.load('training_data.npy')
56 X = torch.tensor(X, dtype=torch.float32)
57 Y = torch.tensor(Y)
58
59 # create dataset and loader for mini batches
60 train_dataset = TensorDataset(X, Y)
61 train_loader = DataLoader(train_dataset, batch_size=32, shuffle=True)
62
63 # set to 'cuda' if gpu is available
64 device = 'cpu'
65
66 # construct neural network and move it to device
67 model = MLP(input_dim=2, output_dim=10, hidden_dim=100, num_layers=4)
68 model.train()
69 model.to(device)
70
71 # construct optimizer
72 opt = torch.optim.Adam(model.parameters(), lr=0.001)
73
```

# Training Loop

```
74 from torch.nn.functional import cross_entropy
75
76 for epoch in range(100):
77     for x, y_true in train_loader:
78         # set gradients to zero
79         opt.zero_grad()
80
81         # move data to device
82         x = x.to(device)
83         y_true = y_true.to(device)
84
85         # forward pass and loss
86         y_pred = model(x)
87         loss = torch.nn.functional.cross_entropy(y_pred, y_true)
88
89         # backward pass and sgd step
90         loss.backward()
91         opt.step()
92
93     # Insert a validation loop and some logging here...
94
```

## References I

- Zonghan Wu, Shirui Pan, Fengwen Chen, Guodong Long, Chengqi Zhang, and S Yu Philip. A comprehensive survey on graph neural networks. *IEEE transactions on neural networks and learning systems*, 2020.
- Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. *arXiv preprint arXiv:1609.02907*, 2016.