Lecture 3: Graph Neural Networks

İsmail İlkan Ceylan

Advanced Topics in Machine Learning, University of Oxford

Relational Learning

25.01.2021

Lecture 1 - 2 (Knowledge graph embeddings): We discussed shallow node embedding models in the context of knowledge graphs, where the idea is to learn low-dimensional embeddings of the nodes in a graph.

Lecture 1 - 2 (Knowledge graph embeddings): We discussed shallow node embedding models in the context of knowledge graphs, where the idea is to learn low-dimensional embeddings of the nodes in a graph.

Lecture 3 - 8 (Graph neural networks): We will study graph neural networks, which are neural networks that learn representations of nodes that depend on the structure of the graph.

Lecture 1 - 2 (Knowledge graph embeddings): We discussed shallow node embedding models in the context of knowledge graphs, where the idea is to learn low-dimensional embeddings of the nodes in a graph.

Lecture 3 - 8 (Graph neural networks): We will study graph neural networks, which are neural networks that learn representations of nodes that depend on the structure of the graph.

Recall key differences between shallow and deep embedding models:

Lecture 1 - 2 (Knowledge graph embeddings): We discussed shallow node embedding models in the context of knowledge graphs, where the idea is to learn low-dimensional embeddings of the nodes in a graph.

Lecture 3 - 8 (Graph neural networks): We will study graph neural networks, which are neural networks that learn representations of nodes that depend on the structure of the graph.

Recall key differences between shallow and deep embedding models:

again on these entities. They are also inherently limited to single-graph tasks.

• Shallow embedding models are transductive in that they do not apply to novel entities unless we train

Lecture 1 - 2 (Knowledge graph embeddings): We discussed shallow node embedding models in the context of knowledge graphs, where the idea is to learn low-dimensional embeddings of the nodes in a graph.

Lecture 3 - 8 (Graph neural networks): We will study graph neural networks, which are neural networks that learn representations of nodes that depend on the structure of the graph.

Recall key differences between shallow and deep embedding models:

- again on these entities. They are also inherently limited to single-graph tasks.
- and can also generalise to novel data points, i.e., inductive models.

• Shallow embedding models are transductive in that they do not apply to novel entities unless we train

• We are interested in more elaborate embedding models which can use any feature information available,

Lecture 1 - 2 (Knowledge graph embeddings): We discussed shallow node embedding models in the context of knowledge graphs, where the idea is to learn low-dimensional embeddings of the nodes in a graph.

Lecture 3 - 8 (Graph neural networks): We will study graph neural networks, which are neural networks that learn representations of nodes that depend on the structure of the graph.

Recall key differences between shallow and deep embedding models:

- again on these entities. They are also inherently limited to single-graph tasks.
- and can also generalise to novel data points, i.e., inductive models.

The remaining lectures on relational learning focus on graph neural networks: We will cover the basics and widely employed models, while also focusing on some limitations and extensions proposed in the literature.

• Shallow embedding models are transductive in that they do not apply to novel entities unless we train

• We are interested in more elaborate embedding models which can use any feature information available,

• Motivation and relational inductive bias

- Motivation and relational inductive bias
- Message passing neural networks

- Motivation and relational inductive bias
- Message passing neural networks
- Graph representation learning tasks

- Motivation and relational inductive bias
- Message passing neural networks
- Graph representation learning tasks
- Perspectives for graph neural networks

- Motivation and relational inductive bias
- Message passing neural networks
- Graph representation learning tasks
- Perspectives for graph neural networks
- Summary

Motivation and Relational Inductive Bias

Beyond Euclidian Spaces









Beyond Euclidian Spaces





Geometric deep learning is an umbrella term for deep learning over non-Euclidian spaces (Bronstein et al.), primarily graphs and manifolds. The focus of this course is on graph representation learning.





The landscape of graphs is rich: Directed, undirected? Weighted graphs? Labelled (multi-relational) graphs? Node/edge features?

The landscape of graphs is rich: Directed, undirected? Weighted graphs? Labelled (multi-relational) graphs? Node/edge features?

We focus on undirected, unweighted, and unlabelled (single-relation) graphs, and assume deterministic node features, unless explicitly stated otherwise. We sometimes speak about different types of graphs, but this shall become clear from the context.

The landscape of graphs is rich: Directed, undirected? Weighted graphs? Labelled (multi-relational) graphs? Node/edge features?

We focus on undirected, unweighted, and unlabelled (single-relation) graphs, and assume deterministic node features, unless explicitly stated otherwise. We sometimes speak about different types of graphs, but this shall become clear from the context.

Remark: This choice is partly for brevity, as most of the concepts can be extended to other types of graphs. Not all extensions are straight-forward, though, e.g., extensions to more general graphs such as hyper-graphs are non-trivial and do not presently admit any efficient models.

The landscape of graphs is rich: Directed, undirected? Weighted graphs? Labelled (multi-relational) graphs? Node/edge features?

We focus on undirected, unweighted, and unlabelled (single-relation) graphs, and assume deterministic node features, unless explicitly stated otherwise. We sometimes speak about different types of graphs, but this shall become clear from the context.

Remark: This choice is partly for brevity, as most of the concepts can be extended to other types of graphs. Not all extensions are straight-forward, though, e.g., extensions to more general graphs such as hyper-graphs are non-trivial and do not presently admit any efficient models.

Notation: We consider graphs of the form G = (V, E) with a feature matrix $\mathbf{X} \in \mathbb{R}^{d \times V_G}$, where d is the embedding dimensionality and V_G denotes the set of vertices/nodes.

The landscape of graphs is rich: Directed, undirected? Weighted graphs? Labelled (multi-relational) graphs? Node/edge features?

We focus on undirected, unweighted, and unlabelled (single-relation) graphs, and assume deterministic node features, unless explicitly stated otherwise. We sometimes speak about different types of graphs, but this shall become clear from the context.

Remark: This choice is partly for brevity, as most of the concepts can be extended to other types of graphs. Not all extensions are straight-forward, though, e.g., extensions to more general graphs such as hyper-graphs are non-trivial and do not presently admit any efficient models.

embedding dimensionality and V_G denotes the set of vertices/nodes.

Intuitively, for each node u, we have a feature vector \mathbf{x}_{u} which can be, e.g., domain-specific attributes, or node degrees, or simply one-hot encodings.

Notation: We consider graphs of the form G = (V, E) with a feature matrix $\mathbf{X} \in \mathbb{R}^{d \times V_G}$, where d is the

The landscape of graphs is rich: Directed, undirected? Weighted graphs? Labelled (multi-relational) graphs? Node/edge features?

We focus on undirected, unweighted, and unlabelled (single-relation) graphs, and assume deterministic node features, unless explicitly stated otherwise. We sometimes speak about different types of graphs, but this shall become clear from the context.

Remark: This choice is partly for brevity, as most of the concepts can be extended to other types of graphs. Not all extensions are straight-forward, though, e.g., extensions to more general graphs such as hyper-graphs are non-trivial and do not presently admit any efficient models.

embedding dimensionality and V_G denotes the set of vertices/nodes.

Intuitively, for each node u, we have a feature vector \mathbf{x}_{u} which can be, e.g., domain-specific attributes, or node degrees, or simply one-hot encodings.

rows of the adjacency matrix.

Notation: We consider graphs of the form G = (V, E) with a feature matrix $\mathbf{X} \in \mathbb{R}^{d \times V_G}$, where d is the

Finally, we write \mathbf{A}^G to denote the adjacency matrix of a graph G = (V, E), and $\mathbf{A}^G_{[i]} \in \mathbb{R}^{V_G}$ to denote the

These assumptions are helpful since

- independence: no need to model the dependencies,
- identical distribution: generalisation guarantees possible to new/unseen data points.

These assumptions are helpful since

- independence: no need to model the dependencies,
- identical distribution: generalisation guarantees possible to new/unseen data points.

These assumptions are unrealistic in the context of graphs.

These assumptions are helpful since

- independence: no need to model the dependencies,
- identical distribution: generalisation guarantees possible to new/unseen data points.

These assumptions are unrealistic in the context of graphs.

Suppose we are given a single graph, and we want to classify the nodes in the given graph with respect to a certain property (i.e., node classification).

These assumptions are helpful since

- independence: no need to model the dependencies,
- identical distribution: generalisation guarantees possible to new/unseen data points.

These assumptions are unrealistic in the context of graphs.

Suppose we are given a single graph, and we want to classify the nodes in the given graph with respect to a certain property (i.e., node classification).

Such properties depend on the other nodes through edges, e.g., imagine functions that rely on node statistics (e.g., #neighbours), or the overall graph structure (e.g., is the node in a cycle).

These assumptions are helpful since

- independence: no need to model the dependencies,
- identical distribution: generalisation guarantees possible to new/unseen data points.

These assumptions are unrealistic in the context of graphs.

Suppose we are given a single graph, and we want to classify the nodes in the given graph with respect to a certain property (i.e., node classification).

Such properties depend on the other nodes through edges, e.g., imagine functions that rely on node statistics (e.g., #neighbours), or the overall graph structure (e.g., is the node in a cycle).

When learning functions over nodes, we cannot treat the nodes independently.



Node degrees?

Contains an odd-length cycle?

Minimum vertex cover size 1, 2?





Node degrees?

Contains an odd-length cycle?

Minimum vertex cover size 1, 2?

Functions over graphs, or nodes, necessarily relate to graph properties, which carry valuable information. This information needs to be taken into account adequately.



Node degrees? Contains an odd-length cycle?

Functions over graphs, or nodes, necessarily relate to graph properties, which carry valuable information. This information needs to be taken into account adequately.

Intuitively, this could be achieved by defining similarity measures, for nodes/graphs, which can then be used for the optimisation task.


Learning over Graphs

Learning over Graphs

• Extracting node, or graph-level statistics or features (indicating, e.g., node/graph similarity), and

- Using these features as input to standard machine learning classifiers.

Learning over Graphs

• Extracting node, or graph-level statistics or features (indicating, e.g., node/graph similarity), and

- Extracting node, or graph-level statistics or features (indicating, e.g., node/graph similarity), and
- Using these features as input to standard machine learning classifiers.

(Kriege et al., 2020) for a recent survey.

Learning over Graphs

Most popular graph similarity functions are studied under the name of graph kernel methods; see, e.g.,

- Extracting node, or graph-level statistics or features (indicating, e.g., node/graph similarity), and
- Using these features as input to standard machine learning classifiers.

(Kriege et al., 2020) for a recent survey.

such as Weisfeiler Lehman kernel (Weisfeiler and Leman, 1968)!

Learning over Graphs

- Most popular graph similarity functions are studied under the name of graph kernel methods; see, e.g.,
- Modern deep learning approaches do not explicitly extract such statistics, but there are nevertheless strong connections between modern graph representation learning and some well-known graph kernel methods,

- Extracting node, or graph-level statistics or features (indicating, e.g., node/graph similarity), and
- Using these features as input to standard machine learning classifiers.

(Kriege et al., 2020) for a recent survey.

such as Weisfeiler Lehman kernel (Weisfeiler and Leman, 1968)!

models, in detail later in the course.

Learning over Graphs

- Most popular graph similarity functions are studied under the name of graph kernel methods; see, e.g.,
- Modern deep learning approaches do not explicitly extract such statistics, but there are nevertheless strong connections between modern graph representation learning and some well-known graph kernel methods,
- We will discuss the Weisfeiler Lehman algorithm, and its connection to popular graph neural network

Earlier deep learning models do not capture graph-structured domains adequately.

layer perceptron:

 $f(G) = MLP(\mathbf{A}_{[}^{G})$

where \oplus is vector concatenation of the rows $\mathbf{A}_{[i]}^G \in \mathbb{R}^{V_G}$ of the adjacency matrix \mathbf{A}^G .

Example: Let us consider multi-layer perceptrons. We can define an embedding of a graph G as a multi-

$$[1] \oplus \ldots \oplus \mathbf{A}^G_{[|V_G|]})$$
,

Earlier deep learning models do not capture graph-structured domains adequately.

layer perceptron:

 $f(G) = \mathsf{MLP}(\mathbf{A}_{[1]}^G \oplus \ldots \oplus \mathbf{A}_{[|V_C|]}^G),$

where \oplus is vector concatenation of the rows $\mathbf{A}_{[i]}^G \in \mathbb{R}^{V_G}$ of the adjacency matrix \mathbf{A}^G .

Problem: This representation depends on the ordering of nodes that we used in the adjacency matrix! Learned representation depends on an ordering that is arbitrary!

Example: Let us consider multi-layer perceptrons. We can define an embedding of a graph G as a multi-

Earlier deep learning models do not capture graph-structured domains adequately.

layer perceptron:

 $f(G) = \mathsf{MLP}(\mathbf{A}_{[1]}^G \oplus \ldots \oplus \mathbf{A}_{[|V_C|]}^G),$

where \oplus is vector concatenation of the rows $\mathbf{A}_{[i]}^G \in \mathbb{R}^{V_G}$ of the adjacency matrix \mathbf{A}^G .

Problem: This representation depends on the ordering of nodes that we used in the adjacency matrix! Learned representation depends on an ordering that is arbitrary!

Other models: Convolutional neural networks are well-defined over grids, but not on graphs. Long shortterm memory networks process sequential (including, e.g., tree-shaped) data, but not graphs: In fact, the above-mentioned problem persists, as by encoding graphs as simple sequences, we lose valuable information.

Example: Let us consider multi-layer perceptrons. We can define an embedding of a graph G as a multi-

Earlier deep learning models do not capture graph-structured domains adequately.

layer perceptron:

where \oplus is vector concatenation of the rows $\mathbf{A}_{[i]}^G \in \mathbb{R}^{V_G}$ of the adjacency matrix \mathbf{A}^G .

Problem: This representation depends on the ordering of nodes that we used in the adjacency matrix! Learned representation depends on an ordering that is arbitrary!

Other models: Convolutional neural networks are well-defined over grids, but not on graphs. Long shortterm memory networks process sequential (including, e.g., tree-shaped) data, but not graphs: In fact, the above-mentioned problem persists, as by encoding graphs as simple sequences, we lose valuable information.

We need a new kind of deep learning framework!

The Quest for a New Framework

Example: Let us consider multi-layer perceptrons. We can define an embedding of a graph G as a multi-

 $f(G) = \mathsf{MLP}(\mathbf{A}_{[1]}^G \oplus \ldots \oplus \mathbf{A}_{[|V_C|]}^G),$

Given a set of graphs \mathscr{G} , we consider functions of the form $f: \mathscr{G} \to \mathbb{R}^{V_G}$.

Given a set of graphs \mathscr{G} , we consider functions of the form $f : \mathscr{G} \to \mathbb{R}^{V_G}$. **Invariance**: We say that a function f is permutation-invariant if for isomorphic graphs $G, H \in \mathscr{G}$ it holds that f(G) = f(H).

Given a set of graphs \mathscr{G} , we consider functions of the form $f: \mathscr{G} \to \mathbb{R}^{V_G}$.

that f(G) = f(H).

permutation π of V_G , it holds that $f(G^{\pi}) = f(G)^{\pi}$.

- **Invariance**: We say that a function f is permutation-invariant if for isomorphic graphs $G, H \in \mathcal{G}$ it holds
- **Equivariance**: We say that a function f is permutation-equivariant if it has the property that for every

Given a set of graphs \mathscr{G} , we consider functions of the form $f: \mathscr{G} \to \mathbb{R}^{V_G}$.

that f(G) = f(H).

permutation π of V_G , it holds that $f(G^{\pi}) = f(G)^{\pi}$.

permute the nodes in the graph.

- **Invariance**: We say that a function f is permutation-invariant if for isomorphic graphs $G, H \in \mathcal{G}$ it holds
- **Equivariance**: We say that a function f is permutation-equivariant if it has the property that for every
- Intuitively, permutation-invariance implies that the function does not depend on the ordering of the nodes in the graph, and permutation-equivariance implies the output of f is permuted in a consistent way when we

Given a set of graphs \mathscr{G} , we consider functions of the form $f: \mathscr{G} \to \mathbb{R}^{V_G}$. that f(G) = f(H).

permutation π of V_G , it holds that $f(G^{\pi}) = f(G)^{\pi}$.

permute the nodes in the graph.

well-defined over these functions, but we cannot speak of permutation-equivariance here.

- **Invariance**: We say that a function f is permutation-invariant if for isomorphic graphs $G, H \in \mathcal{G}$ it holds
- **Equivariance**: We say that a function f is permutation-equivariant if it has the property that for every
- Intuitively, permutation-invariance implies that the function does not depend on the ordering of the nodes in the graph, and permutation-equivariance implies the output of f is permuted in a consistent way when we
- We can also consider functions of different forms, e.g., $f: \mathscr{G} \to \mathbb{R}$. Note that permutation-invariance is still

Given a set of graphs \mathscr{G} , we consider functions of the form $f: \mathscr{G} \to \mathbb{R}^{V_G}$. that f(G) = f(H).

permutation π of V_G , it holds that $f(G^{\pi}) = f(G)^{\pi}$.

permute the nodes in the graph.

well-defined over these functions, but we cannot speak of permutation-equivariance here.

- **Invariance**: We say that a function f is permutation-invariant if for isomorphic graphs $G, H \in \mathcal{G}$ it holds
- **Equivariance**: We say that a function f is permutation-equivariant if it has the property that for every
- Intuitively, permutation-invariance implies that the function does not depend on the ordering of the nodes in the graph, and permutation-equivariance implies the output of f is permuted in a consistent way when we
- We can also consider functions of different forms, e.g., $f: \mathscr{G} \to \mathbb{R}$. Note that permutation-invariance is still
- **Argument**: These properties are important for graph representation learning tasks, as they provide a strong relational inductive bias! The goal is to develop a deep learning framework enhanced with these properties.

Question: Wouldn't it be possible to learn properties such as invariance and equivariance from data during training? Why do we need a framework that ensures such properties?

Question: Wouldn't it be possible to learn properties such as invariance and equivariance from data during training? Why do we need a framework that ensures such properties?

Discussion: Let us loosely denote by M a model that does not have these properties and by M^+ a model that has these properties. We aim to learn a rather simple permutation-invariant or equivariant function f:

Question: Wouldn't it be possible to learn properties such as invariance and equivariance from data during training? Why do we need a framework that ensures such properties?

Discussion: Let us loosely denote by M a model that does not have these properties and by M^+ a model that has these properties. We aim to learn a rather simple permutation-invariant or equivariant function f:

 Theoretically, it possible to learn f using M, but or even approximate this well in practice.

• Theoretically, it possible to learn f using M, but it is non-trivial to ensure e.g., invariance to orderings,

Question: Wouldn't it be possible to learn properties such as invariance and equivariance from data during training? Why do we need a framework that ensures such properties?

Discussion: Let us loosely denote by M a model that does not have these properties and by M^+ a model that has these properties. We aim to learn a rather simple permutation-invariant or equivariant function f:

- Theoretically, it possible to learn f using M, but or even approximate this well in practice.
- Learning a function f will likely require longer t
 M needs to learn "more".

• Theoretically, it possible to learn f using M, but it is non-trivial to ensure e.g., invariance to orderings,

• Learning a function f will likely require longer training time for model M as compared to model M^+ , as

Question: Wouldn't it be possible to learn properties such as invariance and equivariance from data during training? Why do we need a framework that ensures such properties?

Discussion: Let us loosely denote by M a model that does not have these properties and by M^+ a model that has these properties. We aim to learn a rather simple permutation-invariant or equivariant function f:

- or even approximate this well in practice.
- M needs to learn "more".
- e.g., M likely needs more examples (of orderings) so as to learn invariance to them.

• Theoretically, it possible to learn f using M, but it is non-trivial to ensure e.g., invariance to orderings,

• Learning a function f will likely require longer training time for model M as compared to model M^+ , as

• Learning a function f will likely require more training data for model M as compared to model M^+ , as,

Question: Wouldn't it be possible to learn properties such as invariance and equivariance from data during training? Why do we need a framework that ensures such properties?

Discussion: Let us loosely denote by M a model that does not have these properties and by M^+ a model that has these properties. We aim to learn a rather simple permutation-invariant or equivariant function f:

- or even approximate this well in practice.
- M needs to learn "more".
- e.g., M likely needs more examples (of orderings) so as to learn invariance to them.
- translation-invariant).

• Theoretically, it possible to learn f using M, but it is non-trivial to ensure e.g., invariance to orderings,

• Learning a function f will likely require longer training time for model M as compared to model M^+ , as

• Learning a function f will likely require more training data for model M as compared to model M^+ , as,

• Having the right inductive bias for a domain is very important — This has been observed on other deep learning models, and motivated new architectures (e.g., the use of convolutions which are

Relational Inductive Bias: Investigate on a simple function, i.e., on a concrete task of pathfinding.

Relational Inductive Bias: Investigate on a simple function, i.e., on a concrete task of pathfinding.

Background: (Weston et al., 2015) proposed a collection of proxy tasks (bAbI) that are aimed at evaluating certain reasoning capabilities in the context of question answering. (Li et al., 2016) transformed the "bAbl Task 19", a kind of pathfinding, into a symbolic form, and conducted experiments. In this reformulation, we are given a set of connections:

EsA. BnC. EwF. BwE.

Relational Inductive Bias: Investigate on a simple function, i.e., on a concrete task of pathfinding.

Background: (Weston et al., 2015) proposed a collection of proxy tasks (bAbl) that are aimed at evaluating certain reasoning capabilities in the context of question answering. (Li et al., 2016) transformed the "bAbl Task 19", a kind of pathfinding, into a symbolic form, and conducted experiments. In this reformulation, we are given a set of connections:

where, for instance, "E s A" denotes A is reachable from E by going south. The task is to find a path between, e.g., B and A, such as w,s.

EsA. BnC. EwF. BwE.

Relational Inductive Bias: Investigate on a simple function, i.e., on a concrete task of pathfinding.

Background: (Weston et al., 2015) proposed a collection of proxy tasks (bAbl) that are aimed at evaluating certain reasoning capabilities in the context of question answering. (Li et al., 2016) transformed the "bAbl Task 19", a kind of pathfinding, into a symbolic form, and conducted experiments. In this reformulation, we are given a set of connections:

between, e.g., B and A, such as w,s.

This is a rather simple pathfinding problem on graphs defined over edge types s, n, e, w. Can we learn to predict such paths?

EsA. BnC. EwF. BwE,

where, for instance, "E s A" denotes A is reachable from E by going south. The task is to find a path

Relational Inductive Bias: Investigate on a simple function, i.e., on a concrete task of pathfinding.

Background: (Weston et al., 2015) proposed a collection of proxy tasks (bAbI) that are aimed at evaluating certain reasoning capabilities in the context of question answering. (Li et al., 2016) transformed the "bAbl Task 19", a kind of pathfinding, into a symbolic form, and conducted experiments. In this reformulation, we are given a set of connections:

between, e.g., B and A, such as w,s.

This is a rather simple pathfinding problem on graphs defined over edge types s, n, e, w. Can we learn to predict such paths?

Note, for instance, that one would expect the answers to be the same for isomorphic graph instances: permutation-invariance.

EsA. BnC. EwF. BwE.

where, for instance, "E s A" denotes A is reachable from E by going south. The task is to find a path

Pathfinding: How well do earlier deep learning models, e.g., LSTMs, perform on this problem?

Pathfinding: How well do earlier deep learning models, e.g., LSTMs, perform on this problem?

Results: (Li et al., 2016) reports the empirical results relative to LSTMs and gated graph sequence neural networks (GGSNNs), i.e., a graph neural network model proposed in the same paper:

Pathfinding: How well do earlier deep learning models, e.g., LSTMs, perform on this problem?

Results: (Li et al., 2016) reports the empirical results relative to LSTMs and gated graph sequence neural networks (GGSNNs), i.e., a graph neural network model proposed in the same paper:

LSTM

 28.2 ± 1.3 with 950 training samples

Pathfinding: How well do earlier deep learning models, e.g., LSTMs, perform on this problem?

Results: (Li et al., 2016) reports the empirical results relative to LSTMs and gated graph sequence neural networks (GGSNNs), i.e., a graph neural network model proposed in the same paper:

LSTM	28.2 ±
GGSNN	71.1 ±

- 1.3 with 950 training samples
- 14.7 with 50 training samples
- 92.5 ± 5.9 with 100 training samples
- 99.0 ± 1.1 with 250 training samples
Relational Inductive Bias: Pathfinding

Pathfinding: How well do earlier deep learning models, e.g., LSTMs, perform on this problem?

Results: (Li et al., 2016) reports the empirical results relative to LSTMs and gated graph sequence neural networks (GGSNNs), i.e., a graph neural network model proposed in the same paper:

LSTM	28.2 ±
GGSNN	71.1 ±
	92.5 ±
	990+

The relational inductive bias helps a lot, both in terms of accuracy, and in the #samples needed Similar phenomenons have been observed in other papers and on other tasks, and we will discuss some of these in more detail later in the course.

- 1.3 with 950 training samples
- 14.7 with 50 training samples
- 5.9 with 100 training samples
- 99.0 ± 1.1 with 250 training samples





We consider input graphs G = (V, E) with a feature matrix $\mathbf{X} \in \mathbb{R}^{d \times V_G}$, where the features will be used to initialise the representations of the respective nodes.



initialise the representations of the respective nodes.

These representations are updated with the information received from their respective neighbourhoods (i.e., message passing), and this process continues for k iterations, yielding final node representations.

We consider input graphs G = (V, E) with a feature matrix $\mathbf{X} \in \mathbb{R}^{d \times V_G}$, where the features will be used to



initialise the representations of the respective nodes.

These representations are updated with the information received from their respective neighbourhoods (i.e., message passing), and this process continues for k iterations, yielding final node representations.

(Gilmer et al., 2017) defined a single common framework, called Message Passing Neural Networks (MPNNs), which captures a popular family of graph neural network models proposed in the literature.

We consider input graphs G = (V, E) with a feature matrix $\mathbf{X} \in \mathbb{R}^{d \times V_G}$, where the features will be used to

Initialisation: A message passing neural network defines, for every node $u \in V$, an initial vector representation $\mathbf{h}_{\mu}^{(0)} = \mathbf{x}_{\mu}$.

Initialisation: A message passing neural network defines, for every node $u \in V$, an initial vector representation $\mathbf{h}_{\mu}^{(0)} = \mathbf{x}_{\mu}$.

Message passing: Let us denote the representation for each node $u \in V$ at iteration t by $\mathbf{h}_{u}^{(t)}$. The representation \mathbf{h}_{u} for each node $u \in V$ is then iteratively updated with the information received from its neighbourhood as:

$$\mathbf{h}_{u}^{(t)} = combine^{(t)} \Big(\mathbf{h}_{u}^{(t-1)}, aggregate^{(t)} \big(\big\{ \mathbf{h}_{v}^{(t-1)} \mid v \in N(u) \big\} \big) \Big),$$

where $aggregate^{(t)}$ and $combine^{(t)}$ are arbitrary differentiable functions (i.e., neural networks), and functions are mean, sum, or max.

aggregate^(t) is a permutation-invariant function by construction, since its input is a set. Typical aggregate

Initialisation: A message passing neural network defines, for every node $u \in V$, an initial vector representation $\mathbf{h}_{\mu}^{(0)} = \mathbf{x}_{\mu}$.

Message passing: Let us denote the representation for each node $u \in V$ at iteration t by $\mathbf{h}_{u}^{(t)}$. The representation \mathbf{h}_{u} for each node $u \in V$ is then iteratively updated with the information received from its neighbourhood as:

$$\mathbf{h}_{u}^{(t)} = combine^{(t)} \Big(\mathbf{h}_{u}^{(t-1)}, aggregate^{(t)} \Big(\Big\{ \mathbf{h}_{v}^{(t-1)} \mid v \in N(u) \Big\} \Big) \Big),$$

where $aggregate^{(t)}$ and $combine^{(t)}$ are arbitrary differentiable functions (i.e., neural networks), and functions are mean, sum, or max.

across all iterations, it is sometimes called a homogeneous MPNN.

aggregate^(t) is a permutation-invariant function by construction, since its input is a set. Typical aggregate

Observe that we allow different functions at different iterations: if $aggregate^{(t)}$ and $combine^{(t)}$ are the same

Initialisation: A message passing neural network defines, for every node $u \in V$, an initial vector representation $\mathbf{h}_{\mu}^{(0)} = \mathbf{x}_{\mu}$.

Message passing: Let us denote the representation for each node $u \in V$ at iteration t by $\mathbf{h}_{u}^{(t)}$. The representation \mathbf{h}_{u} for each node $u \in V$ is then iteratively updated with the information received from its neighbourhood as:

$$\mathbf{h}_{u}^{(t)} = combine^{(t)} \Big(\mathbf{h}_{u}^{(t-1)}, aggregate^{(t)} \Big(\Big\{ \mathbf{h}_{v}^{(t-1)} \mid v \in N(u) \Big\} \Big) \Big),$$

where $aggregate^{(t)}$ and $combine^{(t)}$ are arbitrary differentiable functions (i.e., neural networks), and functions are mean, sum, or max.

across all iterations, it is sometimes called a homogeneous MPNN.

Final representation: Upon termination, the final node representations will be denoted as $\mathbf{z}_{u} = \mathbf{h}_{u}^{(k)}$.

aggregate^(t) is a permutation-invariant function by construction, since its input is a set. Typical aggregate

Observe that we allow different functions at different iterations: if $aggregate^{(t)}$ and $combine^{(t)}$ are the same

























The *i*-th iteration is also called the *i*-th layer of the MPNN, since each iteration can be seen as an "unrolling" of the network. As usual, the number of layers defines the depth of the network, and the embedding dimensionality is called the width of the network.

We describe a simple, concrete base model. The idea is again to initialise every node $u \in V$ to a vector representation $\mathbf{h}_{\mu}^{(0)} = \mathbf{x}_{\mu}$. Then, the representation updates are given as:

$$\mathbf{h}_{u}^{(t)} = \sigma \Big(\mathbf{W}_{self}^{(t)} \mathbf{h}_{u}^{(t-1)} + \mathbf{W}_{neigh}^{(t)} \sum_{v \in N(x)} \mathbf{h}_{v}^{(t-1)} + \mathbf{b}^{(t)} \Big),$$

where $W_{self}^{(t)}, W_{neigh}^{(t)} \in \mathbb{R}^{d(t) \times d(t-1)}$ are trainable parameter matrices and σ is an element-wise non-linear function (e.g., ReLU), and $\mathbf{b}^{(t)} \in \mathbb{R}^{d(t)}$ is a bias term (which we will omit in the sequel).

We describe a simple, concrete base model. The idea is again to initialise every node $u \in V$ to a vector representation $\mathbf{h}_{\mu}^{(0)} = \mathbf{x}_{\mu}$. Then, the representation updates are given as:

$$\mathbf{h}_{u}^{(t)} = \sigma \Big(\mathbf{W}_{self}^{(t)} \mathbf{h}_{u}^{(t-1)} + \mathbf{W}_{neigh}^{(t)} \sum_{v \in N(x)} \mathbf{h}_{v}^{(t-1)} + \mathbf{b}^{(t)} \Big),$$

where $\mathbf{W}_{self}^{(t)}, \mathbf{W}_{neigh}^{(t)} \in \mathbb{R}^{d(t) \times d(t-1)}$ are trainable parameter matrices and σ is an element-wise non-linear function (e.g., ReLU), and $\mathbf{b}^{(t)} \in \mathbb{R}^{d(t)}$ is a bias term (which we will omit in the sequel).

Intuitively, this base model can be seen as an instance of message passing neural networks, where the aggregation function is the sum, and the combine function updates node embeddings using a linear combination with an element-wise non-linearity applied at the end:

We describe a simple, concrete base model. The idea is again to initialise every node $u \in V$ to a vector representation $\mathbf{h}_{\mu}^{(0)} = \mathbf{x}_{\mu}$. Then, the representation updates are given as:

$$\mathbf{h}_{u}^{(t)} = \sigma \Big(\mathbf{W}_{self}^{(t)} \mathbf{h}_{u}^{(t-1)} + \mathbf{W}_{neigh}^{(t)} \sum_{v \in N(x)} \mathbf{h}_{v}^{(t-1)} + \mathbf{b}^{(t)} \Big),$$

where $\mathbf{W}_{self}^{(t)}, \mathbf{W}_{neigh}^{(t)} \in \mathbb{R}^{d(t) \times d(t-1)}$ are trainable parameter matrices and σ is an element-wise non-linear function (e.g., ReLU), and $\mathbf{b}^{(t)} \in \mathbb{R}^{d(t)}$ is a bias term (which we will omit in the sequel).

Intuitively, this base model can be seen as an instance of message passing neural networks, where the aggregation function is the sum, and the combine function updates node embeddings using a linear combination with an element-wise non-linearity applied at the end:

$$aggregate^{(t)}\left(\left\{\mathbf{h}_{v}^{(t-1)} \mid v \in N(u)\right\}\right) = \sum_{v \in N(u)} \mathbf{h}_{v}^{(t-1)}$$

We describe a simple, concrete base model. The idea is again to initialise every node $u \in V$ to a vector representation $\mathbf{h}_{\mu}^{(0)} = \mathbf{x}_{\mu}$. Then, the representation updates are given as:

$$\mathbf{h}_{u}^{(t)} = \sigma \Big(\mathbf{W}_{self}^{(t)} \mathbf{h}_{u}^{(t-1)} + \mathbf{W}_{neigh}^{(t)} \sum_{v \in N(x)} \mathbf{h}_{v}^{(t-1)} + \mathbf{b}^{(t)} \Big),$$

where $\mathbf{W}_{self}^{(t)}, \mathbf{W}_{neigh}^{(t)} \in \mathbb{R}^{d(t) \times d(t-1)}$ are trainable parameter matrices and σ is an element-wise non-linear function (e.g., ReLU), and $\mathbf{b}^{(t)} \in \mathbb{R}^{d(t)}$ is a bias term (which we will omit in the sequel).

Intuitively, this base model can be seen as an instance of message passing neural networks, where the aggregation function is the sum, and the combine function updates node embeddings using a linear combination with an element-wise non-linearity applied at the end:

$$aggregate^{(t)}\left(\left\{\mathbf{h}_{v}^{(t-1)} \mid v \in N(u)\right\}\right) = \sum_{v \in N(u)} \mathbf{h}_{v}^{(t-1)}$$

 $combine^{(t)}\left(\mathbf{h}_{u}^{(t-1)}, aggregate^{(t)}\left(\left\{\mathbf{h}_{v}^{(t-1)} \mid v \in N(u)\right\}\right)\right) = \sigma\left(\mathbf{W}_{self}^{(t)}\mathbf{h}_{u}^{(t-1)} + \mathbf{W}_{neigh}^{(t)}aggregate^{(t)}\left(\left\{\mathbf{h}_{v}^{(t-1)} \mid v \in N(u)\right\}\right)\right)$

Conceptually, it might appear more convenient to unify combine and aggregate functions, by (implicitly) adding self-loops to the nodes.

Conceptually, it might appear more convenient to unify combine and aggregate functions, by (implicitly) adding self-loops to the nodes.

That is, we can define message passing solely based on aggregate function as:

$$\mathbf{h}_{u}^{(t)} = aggregate^{(t)} \left(\left\{ \mathbf{h}_{v}^{(t-1)} \right\} \right)$$

 $(-1) | v \in N(u) \} \cup \{\mathbf{h}_{u}^{(t-1)}\}$

Conceptually, it might appear more convenient to unify combine and aggregate functions, by (implicitly) adding self-loops to the nodes.

That is, we can define message passing solely based on aggregate function as:

$$\mathbf{h}_{u}^{(t)} = aggregate^{(t)} \left(\left\{ \mathbf{h}_{v}^{(t-1)} \right\} \right)^{(t-1)}$$

Here, the aggregation is not only over the set the node's neighbours, but also the node itself.

- $(1) \mid v \in N(u) \} \cup \left\{ \mathbf{h}_{u}^{(t-1)} \right\}$

Conceptually, it might appear more convenient to unify combine and aggregate functions, by (implicitly) adding self-loops to the nodes.

That is, we can define message passing solely based on aggregate function as:

$$\mathbf{h}_{u}^{(t)} = aggregate^{(t)} \left(\left\{ \mathbf{h}_{v}^{(t-1)} \right\} \right)$$

Here, the aggregation is not only over the set the node's neighbours, but also the node itself.

This simplifies the message passing. For example, the base model will then simplify to a model where we do not distinguish between trainable matrices any more:

$$\mathbf{h}_{u}^{(t)} = \sigma \left(\mathbf{W}^{(t)} \sum_{v \in N(x)} \mathbf{h}_{v}^{(t-1)} + \mathbf{h}_{u}^{(t-1)} \right)$$

 $(1) \mid v \in N(u) \} \cup \left\{ \mathbf{h}_{u}^{(t-1)} \right\}$

Conceptually, it might appear more convenient to unify combine and aggregate functions, by (implicitly) adding self-loops to the nodes.

That is, we can define message passing solely based on aggregate function as:

$$\mathbf{h}_{u}^{(t)} = aggregate^{(t)} \left(\left\{ \mathbf{h}_{v}^{(t-1)} \right\} \right)$$

Here, the aggregation is not only over the set the node's neighbours, but also the node itself.

This simplifies the message passing. For example, the base model will then simplify to a model where we do not distinguish between trainable matrices any more:

$$\mathbf{h}_{u}^{(t)} = \sigma \left(\mathbf{W}^{(t)} \sum_{v \in N(x)} \mathbf{h}_{v}^{(t-1)} + \mathbf{h}_{u}^{(t-1)} \right)$$

However, this also severely limits the expressivity of the MPNN, as the information coming from the node's neighbours cannot be differentiated from the information from the node itself!

 $(1) \mid v \in N(u) \} \cup \left\{ \mathbf{h}_{u}^{(t-1)} \right\}$

Conceptually, it might appear more convenient to unify combine and aggregate functions, by (implicitly) adding self-loops to the nodes.

That is, we can define message passing solely based on aggregate function as:

$$\mathbf{h}_{u}^{(t)} = aggregate^{(t)} \left(\left\{ \mathbf{h}_{v}^{(t-1)} \right\} \right)$$

Here, the aggregation is not only over the set the node's neighbours, but also the node itself.

This simplifies the message passing. For example, the base model will then simplify to a model where we do not distinguish between trainable matrices any more:

$$\mathbf{h}_{u}^{(t)} = \sigma \left(\mathbf{W}^{(t)} \sum_{v \in N(x)} \mathbf{h}_{v}^{(t-1)} + \mathbf{h}_{u}^{(t-1)} \right)$$

However, this also severely limits the expressivity of the MPNN, as the information coming from the node's neighbours cannot be differentiated from the information from the node itself!

We focus on the message passing approach without self-loops, unless explicitly mentioned otherwise.

 $(1) \mid v \in N(u) \} \cup \left\{ \mathbf{h}_{u}^{(t-1)} \right\}$

Graph Pooling
We have defined a message passing approach to produce a set of node embeddings, but we are also interested in making predictions at the graph level.

We have defined a message passing approach to produce a set of node embeddings, but we are also interested in making predictions at the graph level.

Embedding for the entire graph: This task is often referred to as graph pooling, since we "pool" together the node embeddings in order to learn an embedding of the entire graph (Hamilton, 2020).

We have defined a message passing approach to produce a set of node embeddings, but we are also interested in making predictions at the graph level.

the node embeddings in order to learn an embedding of the entire graph (Hamilton, 2020).

are various ways of defining \mathbf{z}_{G} which eventually lead to different models.

- **Embedding for the entire graph**: This task is often referred to as graph pooling, since we "pool" together
- **Final representation**: For a given graph G, let us denote the final graph representation as $\mathbf{z}_G = \mathbf{h}_G^{(k)}$. There

We have defined a message passing approach to produce a set of node embeddings, but we are also interested in making predictions at the graph level.

the node embeddings in order to learn an embedding of the entire graph (Hamilton, 2020).

are various ways of defining \mathbf{z}_{G} which eventually lead to different models.

that aggregation operates over the local neighbourhood only.

- **Embedding for the entire graph**: This task is often referred to as graph pooling, since we "pool" together
- **Final representation**: For a given graph G, let us denote the final graph representation as $\mathbf{z}_G = \mathbf{h}_G^{(k)}$. There
- Specifically, we need a mapping from the set of all the node embeddings $\{\mathbf{z}_{u_1}...\mathbf{z}_{u_n}\}$ to \mathbf{z}_G . This is very similar to the aggregate function (whose domain is also a set of embeddings) with the only difference being

We have defined a message passing approach to produce a set of node embeddings, but we are also interested in making predictions at the graph level.

the node embeddings in order to learn an embedding of the entire graph (Hamilton, 2020).

are various ways of defining \mathbf{z}_{G} which eventually lead to different models.

that aggregation operates over the local neighbourhood only.

mean, which are then typically normalised with respect to, e.g., the size of the nodes.

- **Embedding for the entire graph**: This task is often referred to as graph pooling, since we "pool" together
- **Final representation**: For a given graph G, let us denote the final graph representation as $\mathbf{z}_G = \mathbf{h}_G^{(k)}$. There
- Specifically, we need a mapping from the set of all the node embeddings $\{\mathbf{z}_{u_1}...\mathbf{z}_{u_n}\}$ to \mathbf{z}_G . This is very similar to the aggregate function (whose domain is also a set of embeddings) with the only difference being
- In principle, any aggregation function can also be used to generate \mathbf{z}_G , and common choices are sum, or

We have defined a message passing approach to produce a set of node embeddings, but we are also interested in making predictions at the graph level.

the node embeddings in order to learn an embedding of the entire graph (Hamilton, 2020).

are various ways of defining \mathbf{z}_{G} which eventually lead to different models.

that aggregation operates over the local neighbourhood only.

mean, which are then typically normalised with respect to, e.g., the size of the nodes.

There are various methods for graph, or more generally, relational pooling (Murphy et al., 2019).

- **Embedding for the entire graph**: This task is often referred to as graph pooling, since we "pool" together
- **Final representation**: For a given graph G, let us denote the final graph representation as $\mathbf{z}_G = \mathbf{h}_G^{(k)}$. There
- Specifically, we need a mapping from the set of all the node embeddings $\{\mathbf{z}_{u_1}...\mathbf{z}_{u_n}\}$ to \mathbf{z}_G . This is very similar to the aggregate function (whose domain is also a set of embeddings) with the only difference being
- In principle, any aggregation function can also be used to generate \mathbf{z}_{G} , and common choices are sum, or

We have defined a message passing approach to produce a set of node embeddings, but we are also interested in making predictions at the graph level.

the node embeddings in order to learn an embedding of the entire graph (Hamilton, 2020).

are various ways of defining \mathbf{z}_{G} which eventually lead to different models.

that aggregation operates over the local neighbourhood only.

mean, which are then typically normalised with respect to, e.g., the size of the nodes.

There are various methods for graph, or more generally, relational pooling (Murphy et al., 2019).

- **Embedding for the entire graph**: This task is often referred to as graph pooling, since we "pool" together
- **Final representation**: For a given graph G, let us denote the final graph representation as $\mathbf{z}_G = \mathbf{h}_G^{(k)}$. There
- Specifically, we need a mapping from the set of all the node embeddings $\{\mathbf{z}_{u_1}...\mathbf{z}_{u_n}\}$ to \mathbf{z}_G . This is very similar to the aggregate function (whose domain is also a set of embeddings) with the only difference being
- In principle, any aggregation function can also be used to generate \mathbf{z}_{G} , and common choices are sum, or



Problem: We have defined MPNNs which learn embeddings of nodes and graphs. The presented message passing approach, however, is local, e.g., for a disconnected graph no information will flow across disjoint subgraphs. This is a serious limitation!

Problem: We have defined MPNNs which learn embeddings of nodes and graphs. The presented message passing approach, however, is local, e.g., for a disconnected graph no information will flow across disjoint subgraphs. This is a serious limitation!

Remark: A graph embedding is necessarily global in the sense that it is composed of all nodes, but it is not a solution to the above-mentioned problem: During message passing there are still no messages between disjoint subgraphs and so the learned node embeddings are "blind" to other embeddings in disjoint subgraphs. The graph embedding, which is based on node embeddings, will necessarily have induced limitations.

Problem: We have defined MPNNs which learn embeddings of nodes and graphs. The presented message passing approach, however, is local, e.g., for a disconnected graph no information will flow across disjoint subgraphs. This is a serious limitation!

Remark: A graph embedding is necessarily global in the sense that it is composed of all nodes, but it is not a solution to the above-mentioned problem: During message passing there are still no messages between disjoint subgraphs and so the learned node embeddings are "blind" to other embeddings in disjoint subgraphs. The graph embedding, which is based on node embeddings, will necessarily have induced limitations.

Solution: To break this behaviour, a standard approach is to use a global feature computation on each layer of the MPNN (Battaglia et al., 2018), also called a global attribute computation, or global readout.

We present a simple extension to MPNNs by allowing global readouts, where in each layer we also compute a feature vector for the whole graph G and combine it with local aggregations.

Message passing with global readout: The representation \mathbf{h}_u for each node $u \in V$ is then iteratively updated with the information received from its neighbourhood as well as a global feature vector as:

$$\mathbf{h}_{u}^{(t)} = combine^{(t)} \Big(\mathbf{h}_{u}^{(t-1)}, aggregate^{(t)} \Big(\Big\{ \mathbf{h}_{v}^{(t-1)} \mid v \in N(u) \Big\} \Big), read^{(t)} \Big(\Big\{ \mathbf{h}_{w}^{(t-1)} \mid w \in G \Big\} \Big) \Big),$$

where $read^{(t)}$ is a differentiable function. Similarly to $aggregate^{(t)}$, $read^{(t)}$ is permutation-invariant by construction, and all aggregate functions are typical candidates also for $read^{(t)}$.

We present a simple extension to MPNNs by allowing global readouts, where in each layer we also compute a feature vector for the whole graph G and combine it with local aggregations.

Message passing with global readout: The representation \mathbf{h}_{u} for each node $u \in V$ is then iteratively updated with the information received from its neighbourhood as well as a global feature vector as:

$$\mathbf{h}_{u}^{(t)} = combine^{(t)} \Big(\mathbf{h}_{u}^{(t-1)}, aggregate^{(t)} \Big(\Big\{ \mathbf{h}_{v}^{(t-1)} \mid v \in N(u) \Big\} \Big), read^{(t)} \Big(\Big\{ \mathbf{h}_{w}^{(t-1)} \mid w \in G \Big\} \Big) \Big),$$

where $read^{(t)}$ is a differentiable function. Similarly to $aggregate^{(t)}$, $read^{(t)}$ is permutation-invariant by construction, and all aggregate functions are typical candidates also for $read^{(t)}$.

(Battaglia et al., 2018) defines a generalised message passing framework for relational reasoning over graph representations, and message passing with global readout can be seen as a special case of this framework.

We present a simple extension to MPNNs by allowing global readouts, where in each layer we also compute a feature vector for the whole graph G and combine it with local aggregations.

Message passing with global readout: The representation \mathbf{h}_{u} for each node $u \in V$ is then iteratively updated with the information received from its neighbourhood as well as a global feature vector as:

$$\mathbf{h}_{u}^{(t)} = combine^{(t)} \Big(\mathbf{h}_{u}^{(t-1)}, aggregate^{(t)} \Big(\Big\{ \mathbf{h}_{v}^{(t-1)} \mid v \in N(u) \Big\} \Big), read^{(t)} \Big(\Big\{ \mathbf{h}_{w}^{(t-1)} \mid w \in G \Big\} \Big) \Big),$$

where $read^{(t)}$ is a differentiable function. Similarly to $aggregate^{(t)}$, $read^{(t)}$ is permutation-invariant by construction, and all aggregate functions are typical candidates also for $read^{(t)}$.

This seemingly simple reformulation makes a significant difference in terms of the expressive power of MPNNs (Barcelo et al., 2020), as we shall see later in the course.

(Battaglia et al., 2018) defines a generalised message passing framework for relational reasoning over graph representations, and message passing with global readout can be seen as a special case of this framework.

Generalised message passing (Battaglia et al., 2018): The main idea is to define a message passing protocol that takes into account the internal representations for edges, nodes, and the graph, respectively:

$$\begin{aligned} \mathbf{h}_{(u,v)}^{(t)} &= combine_{e} \Big(\mathbf{h}_{(u,v)}^{(t-1)}, \mathbf{h}_{u}^{(t-1)}, \mathbf{h}_{v}^{(t-1)}, \mathbf{h}_{G}^{(t-1)} \Big), \\ \mathbf{h}_{u}^{(t)} &= combine_{n} \Big(\mathbf{h}_{u}^{(t-1)}, aggregate^{(t)} \Big(\big\{ \mathbf{h}_{v}^{(t-1)} \mid v \in N(u) \big\} \Big), \mathbf{h}_{G}^{(t-1)} \Big), \\ \mathbf{h}_{G}^{(t)} &= combine_{G} \Big(\mathbf{h}_{G}^{(t-1)}, \big\{ \mathbf{h}_{u}^{(t)} \mid u \in V_{G} \big\}, \big\{ \mathbf{h}_{(u,v)}^{(t)} \mid (u,v) \in E \big\} \Big), \end{aligned}$$

where at each iteration, each update happens in the given equation order.

Generalised message passing (Battaglia et al., 2018): The main idea is to define a message passing protocol that takes into account the internal representations for edges, nodes, and the graph, respectively:

$$\begin{aligned} \mathbf{h}_{(u,v)}^{(t)} &= combine_{e} \Big(\mathbf{h}_{(u,v)}^{(t-1)}, \mathbf{h}_{u}^{(t-1)}, \mathbf{h}_{V}^{(t-1)}, \mathbf{h}_{G}^{(t-1)} \Big), \\ \mathbf{h}_{u}^{(t)} &= combine_{n} \Big(\mathbf{h}_{u}^{(t-1)}, aggregate^{(t)} \Big(\big\{ \mathbf{h}_{v}^{(t-1)} \mid v \in N(u) \big\} \Big), \mathbf{h}_{G}^{(t-1)} \Big), \\ \mathbf{h}_{G}^{(t)} &= combine_{G} \Big(\mathbf{h}_{G}^{(t-1)}, \big\{ \mathbf{h}_{u}^{(t)} \mid u \in V_{G} \big\}, \big\{ \mathbf{h}_{(u,v)}^{(t)} \mid (u,v) \in E \big\} \Big), \end{aligned}$$

where at each iteration, each update happens in the given equation order.

The idea is to generate hidden embeddings $\mathbf{h}_{(u,v)}$ for each edge (u,v) in the graph G, as well as an embedding \mathbf{h}_G corresponding to the entire graph.

Generalised message passing (Battaglia et al., 2018): The main idea is to define a message passing protocol that takes into account the internal representations for edges, nodes, and the graph, respectively:

$$\begin{split} \mathbf{h}_{(u,v)}^{(t)} &= combine_{e} \Big(\mathbf{h}_{(u,v)}^{(t-1)}, \mathbf{h}_{u}^{(t-1)}, \mathbf{h}_{v}^{(t-1)}, \mathbf{h}_{G}^{(t-1)} \Big), \\ \mathbf{h}_{u}^{(t)} &= combine_{n} \Big(\mathbf{h}_{u}^{(t-1)}, aggregate^{(t)} \Big(\big\{ \mathbf{h}_{v}^{(t-1)} \mid v \in N(u) \big\} \Big), \mathbf{h}_{G}^{(t-1)} \Big), \\ \mathbf{h}_{G}^{(t)} &= combine_{G} \Big(\mathbf{h}_{G}^{(t-1)}, \big\{ \mathbf{h}_{u}^{(t)} \mid u \in V_{G} \big\}, \big\{ \mathbf{h}_{(u,v)}^{(t)} \mid (u,v) \in E \big\} \Big), \end{split}$$

where at each iteration, each update happens in the given equation order.

The idea is to generate hidden embeddings $\mathbf{h}_{(u,v)}$ for each edge (u, v) in the graph G, as well as an embedding \mathbf{h}_G corresponding to the entire graph.

This allows the message passing model to easily integrate edge and graph-level features, especially in the multi-relational context, where each edge can be labelled differently.

Graph Representation Learning Tasks

graph, i.e., y_v for all $v \in V \setminus V_{tr}$.

Task: Given a single graph G = (V, E) with a feature matrix $\mathbf{X} \in \mathbb{R}^{d \times V_G}$, where a subset of the nodes $\{(u, y_u) \mid u \in V_{tr} \subset V\}$ are labeled with a class, predict the labels of the remaining nodes, i.e., test nodes in the

Task: Given a single graph G = (V, E) with a feature matrix $\mathbf{X} \in \mathbb{R}^{d \times V_G}$, where a subset of the nodes $\{(u, y_u) \mid u \in V_{tr} \subset V\}$ are labeled with a class, predict the labels of the remaining nodes, i.e., test nodes in the graph, i.e., y_v for all $v \in V \setminus V_{tr}$.

ML). The task is to predict the category (or, categories) of the remaining papers.

Example (Kipf and Welling, 2017): Consider citation networks such as Citeseer, where nodes represent papers, and edges denote citation links, and a subset of the nodes/papers are labelled with a paper category (e.g., AI,

Task: Given a single graph G = (V, E) with a feature matrix $\mathbf{X} \in \mathbb{R}^{d \times V_G}$, where a subset of the nodes $\{(u, y_u) \mid u \in V_{tr} \subset V\}$ are labeled with a class, predict the labels of the remaining nodes, i.e., test nodes in the graph, i.e., y_v for all $v \in V \setminus V_{tr}$.

ML). The task is to predict the category (or, categories) of the remaining papers.

Training: Based on an appropriate loss function, e.g., negative log-likelihood loss:

$$\mathscr{L} = \sum_{u \in V_{tr}} -\log(\operatorname{softr}$$

where \mathbf{y}_u is a one-hot vector indicating the class y_u of the training node $u \in V_{tr}$, and the softmax $(\mathbf{z}_u, \mathbf{y}_u)$ function denotes the predicted probability that the node u belongs to the class y_{μ} .

As usual, the gradient of the loss is backpropagated through the parameters of the GNN using stochastic gradient descent or one of its variants.

Example (Kipf and Welling, 2017): Consider citation networks such as Citeseer, where nodes represent papers, and edges denote citation links, and a subset of the nodes/papers are labelled with a paper category (e.g., AI,

 $\max(\mathbf{z}_{u}, \mathbf{y}_{u})),$

Test nodes: The given loss function uses only nodes from the training set V_{tr} , but test nodes $V_{test} = V \setminus V_{tr}$ may still be observed during training, i.e., they can take part in the message passing.

may still be observed during training, i.e., they can take part in the message passing.

This yields two regimes for learning (Hamilton, 2020):

Test nodes: The given loss function uses only nodes from the training set V_{tr} , but test nodes $V_{test} = V \setminus V_{tr}$

may still be observed during training, i.e., they can take part in the message passing.

This yields two regimes for learning (Hamilton, 2020):

nodes.

Test nodes: The given loss function uses only nodes from the training set V_{tr} , but test nodes $V_{test} = V \setminus V_{tr}$

• Transductive: In this regime all nodes, including test nodes, are observed during training, i.e., their representations are computed during message passing and they also affect the representation of other

may still be observed during training, i.e., they can take part in the message passing.

This yields two regimes for learning (Hamilton, 2020):

- nodes.

Test nodes: The given loss function uses only nodes from the training set V_{tr} , but test nodes $V_{test} = V \setminus V_{tr}$

• Transductive: In this regime all nodes, including test nodes, are observed during training, i.e., their representations are computed during message passing and they also affect the representation of other

• Inductive: In this regime, not all test nodes are observed during training. This means that for some test nodes, neither the nodes themselves nor their edges (and hence their relation to other nodes) are known.

may still be observed during training, i.e., they can take part in the message passing.

This yields two regimes for learning (Hamilton, 2020):

- nodes.

Transductive classification example: Suppose we have access to the full citation graph at training time in the citation networks task, and define a subset of the nodes as the test nodes. This is a transductive task, as the structural information of the nodes is available, and is used during the message passing.

Test nodes: The given loss function uses only nodes from the training set V_{tr} , but test nodes $V_{test} = V \setminus V_{tr}$

• Transductive: In this regime all nodes, including test nodes, are observed during training, i.e., their representations are computed during message passing and they also affect the representation of other

• Inductive: In this regime, not all test nodes are observed during training. This means that for some test nodes, neither the nodes themselves nor their edges (and hence their relation to other nodes) are known.

may still be observed during training, i.e., they can take part in the message passing.

This yields two regimes for learning (Hamilton, 2020):

- nodes.

Transductive classification example: Suppose we have access to the full citation graph at training time in the citation networks task, and define a subset of the nodes as the test nodes. This is a transductive task, as the structural information of the nodes is available, and is used during the message passing.

Inductive classification example: Suppose we only see a subgraph of the citation graph at training time in the citation networks task, and we define a set of nodes from a disjoint subgraph as test nodes. This is an inductive task, as we have not observed the test nodes during training.

Test nodes: The given loss function uses only nodes from the training set V_{tr} , but test nodes $V_{test} = V \setminus V_{tr}$

• Transductive: In this regime all nodes, including test nodes, are observed during training, i.e., their representations are computed during message passing and they also affect the representation of other

• Inductive: In this regime, not all test nodes are observed during training. This means that for some test nodes, neither the nodes themselves nor their edges (and hence their relation to other nodes) are known.



Node classification is usually viewed as a semi-supervised learning task, and this should be understood in reference to transductive node classification. It is called semi-supervised since:



Node classification is usually viewed as a semi-supervised learning task, and this should be understood in reference to transductive node classification. It is called semi-supervised since:

• We train using training labels, as standard in supervised learning



reference to transductive node classification. It is called semi-supervised since:

- We train using training labels, as standard in supervised learning
- We additionally have access to the structural information of unlabelled test nodes

Node classification is usually viewed as a semi-supervised learning task, and this should be understood in


reference to transductive node classification. It is called semi-supervised since:

- We train using training labels, as standard in supervised learning
- We additionally have access to the structural information of unlabelled test nodes

This means that we trained using both labelled and unlabelled data.



reference to transductive node classification. It is called semi-supervised since:

- We train using training labels, as standard in supervised learning
- We additionally have access to the structural information of unlabelled test nodes

This means that we trained using both labelled and unlabelled data.

requires i.i.d. assumption, which doesn't hold for node classification).

- In transductive node classification, which is by far the most common node classification task on graphs, we typically have access to the full graph, including all the test nodes (and, e.g., their neighbourhood), and this makes the term semi-supervised somewhat more appropriate (though, the standard semi-supervised setting also



reference to transductive node classification. It is called semi-supervised since:

- We train using training labels, as standard in supervised learning
- We additionally have access to the structural information of unlabelled test nodes

This means that we trained using both labelled and unlabelled data.

requires i.i.d. assumption, which doesn't hold for node classification).

Inductive node classification, on the other hand, can be seen as supervised learning.

- In transductive node classification, which is by far the most common node classification task on graphs, we typically have access to the full graph, including all the test nodes (and, e.g., their neighbourhood), and this makes the term semi-supervised somewhat more appropriate (though, the standard semi-supervised setting also



reference to transductive node classification. It is called semi-supervised since:

- We train using training labels, as standard in supervised learning
- We additionally have access to the structural information of unlabelled test nodes

This means that we trained using both labelled and unlabelled data.

requires i.i.d. assumption, which doesn't hold for node classification).

Inductive node classification, on the other hand, can be seen as supervised learning.

and terminology in machine learning.

- In transductive node classification, which is by far the most common node classification task on graphs, we typically have access to the full graph, including all the test nodes (and, e.g., their neighbourhood), and this makes the term semi-supervised somewhat more appropriate (though, the standard semi-supervised setting also
- This can be noted as another instance where machine learning on graphs tends to differ from standard practices



Task: Given a set of graphs $\mathscr{G} = \{G_1, \dots, G_n\}$, where a subset of the graphs $\{(G_i, y_i) \mid G_i \in \mathscr{G}_{tr} \subset \mathscr{G}\}$ are labeled with a class, predict the labels of the remaining (test) graphs, i.e., y_j for all $G_j \in \mathscr{G} \setminus \mathscr{G}_{tr}$.

labeled with a class, predict the labels of the remaining (test) graphs, i.e., y_i for all $G_i \in \mathscr{G} \setminus \mathscr{G}_{tr}$.

Example: The IMDB datasets (Morris et al., 2020) consist of the so-called ego-networks for each movie, and contains information such as actor collaborations for each movie. The task is to predict the genre (e.g., action, horror) of the movie.

Task: Given a set of graphs $\mathcal{G} = \{G_1, \dots, G_n\}$, where a subset of the graphs $\{(G_i, y_i) \mid G_i \in \mathcal{G}_{tr} \subset \mathcal{G}\}$ are

Task: Given a set of graphs $\mathcal{G} = \{G_1, \dots, G_n\}$, where a subset of the graphs $\{(G_i, y_i) \mid G_i \in \mathcal{G}_{tr} \subset \mathcal{G}\}$ are labeled with a class, predict the labels of the remaining (test) graphs, i.e., y_i for all $G_i \in \mathscr{G} \setminus \mathscr{G}_{tr}$.

horror) of the movie.

Training is similar to node classification, except that we use the final embedding of the graph instead of nodes.

Example: The IMDB datasets (Morris et al., 2020) consist of the so-called ego-networks for each movie, and contains information such as actor collaborations for each movie. The task is to predict the genre (e.g., action,

labeled with a class, predict the labels of the remaining (test) graphs, i.e., y_i for all $G_i \in \mathcal{G} \setminus \mathcal{G}_{tr}$.

horror) of the movie.

label, and the goal is to use a labeled set of graphs to learn a mapping from graphs to class labels.

- **Task**: Given a set of graphs $\mathscr{G} = \{G_1, \dots, G_n\}$, where a subset of the graphs $\{(G_i, y_i) \mid G_i \in \mathscr{G}_{tr} \subset \mathscr{G}\}$ are
- **Example**: The IMDB datasets (Morris et al., 2020) consist of the so-called ego-networks for each movie, and contains information such as actor collaborations for each movie. The task is to predict the genre (e.g., action,
- Training is similar to node classification, except that we use the final embedding of the graph instead of nodes.
- Note that graph classification is a supervised learning task: Each graph is an i.i.d. data point associated with a

Intuitively, graph neural networks can be seen as encoders, which, for a given graph (or set of graphs), learn an embedding for each node (respectively, each graph).

Intuitively, graph neural networks can be seen as encoders, which, for a given graph (or set of graphs), learn an embedding for each node (respectively, each graph).

Taking this broader perspective, we have the liberty to choose how to use the learned embeddings, i.e., they can be used for any standard machine learning task, e.g., classification, regression, clustering, etc. and so they have natural counterparts both on node and graph-level tasks.

Intuitively, graph neural networks can be seen as encoders, which, for a given graph (or set of graphs), learn an embedding for each node (respectively, each graph).

Taking this broader perspective, we have the liberty to choose how to use the learned embeddings, i.e., they can be used for any standard machine learning task, e.g., classification, regression, clustering, etc. and so they have natural counterparts both on node and graph-level tasks.

Example (Gilmer et al., 2017): Given a graph representing the structure of a molecule, an interesting application is to build a regression model that could predict that molecule's toxicity or solubility, which is an application of graph regression.

Intuitively, graph neural networks can be seen as encoders, which, for a given graph (or set of graphs), learn an embedding for each node (respectively, each graph).

Taking this broader perspective, we have the liberty to choose how to use the learned embeddings, i.e., they can be used for any standard machine learning task, e.g., classification, regression, clustering, etc. and so they have natural counterparts both on node and graph-level tasks.

Example (Gilmer et al., 2017): Given a graph representing the structure of a molecule, an interesting application is to build a regression model that could predict that molecule's toxicity or solubility, which is an application of graph regression.

Example: Similarly, if we are interested in applications such as community detection, this can be framed as a node clustering problem. In a similar vein, we may want to cluster a set of graphs, i.e., graph clustering.

Intuitively, graph neural networks can be seen as encoders, which, for a given graph (or set of graphs), learn an embedding for each node (respectively, each graph).

Taking this broader perspective, we have the liberty to choose how to use the learned embeddings, i.e., they can be used for any standard machine learning task, e.g., classification, regression, clustering, etc. and so they have natural counterparts both on node and graph-level tasks.

Example (Gilmer et al., 2017): Given a graph representing the structure of a molecule, an interesting application is to build a regression model that could predict that molecule's toxicity or solubility, which is an application of graph regression.

Example: Similarly, if we are interested in applications such as community detection, this can be framed as a node clustering problem. In a similar vein, we may want to cluster a set of graphs, i.e., graph clustering.

Graph clustering is an unsupervised learning task, a natural extension of standard clustering.

Perspectives on Graph Neural Networks

Convolutional neural networks (CNNs) have been very successful in machine learning problems on gridstructured data, especially for data that requires translational equivariance/invariance with respect to this grid.

Convolutional neural networks (CNNs) have been very successful in machine learning problems on grid-

generalise Euclidean convolutions to the graph domain (Bruna et al., 2014)!

- structured data, especially for data that requires translational equivariance/invariance with respect to this grid.
- This is one of the historical perspectives behind the development of graph neural networks: The main idea is to

Convolutional neural networks (CNNs) have been very successful in machine learning problems on grid-

generalise Euclidean convolutions to the graph domain (Bruna et al., 2014)!

Welling, 2016), which also falls into the class of MPNNs.

- structured data, especially for data that requires translational equivariance/invariance with respect to this grid.
- This is one of the historical perspectives behind the development of graph neural networks: The main idea is to
- This line of research has led to graph convolutional networks (GCNs), a very popular GNN model (Kipf and

Convolutional neural networks (CNNs) have been very successful in machine learning problems on grid-

generalise Euclidean convolutions to the graph domain (Bruna et al., 2014)!

Welling, 2016), which also falls into the class of MPNNs.

generalise the convolution operator to graph domains.

- structured data, especially for data that requires translational equivariance/invariance with respect to this grid.
- This is one of the historical perspectives behind the development of graph neural networks: The main idea is to
- This line of research has led to graph convolutional networks (GCNs), a very popular GNN model (Kipf and
- Intuitively, the graph Laplacian allows a natural link between graphs and vector spaces, and so the structure of the graph can be exploited with the spectrum of its graph-Laplacian. The main idea is then to use this to

Convolutional neural networks (CNNs) have been very successful in machine learning problems on grid-

generalise Euclidean convolutions to the graph domain (Bruna et al., 2014)!

Welling, 2016), which also falls into the class of MPNNs.

generalise the convolution operator to graph domains.

We will discuss and introduce graph convolutional networks in Lecture 4!

- structured data, especially for data that requires translational equivariance/invariance with respect to this grid.
- This is one of the historical perspectives behind the development of graph neural networks: The main idea is to
- This line of research has led to graph convolutional networks (GCNs), a very popular GNN model (Kipf and
- Intuitively, the graph Laplacian allows a natural link between graphs and vector spaces, and so the structure of the graph can be exploited with the spectrum of its graph-Laplacian. The main idea is then to use this to

Another historical perspective is offered through the connection to the fundamental problem of graph isomorphism.

Another historical perspective is offered through the connection to the fundamental problem of graph isomorphism.

networks cannot distinguish all graphs, and so they have limited expressive power!

- Intuitively, classifying graph-structured data requires the ability to distinguish graphs; and message passing neural

Another historical perspective is offered through the connection to the fundamental problem of graph isomorphism.

networks cannot distinguish all graphs, and so they have limited expressive power!

The connection to graph isomorphism testing offers many interesting and theoretical insights, and has power (Morris et al., 2018, Xu et al, 2019).

- Intuitively, classifying graph-structured data requires the ability to distinguish graphs; and message passing neural
- contributed a great deal to our understanding of graph neural networks, especially in terms of their expressive

Another historical perspective is offered through the connection to the fundamental problem of graph isomorphism.

networks cannot distinguish all graphs, and so they have limited expressive power!

The connection to graph isomorphism testing offers many interesting and theoretical insights, and has power (Morris et al., 2018, Xu et al, 2019).

We will cover these aspects in Lecture 5.

- Intuitively, classifying graph-structured data requires the ability to distinguish graphs; and message passing neural
- contributed a great deal to our understanding of graph neural networks, especially in terms of their expressive

It is also possible to draw connections from traditional machine learning models, i.e., through probabilistic graphical models.

It is also possible to draw connections from traditional machine learning models, i.e., through probabilistic graphical models.

Message passing algorithms are used in the context of probabilistic graphical models, and can be dated back to belief propagation in Bayesian networks (Pearl, 82).

It is also possible to draw connections from traditional machine learning models, i.e., through probabilistic graphical models.

Message passing algorithms are used in the context of probabilistic graphical models, and can be dated back to belief propagation in Bayesian networks (Pearl, 82).

(Dai et al., 2016) offers an alternative perspective for graph neural networks based on probabilistic graphical models: The message passing operation that underlies MPNNs can be viewed as a neural network analogue of certain message passing algorithms that are commonly used for variational inference in probabilistic models to infer distributions over latent variables.

It is also possible to draw connections from traditional machine learning models, i.e., through probabilistic graphical models.

Message passing algorithms are used in the context of probabilistic graphical models, and can be dated back to belief propagation in Bayesian networks (Pearl, 82).

(Dai et al., 2016) offers an alternative perspective for graph neural networks based on probabilistic graphical models: The message passing operation that underlies MPNNs can be viewed as a neural network analogue of certain message passing algorithms that are commonly used for variational inference in probabilistic models to infer distributions over latent variables.

We will not cover this aspect, but you will encounter, e.g., probabilistic graphical models, or variational inference as part of the Bayesian machine learning lectures.

Relational inductive bias is crucial

- Relational inductive bias is crucial
- Message passing neural networks as a framework

- Relational inductive bias is crucial
- Message passing neural networks as a framework
- A basic message passing neural network model and its extensions
- Relational inductive bias is crucial
- Message passing neural networks as a framework
- A basic message passing neural network model and its extensions
- Node and graph classification and beyond

- Relational inductive bias is crucial
- Message passing neural networks as a framework
- A basic message passing neural network model and its extensions
- Node and graph classification and beyond
- Historical perspectives for graph neural networks

- Relational inductive bias is crucial
- Message passing neural networks as a framework
- A basic message passing neural network model and its extensions
- Node and graph classification and beyond
- Historical perspectives for graph neural networks
- We have not covered any concrete model beyond the very basic one: Lecture 4!

- Relational inductive bias is crucial
- Message passing neural networks as a framework
- A basic message passing neural network model and its extensions
- Node and graph classification and beyond
- Historical perspectives for graph neural networks
- We have not covered any concrete model beyond the very basic one: Lecture 4!

• Additional reading material: This lecture is partially based on Chapters 5 - 7 of (Hamilton, 2020)

References

- Euclidean data, IEEE Signal Processing Magazine, 2017.
- N. Kriege, F. Johansson, and C. Morris. A survey on graph kernels. *Appl. Netw. Sci.*, 2020.
- reduction. Nauchno-Technicheskaya Informatsia, 1968.
- toy tasks. arXiv preprint arXiv:1502.05698, 2015.
- chemistry. ICML, 2017.
- datasets for learning with graphs. ICML Workshop on Graph Representation Learning and Beyond, 2020.

• M. M. Bronstein, J. Bruna, Y. LeCun, A. Szlam, P. Vandergheynst, Geometric deep learning: going beyond

• B. Weisfeiler and A. Lehman. A reduction of a graph to a canonical form and an algebra arising during this

• J. Weston, A. Bordes, S. Chopra, and T. Mikolov. Towards Al-complete question answering: a set of prerequisite

• J. Gilmer, S.S. Schoenholz, P. F. Riley, O. Vinyals, and G. E. Dahl. 2017. Neural message passing for Quantum

• C. Morris, N. M. Kriege, F. Bause, K. Kersting, P. Mutzel, and M. Neumann. Tudataset: A collection of benchmark

• R.L.Murphy, B. Srinivasan, V.A. Rao, and B. Ribeiro, Relational Pooling for Graph Representations. ICML, 2019.

References

- 1806.01261, 2018.
- quantum chemistry. ICML, 2017.
- *ICLR*, 2014.
- Y. Li, D. Tarlow, M. Brockschmidt, and R.S. Zemel, Gated graph sequence neural networks. *ICLR*, 2016.

• Peter W. Battaglia, Jessica B. Hamrick, Victor Bapst, Alvaro Sanchez-Gonzalez, Vin ícius Flores Zambaldi, Mateusz Malinowski, Andrea Tacchetti, David Raposo, Adam Santoro, Ryan Faulkner, C, aglar Gulc, ehre, H. Francis Song, Andrew J. Ballard, Justin Gilmer, George E. Dahl, Ashish "Vaswani, Kelsey R. Allen, Charles Nash, Victoria Langston, Chris Dyer, Nicolas Heess, Daan Wierstra, Pushmeet Kohli, Matthew Botvinick, Oriol Vinyals, Yujia Li, and Razvan Pascanu. Relational inductive biases, deep learning, and graph networks. CoRR, abs/

• Neural MP: J. Gilmer, S. S. Schoenholz, P. F. Riley, O. Vinyals, and G. E. Dahl, Neural message passing for

• J. Bruna, W. Zaremba, A. Szlam, Y. LeCun. Spectral Networks and Locally Connected Networks on Graphs.

• T. Kipf and M. Welling, Semi-supervised classification with graph convolutional networks. ICLR, 2017.

References

- C. Morris, M. Ritzert, M. Fey, W. Hamilton, J. E. Lenssen, G. Rattan, and M. Grohe, Weisfeiler and Leman go neural: Higher-order graph neural networks. *AAAI*, 2018.
- K. Xu, W. Hu, J. Leskovec, and S. Jegelka. How powerful are graph neural networks? ICLR, 2019.
- P. Barcelo, E. Kostylev, M. Monet, J. Perez, J. Reutter, and J. Silva. The logical expressiveness of graph neural networks. *ICLR*, 2020.
- Pearl, Judea. Reverend Bayes on inference engines: A distributed hierarchical approach. AAAI, 1982.