

The background of the slide features a 3D visualization of a graph structure. It consists of numerous nodes and edges, with some nodes highlighted in a light blue color. The graph is overlaid on a background that appears to be a globe or a similar spherical object, with a warm, golden-yellow glow emanating from the center. The overall aesthetic is technical and modern.

A Query Language Perspective on Graph Learning

Floris Geerts (University of Antwerp)

A small graphic element in the bottom right corner of the slide, consisting of a horizontal row of five small, colored squares: yellow, cyan, blue, green, and red.

- ▶ Database (DB) theoreticians love graphs.

- ▶ Database (DB) theoreticians love graphs.
- ▶ But so do machine learners (ML).

- ▶ Database (DB) theoreticians love graphs.
- ▶ But so do machine learners (ML).

How to understand what ML folks are doing with graphs from a DB perspective?

- ▶ Graph learning methods can be **expressed** in specialized **graph embedding languages**.

- ▶ Graph learning methods can be **expressed** in specialized **graph embedding languages**.
- ▶ These languages can be **analyzed** with regards to **expressive power** using familiar DB techniques.

- ▶ Graph learning methods can be **expressed** in specialized **graph embedding languages**.
 - ▶ These languages can be **analyzed** with regards to **expressive power** using familiar DB techniques.
- ▶ This results in a **better understanding** of graph learning methods; and
 - ▶ forms a **bridge** between graph learning and DB theory.

Previous keynotes at PODS on graph learning:

- ▶ word2vec, node2vec, graph2vec, X2vec: Towards a Theory of Vector Embeddings of Structured Data, by [Martin Grohe](#) (PODS 2020)
- ▶ Databases as Graphs: Predictive Queries for Declarative Machine Learning by [Jure Leskovec](#).

I nevertheless hope to convey some **alternative point of view**.

1. Preliminaries

Graphs, embeddings and graph learning

Graphs: One definition to rule them all

- ▶ A **graph** $G = (V_G, E_G, L_G)$ with **vertex** set V_G , **edge** set $E_G \subseteq V_G \times V_G$, and **vertex labelling** $L_G : V_G \rightarrow \Sigma$ for some set Σ of **labels**.
- ▶ We often assume $\Sigma = \mathbb{R}^d$ for some dimension $d \in \mathbb{N}$.

Finite set of labels \mapsto **hot-one encoding**, e.g., labels a , b and c : $\left(\underbrace{1}_a, \underbrace{0}_b, \underbrace{0}_c \right)$.



Event Graphs



Computer Networks



Disease Pathways



Social Networks



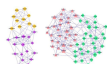
Economic Networks



Communication Networks



Knowledge Graphs



Regulatory Networks



Scene Graphs



Food Webs



Particle Networks



Underground Networks



Citation Networks



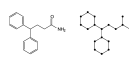
Internet



Networks of Neurons



Code Graphs



Molecules



3D Shapes

In this tutorial, **graph learning** is about learning (partially) unknown graph embeddings.

- ▶ Let \mathcal{G} be the class of **all graphs**.
- ▶ Let \mathbb{Y} be an **output space**.
- ▶ A **graph embedding** is a function of the form

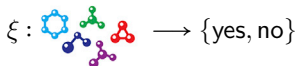
$$\xi : \mathcal{G} \rightarrow \mathbb{Y}.$$

In this tutorial, **graph learning** is about learning (partially) unknown graph embeddings.

- ▶ Let \mathcal{G} be the class of **all graphs**.
- ▶ Let \mathbb{Y} be an **output space**.
- ▶ A **graph embedding** is a function of the form

$$\xi : \mathcal{G} \rightarrow \mathbb{Y}.$$

For example: **prediction** of chemical/medical property of molecules



Also, graph learning is about learning **unknown vertex embeddings**.

- ▶ Let \mathcal{G} be the class of all graphs.
- ▶ Let \mathcal{V} be the class of **all vertices**.
- ▶ Let \mathbb{Y} be an output space.
- ▶ A **vertex embedding** is a function of the form

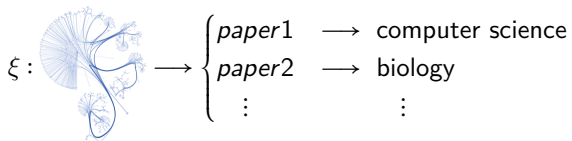
$$\xi : \mathcal{G} \rightarrow (\mathcal{V} \rightarrow \mathbb{Y}).$$

Also, graph learning is about learning **unknown vertex embeddings**.

- ▶ Let \mathcal{G} be the class of all graphs.
- ▶ Let \mathcal{V} be the class of **all vertices**.
- ▶ Let \mathbb{Y} be an output space.
- ▶ A **vertex embedding** is a function of the form

$$\xi : \mathcal{G} \rightarrow (\mathcal{V} \rightarrow \mathbb{Y}).$$

For example: **prediction of the subject of papers** in citation network



More generally, graph learning is about learning **unknown p -vertex embeddings**.

- ▶ Let \mathcal{G} be the class of all graphs.
- ▶ Let \mathcal{V} be the class of all vertices.
- ▶ Let \mathbb{Y} be an output space.
- ▶ A **p -vertex embedding** is a function of the form

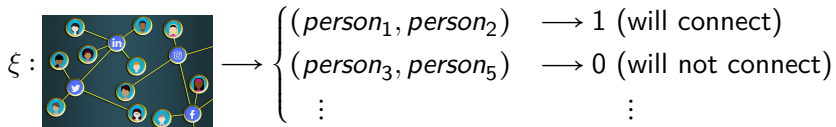
$$\xi : \mathcal{G} \rightarrow (\mathcal{V}^p \rightarrow \mathbb{Y}).$$

More generally, graph learning is about learning **unknown p -vertex embeddings**.

- ▶ Let \mathcal{G} be the class of all graphs.
- ▶ Let \mathcal{V} be the class of all vertices.
- ▶ Let \mathbb{Y} be an output space.
- ▶ A **p -vertex embedding** is a function of the form

$$\xi : \mathcal{G} \rightarrow (\mathcal{V}^p \rightarrow \mathbb{Y}).$$

For example: **Link prediction** in social networks ($p = 2$)



Embeddings will be at the core of this tutorial.

An **important** requirement is that embeddings should be **invariant**, i.e., **independent** of the **chosen graph representation**.

- ▶ A p -vertex embedding ξ is called **invariant** if for any two graphs G and H in \mathcal{G} , for any **graph isomorphism** $\pi : V_G \rightarrow V_H$ from G to H and any p -tuple of vertices $\mathbf{v} \in V_G^p$

$$\xi(G, \mathbf{v}) = \xi(\pi(G), \pi(\mathbf{v})).$$

- ▶ Similar to the **genericity** requirement for query languages.

An **important** requirement is that embeddings should be **invariant**, i.e., **independent** of the **chosen graph representation**.

- ▶ A p -vertex embedding ξ is called **invariant** if for any two graphs G and H in \mathcal{G} , for any **graph isomorphism** $\pi : V_G \rightarrow V_H$ from G to H and any p -tuple of vertices $\mathbf{v} \in V_G^p$

$$\xi(G, \mathbf{v}) = \xi(\pi(G), \pi(\mathbf{v})).$$

- ▶ Similar to the **genericity** requirement for query languages.

How are embedding methods specified?

- ▶ In the ML community, embedding methods are described by their **implementations** using **linear algebra** and other **computations on real numbers**.
- ▶ Crucially, these implementation have **learnable parameters/weights**.
- ▶ Typically, embeddings are defined **layer-wise** (deep architectures).

- ▶ Let σ a **non-linear activation function** $\mathbb{R} \rightarrow \mathbb{R}$ (ReLU, sigmoid, sign,...).
- ▶ Vertex set of graph identified with $[n] := \{1, 2, \dots, n\}$ for some $n \in \mathbb{N}$.
- ▶ Output space $\mathbb{Y} = \mathbb{R}^d$ for some $d \in \mathbb{N}$.
- ▶ Matrix $\mathbf{F}^{(t)}$ in $\mathbb{R}^{n \times d}$ represents **vertex feature** computed in **layer t** .
- ▶ In particular, $\mathbf{F}_{v \bullet}^{(t)}$ in $\mathbb{R}^{1 \times d}$ denotes **embedding of vertex v** .
- ▶ “learnable” **weight matrices** $\mathbf{W}_1^{(t)}, \mathbf{W}_2^{(t)} \in \mathbb{R}^{d \times d}$ and **bias** $\mathbf{b}^{(t)} \in \mathbb{R}^{1 \times d}$.

$$\mathbf{F}_{v \bullet}^{(0)} := L_G(v) \quad \mathbf{F}_{v \bullet}^{(t)} := \sigma \left(\mathbf{F}_{v \bullet}^{(t-1)} \mathbf{W}_1^{(t)} + \sum_{u \in N_G(v)} \mathbf{F}_{u \bullet}^{(t-1)} \mathbf{W}_2^{(t)} + \mathbf{b}^{(t)} \right)$$

- ▶ By varying weights and biases, an **infinite family** of vertex embeddings is obtained.

- ▶ Let σ a **non-linear activation function** $\mathbb{R} \rightarrow \mathbb{R}$ (ReLU, sigmoid, sign,...).
- ▶ Vertex set of graph identified with $[n] := \{1, 2, \dots, n\}$ for some $n \in \mathbb{N}$.
- ▶ Output space $\mathbb{Y} = \mathbb{R}^d$ for some $d \in \mathbb{N}$.
- ▶ Matrix $\mathbf{F}^{(t)}$ in $\mathbb{R}^{n \times d}$ represents **vertex feature** computed in **layer t** .
- ▶ In particular, $\mathbf{F}_{v \bullet}^{(t)}$ in $\mathbb{R}^{1 \times d}$ denotes **embedding of vertex v** .
- ▶ “learnable” **weight matrices** $\mathbf{W}_1^{(t)}, \mathbf{W}_2^{(t)} \in \mathbb{R}^{d \times d}$ and **bias** $\mathbf{b}^{(t)} \in \mathbb{R}^{1 \times d}$.

$$\mathbf{F}_{v \bullet}^{(0)} := L_G(v) \quad \mathbf{F}_{v \bullet}^{(t)} := \sigma \left(\mathbf{F}_{v \bullet}^{(t-1)} \mathbf{W}_1^{(t)} + \sum_{u \in N_G(v)} \mathbf{F}_{u \bullet}^{(t-1)} \mathbf{W}_2^{(t)} + \mathbf{b}^{(t)} \right)$$

- ▶ By varying weights and biases, an **infinite family** of vertex embeddings is obtained.

- ▶ Let σ a **non-linear activation function** $\mathbb{R} \rightarrow \mathbb{R}$ (ReLU, sigmoid, sign,...).
- ▶ Vertex set of graph identified with $[n] := \{1, 2, \dots, n\}$ for some $n \in \mathbb{N}$.
- ▶ Output space $\mathbb{Y} = \mathbb{R}^d$ for some $d \in \mathbb{N}$.
- ▶ Matrix $\mathbf{F}^{(t)}$ in $\mathbb{R}^{n \times d}$ represents **vertex feature** computed in **layer** t .
- ▶ In particular, $\mathbf{F}_{v \bullet}^{(t)}$ in $\mathbb{R}^{1 \times d}$ denotes **embedding of vertex** v .
- ▶ “learnable” **weight** matrices $\mathbf{W}_1^{(t)}, \mathbf{W}_2^{(t)} \in \mathbb{R}^{d \times d}$ and **bias** $\mathbf{b}^{(t)} \in \mathbb{R}^{1 \times d}$.

$$\mathbf{F}_{v \bullet}^{(0)} := L_G(v) \quad \mathbf{F}_{v \bullet}^{(t)} := \sigma \left(\mathbf{F}_{v \bullet}^{(t-1)} \mathbf{W}_1^{(t)} + \sum_{u \in N_G(v)} \mathbf{F}_{u \bullet}^{(t-1)} \mathbf{W}_2^{(t)} + \mathbf{b}^{(t)} \right)$$

- ▶ By varying weights and biases, an **infinite family** of vertex embeddings is obtained.

- ▶ Let σ a **non-linear activation function** $\mathbb{R} \rightarrow \mathbb{R}$ (ReLU, sigmoid, sign,...).
- ▶ Vertex set of graph identified with $[n] := \{1, 2, \dots, n\}$ for some $n \in \mathbb{N}$.
- ▶ Output space $\mathbb{Y} = \mathbb{R}^d$ for some $d \in \mathbb{N}$.
- ▶ Matrix $\mathbf{F}^{(t)}$ in $\mathbb{R}^{n \times d}$ represents **vertex feature** computed in **layer** t .
- ▶ In particular, $\mathbf{F}_{v\bullet}^{(t)}$ in $\mathbb{R}^{1 \times d}$ denotes **embedding of vertex** v .
- ▶ “learnable” **weight** matrices $\mathbf{W}_1^{(t)}, \mathbf{W}_2^{(t)} \in \mathbb{R}^{d \times d}$ and **bias** $\mathbf{b}^{(t)} \in \mathbb{R}^{1 \times d}$.

$$\mathbf{F}_{v\bullet}^{(0)} := L_G(v) \quad \mathbf{F}_{v\bullet}^{(t)} := \sigma \left(\mathbf{F}_{v\bullet}^{(t-1)} \mathbf{W}_1^{(t)} + \sum_{u \in N_G(v)} \mathbf{F}_{u\bullet}^{(t-1)} \mathbf{W}_2^{(t)} + \mathbf{b}^{(t)} \right)$$

- ▶ By varying weights and biases, an **infinite family** of vertex embeddings is obtained.

- ▶ We can also define a **graph embedding**
- ▶ “learnable” **weight** matrix $\mathbf{W} \in \mathbb{R}^{d \times d}$ and **bias** $\mathbf{b} \in \mathbb{R}^{1 \times d}$.
- ▶ L is **number of layers**.

$$\mathbf{F} := \sigma \left(\sum_{v \in V_G} \mathbf{F}_{v \bullet}^{(L)} \mathbf{W} + \mathbf{b} \right)$$

- ▶ Easy to see that these GNNs define **invariant** embeddings.

- ▶ But what does “learning an unknown embedding” mean?
- ▶ We briefly discuss this in the [semi-supervised](#) setting.

Ingredient #1: Training set

We want to learn $\Psi : \mathcal{G} \rightarrow (\mathcal{V}^p \rightarrow \mathbb{Y})$ but we may only **partially** know this embedding ...

- ▶ Partial knowledge of Ψ is revealed through a **training set**

$$\mathcal{T} = \left\{ (G_1, \mathbf{v}_1, \Psi(G_1, \mathbf{v}_1)), \dots, (G_\ell, \mathbf{v}_\ell, \Psi(G_\ell, \mathbf{v}_\ell)) \right\} \subseteq \mathcal{G} \times \mathcal{V}^p \times \mathbb{Y},$$

with **graphs** $G_i \in \mathcal{G}$ and **p -vertex tuples** \mathbf{v}_i in G_i .

Ingredient #1: Training set

We want to learn $\Psi : \mathcal{G} \rightarrow (\mathcal{V}^p \rightarrow \mathbb{Y})$ but we may only **partially** know this embedding ...

- ▶ Partial knowledge of Ψ is revealed through a **training set**

$$\mathcal{T} = \left\{ (G_1, \mathbf{v}_1, \Psi(G_1, \mathbf{v}_1)), \dots, (G_\ell, \mathbf{v}_\ell, \Psi(G_\ell, \mathbf{v}_\ell)) \right\} \subseteq \mathcal{G} \times \mathcal{V}^p \times \mathbb{Y},$$

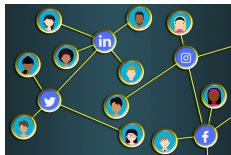
with **graphs** $G_i \in \mathcal{G}$ and **p -vertex tuples** \mathbf{v}_i in G_i .



(molecule, yes/no)



(cora, paper, topic)



(social, p_x, p_y , yes/no)

Ingredient #2: Hypothesis class

The partially known embedding Ψ will be learned by selecting a good candidate from a class of embeddings.

- ▶ An hypothesis class is a collection \mathcal{F} of invariant p -vertex embeddings:

$$\mathcal{F} \subseteq \{\text{all invariant } p\text{-vertex embeddings}\}.$$

Ingredient #2: Hypothesis class

The partially known embedding Ψ will be learned by selecting a good candidate from a class of embeddings.

- ▶ An hypothesis class is a collection \mathcal{F} of invariant p -vertex embeddings:

$$\mathcal{F} \subseteq \{\text{all invariant } p\text{-vertex embeddings}\}.$$

For example, \mathcal{F} can be the collection of

- ▶ GNN 101's
- ▶ Graph kernel methods
- ▶ Message-Passing Neural Networks
- ▶ Invariant Graph Networks
- ▶ Subgraph Networks
- ▶

Ingredient #3: Loss function

How to compare Ψ (embedding to be learned) with embeddings ξ from \mathcal{F} ?

- ▶ This is done using a **loss function** $\mathbf{L} : \mathbb{Y}^2 \rightarrow \mathbb{R}$.
- ▶ Given graph G , p -vertex tuple \mathbf{v} in our training set \mathcal{T} and embedding ξ in our hypothesis class \mathcal{F} ,

$$\mathbf{L}\left(\underbrace{\xi(G, \mathbf{v})}_{\in \mathbb{Y}}, \underbrace{\Psi(G, \mathbf{v})}_{\in \mathbb{Y}}\right) \in \mathbb{R}$$

measures **quality** of ξ on the training example $(G, \mathbf{v}, \Psi(G, \mathbf{v}))$.

Ingredient #3: Loss function

How to **compare** Ψ (embedding to be learned) with embeddings ξ from \mathcal{F} ?

- ▶ This is done using a **loss function** $\mathbf{L} : \mathbb{Y}^2 \rightarrow \mathbb{R}$.
- ▶ Given graph G , p -vertex tuple \mathbf{v} in our training set \mathcal{T} and embedding ξ in our hypothesis class \mathcal{F} ,

$$\mathbf{L}\left(\underbrace{\xi(G, \mathbf{v})}_{\in \mathbb{Y}}, \underbrace{\Psi(G, \mathbf{v})}_{\in \mathbb{Y}}\right) \in \mathbb{R}$$

measures **quality** of ξ on the training example $(G, \mathbf{v}, \Psi(G, \mathbf{v}))$.

Example loss functions: cross entropy, least squares, ...

Graph learning: Empirical risk minimization

- ▶ Given training data

$$\mathcal{T} = \left\{ (G_1, \mathbf{v}_1, \Psi(G_1, \mathbf{v}_1)), \dots, (G_\ell, \mathbf{v}_\ell, \Psi(G_\ell, \mathbf{v}_\ell)) \right\} \subseteq \mathcal{G} \times \mathcal{V}^p \times \mathbb{Y},$$

- ▶ hypothesis class \mathcal{F} , and
- ▶ loss function \mathbf{L} , return

$$\hat{\xi} := \arg \min_{\xi \in \mathcal{F}} \frac{1}{|\mathcal{T}|} \sum_{(G_i, \mathbf{v}_i, \Psi(G_i, \mathbf{v}_i)) \in \mathcal{T}} \mathbf{L}(\xi(G_i, \mathbf{v}_i), \Psi(G_i, \mathbf{v}_i)).$$

- ▶ In other words, find a graph, vertex or p -vertex embedding in \mathcal{F} which **minimizes the risk** (measured by the loss function) on the training data.

Graph learning: Empirical risk minimization

$$\hat{\xi} := \arg \min_{\xi \in \mathcal{F}} \frac{1}{|\mathcal{T}|} \sum_{(G_i, \mathbf{v}_i, \Psi(G, \mathbf{v}_i)) \in \mathcal{T}} \mathbf{L}(\xi(G_i, \mathbf{v}_i), \Psi(G_i, \mathbf{v}_i)).$$

- ▶ Graph learning systems: **optimization techniques** for finding best hypothesis.
- ▶ Typically based on **back propagation** and **gradient descent** like methods.

Graph learning: Empirical risk minimization

$$\hat{\xi} := \arg \min_{\xi \in \mathcal{F}} \frac{1}{|\mathcal{T}|} \sum_{(G_i, \mathbf{v}_i, \Psi(G, \mathbf{v}_i)) \in \mathcal{T}} \mathbf{L}(\xi(G_i, \mathbf{v}_i), \Psi(G_i, \mathbf{v}_i)).$$

- ▶ Graph learning systems: **optimization techniques** for finding best hypothesis.
- ▶ Typically based on **back propagation** and **gradient descent** like methods.

- ▶ We will be focussing on:

Expressivity of classes \mathcal{F} of embeddings.

2. Expressive Power

What graph information can be extracted by embedding methods?

Recall our GNN 101's.

- ▶ Which graph or vertex embeddings can they **express**?
- ▶ Which graph or vertex embeddings can they **approximate**?
- ▶ Which graphs or vertices can be **discriminated/distinguished**?

Answers to these questions may reveal

- ▶ what **graph information** is used by embedding methods;
- ▶ which embeddings could – in principle – be **learned**; and
- ▶ whether **more powerful** embedding methods may be needed for the application at hand.

Let $\Psi : \mathcal{G} \rightarrow (\mathcal{V}^p \rightarrow \mathbb{Y})$ be a p -vertex embedding and let \mathcal{F} be a class of embeddings and let \mathcal{C} be a subset of \mathcal{G} .

- ▶ \mathcal{F} can \mathcal{C} -express Ψ if

$$\exists \xi \in \mathcal{F}, \forall G \in \mathcal{C}, \forall \mathbf{v} \in V_G^p : \Psi(G, \mathbf{v}) = \xi(G, \mathbf{v}).$$

- ▶ \mathcal{F} can \mathcal{C} -approximate Ψ if

$$\forall \epsilon > 0, \exists \xi_\epsilon \in \mathcal{F}, \forall G \in \mathcal{C}, \forall \mathbf{v} \in V_G^p : \|\Psi(G, \mathbf{v}) - \xi_\epsilon(G, \mathbf{v})\| < \epsilon.$$

for some norm $\|\cdot\|$ on \mathbb{Y} .

If $\mathcal{C} = \mathcal{G}$ we just say express or approximate.

Separation power measures how well \mathcal{F} can separate different inputs.

- ▶ As before, \mathcal{F} be a class of p -vertex embeddings.
- ▶ The separation power of \mathcal{F} is captured by equivalence relation $\rho(\mathcal{F})$ on $\mathcal{G} \times \mathcal{V}^p$:

$$(G, \mathbf{v}; H, \mathbf{w}) \in \rho(\mathcal{F}) \iff \forall \xi \in \mathcal{F} : \xi(G, \mathbf{v}) = \xi(H, \mathbf{w}).$$

- ▶ In other words, (G, \mathbf{v}) and (H, \mathbf{w}) are in $\rho(\mathcal{F})$ when these cannot be separated by any embedding in \mathcal{F}
- ▶ Similar to the notion of indistinguishability for logics and query languages.

Strongest power

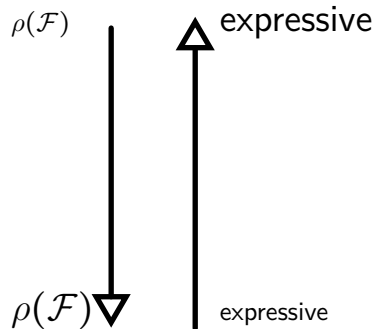
- ▶ \mathcal{F} powerful enough to distinguish non-isomorphic graphs:

$$\rho(\mathcal{F}) = \{\text{all pairs of isomorphic graphs}\}.$$

Weakest power

- ▶ \mathcal{F} consisting of constant functions.

$$\rho(\mathcal{F}) = \{\text{all pairs of graphs}\}.$$



Separation power allows for **comparing totally different embedding methods** by means of subset relationship of their separation power.

Expressive power in ML community

- ▶ Primary focus has been on **separation power**.
- ▶ Aim is to provide a **characterization** of when $(G, H) \in \rho(\mathcal{F})$ holds.
- ▶ For example,

Theorem

$$\rho(\text{GNNs 101}) = \rho(\text{color refinement}).$$

Shown in the - by now - seminal paper in the area of graph learning¹.

¹ Morris, Ritzert, Fey, Hamilton, Lenssen, Rattan, Grohe. Weisfeiler and Leman Go Neural: Higher-Order Graph Neural Networks. AAAI (2019)

Expressive power in ML community

- ▶ Characterizing $\rho(\mathcal{F})$ is a bit of **hot potato game**.
- ▶ Try to find characterizations of $\rho(\mathcal{F})$ that are **insightful**.
- ▶ For example, **color refinement** has been well-studied, many properties thereof are known.
- ▶ In particular, G and H are in $\rho(\text{color refinement})$ if and only if

$$\text{hom}(T, G) = \text{hom}(T, H)$$

for all **trees** T . Here, $\text{hom}(T, G)$ counts **homomorphisms** from T to G .²

- ▶ GNNs 101 can only leverage **tree-based information** present in the graphs.

²  Dell, Grohe and Rattan: Lovász Meets Weisfeiler and Leman. ICALP 2018

Expressive power in ML community

- ▶ Characterizing $\rho(\mathcal{F})$ is a bit of **hot potato game**.
- ▶ Try to find characterizations of $\rho(\mathcal{F})$ that are **insightful**.
- ▶ For example, **color refinement** has been well-studied, many properties thereof are known.
- ▶ In particular, G and H are in $\rho(\text{color refinement})$ if and only if

$$\text{hom}(T, G) = \text{hom}(T, H)$$

for all **trees** T . Here, $\text{hom}(T, G)$ **counts homomorphisms** from T to G .²

- ▶ GNNs 101 can only leverage **tree-based information** present in the graphs.

²  Dell, Grohe and Rattan: Lovász Meets Weisfeiler and Leman. ICALP 2018

Two additional reasons for studying separation power:

1. Close connection between **separation power** $\rho(\mathcal{F})$ and ability for \mathcal{F} to **approximate** functions.
2. Close connection between $\rho(\mathcal{F})$ and **Vapnik-Chervonenkis (VC) dimension** of \mathcal{F} . This implies properties of generalization aspects of \mathcal{F} .³

³ WL meet VC. Morris, G*, Tönshoff and Grohe. ICML 2023

- ▶ \mathcal{C} be a **compact** subset of \mathcal{G} .
- ▶ Assume embeddings in \mathcal{F} are **continuous**.
- ▶ We can use tools from **analysis** and **topology**.


Theorem (General version of Stone-Weierstrass)

If \mathcal{F} is closed under linear combinations and product, then \mathcal{F} can **\mathcal{C} -approximate** any continuous embedding $\Psi : \mathcal{G} \rightarrow \mathbb{R}$ such that

$$\rho(\mathcal{F}) \subseteq \rho(\{\Psi\})$$

holds.

- ▶ Can be generalized to vertex embeddings and output spaces $\mathbb{Y} = \mathbb{R}^d$.^{4,5}

⁴  Azizian and Lelarge. Characterizing the Expressive Power of Invariant and Equivariant GNNs, ICLR 2021

⁵  G* and Reutter. Expressiveness and Approximation Properties of GNNs. ICLR 2022

Theorem

On compact sets of graphs, GNNs 101 can approximate any continuous embedding Ψ whose separation power is bounded by color refinement.

Follows from $\rho(\text{GNNs 101}) = \rho(\text{color refinement})$, universality theorem of neural networks (to approximate product), and Stone-Weierstrass. (Alternative proof based on homomorphism counts.⁶)

⁶  Nguyen and Maehara. Graph Homomorphism Convolution. ICML 2020.

Consequence of Stone-Weierstrass:⁷

Theorem

*For a class \mathcal{F} to be able to approximate **any** invariant embedding on a compact set of graphs, \mathcal{F} needs to be able to separate any two non-isomorphic graphs.*

⁷  Chen, Villar, Chen and Bruna. 2019. On the Equivalence Between Graph Isomorphism Testing and Function Approximation With GNNs. Neurips 2019.

Promised query language perspective is coming up in a few moments.

Expressive power in ML community

- ▶ Every week new embedding methods are being proposed.
- ▶ Continuous stream of papers on arxiv.
- ▶ Has become standard to analyze separation power of new methods.
- ▶ This is done often in an ad hoc way.

This is where the [language approach](#) comes in the picture.

A small selection of methods...

GraphSage GINs GCNs SGNs
GATs GatedGCNs extended GINs 2-IGNs ChebNet ...
Walk GNNs 2WL-GNNs ring-GNNs 1-Dropout GNNs
Id-aware GNNs CayleyNet 3-IGNs 2-FGNNs ...
kWL-GNNs k-FGNNs (k+1)-IGNs GSNs k-Dropout GNNs ...

1. View embedding methods as **queries in some graph embedding language**
2. **Transfer** our understanding of **separation power** of these languages **back** to embedding methods.

Recipe

- ▶ A new embedding method just needs to be cast in the embedding language to know a bound on its expressive power.

3. Embedding Language #1: Message Passing Neural Networks

Message Passing Neural Networks (MPNNs)

- ▶ We go **back in time** (around 2016) when embedding methods like
 - ▶ Graph convolutional networks (Duvenaud et al. 2016, Kearnes et al. 2016),
 - ▶ Gated GNNs (Li et al. 2016),
 - ▶ Interaction Networks (Battaglia et al. 2016),
 - ▶ Deep tensor neural networks (Schütt et al. 2017), and
 - ▶ Laplacian based graph convolutional networks (Bruna et al. 2013, Defferrard et al. 2016, Kipf & Welling 2016)

were “**hot**”.

Message Passing Neural Networks (MPNNs)

- In 2017, Gilmer et al.⁸ looked at the specifications in those papers

- 1: **Input:** molecule, radius R , $H_1^1 \dots H_R^5$, output weights V
- 2: **Initialize:** fingerprint vector
- 3: **for** each atom a in molecule
- 4: $\mathbf{r}_a \leftarrow g(a)$ ▷ lookup atom feat
- 5: **for** $L = 1$ to R ▷ for each la
- 6: **for** each atom a in molecule
- 7: $\mathbf{r}_1 \dots \mathbf{r}_N = \text{neighbors}(a)$
- 8: $\mathbf{v} \leftarrow \mathbf{r}_a + \sum_{i=1}^N \mathbf{r}_i$ ▷ s
- 9: $\mathbf{r}_a \leftarrow \sigma(\mathbf{v} H_L^N)$ ▷ smooth function
- 10: $\mathbf{i} \leftarrow \text{softmax}(\mathbf{r}_a W_L)$ ▷ sparsify
- 1: $\mathbf{f} \leftarrow \mathbf{f} + \mathbf{i}$ ▷
- 2: **Return:** real-valued vector

$$\mathbf{h}_v^{(1)} = [\mathbf{x}_v^\top, \mathbf{0}]^\top \quad (1)$$

$$\mathbf{a}_v^{(t)} = \mathbf{A}_v^\top [\mathbf{h}_1^{(t-1)\top} \dots \mathbf{h}_{|\mathcal{V}|}^{(t-1)\top}]^\top + \mathbf{b} \quad (2)$$

$$\mathbf{z}_v^t = \sigma(\mathbf{W}^z \mathbf{a}_v^{(t)} + \mathbf{U}^z \mathbf{h}_v^{(t-1)}) \quad (3)$$

$$\mathbf{r}_v^t = \sigma(\mathbf{W}^r \mathbf{a}_v^{(t)} + \mathbf{U}^r \mathbf{h}_v^{(t-1)}) \quad (4)$$

$$\widetilde{\mathbf{h}}_v^{(t)} = \tanh(\mathbf{W} \mathbf{a}_v^{(t)} + \mathbf{U}(\mathbf{r}_v^t \odot \mathbf{h}_v^{(t-1)})) \quad (5)$$

$$\mathbf{h}_v^{(t)} = (1 - \mathbf{z}_v^t) \odot \mathbf{h}_v^{(t-1)} + \mathbf{z}_v^t \odot \widetilde{\mathbf{h}}_v^{(t)}. \quad (6)$$

$$\text{IN}(G) = \phi_O(a(G, X, \phi_R(m(G)))) \quad (1)$$

$$\begin{aligned} m(G) &= B = \{b_k\}_{k=1 \dots N_R} & a(G, X, E) &= C = \{c_j\}_{j=1 \dots N_O} \\ f_R(b_k) &= e_k & f_O(c_j) &= p_j \\ \phi_R(B) &= E = \{e_k\}_{k=1 \dots N_R} & \phi_O(C) &= P = \{p_j\}_{j=1 \dots N_O} \end{aligned} \quad (2)$$

The marshalling function m rearranges the objects and relations into interaction terms $b_k =$

$$x_{k+1,j} = h\left(V \sum_{i=1}^{f_k-1} F_{k,i,j} V^T x_{k,i}\right) \quad (j = 1 \dots f_k), \quad (3.2)$$

network (GCN) with the following layer-wise propagation rule:

$$H^{(l+1)} = \sigma\left(\tilde{D}^{-\frac{1}{2}} \tilde{A} \tilde{D}^{-\frac{1}{2}} H^{(l)} W^{(l)}\right). \quad (2)$$

Here, $\tilde{A} = A + I_N$ is the adjacency matrix of the undirected graph \mathcal{G} with added self-connections.

and proposed a first unifying framework for specifying embedding methods: **Message Passing Neural Networks**

⁸ Gilmer, Schoenholz, Riley, Vinyals, Dahl. Neural Message Passing for Quantum Chemistry, Neurips, 1263–1272 (2017)

Message Passing Neural Networks (MPNNs)

- ▶ Liberally interpreted, Gilmer et al. proposed an **inductive** way of defining **vertex** and **graph embeddings**.
- ▶ Indeed, one has **initial vertex embeddings**

$$\xi : \mathcal{G} \rightarrow (\mathcal{V} \rightarrow \mathbb{R}^d) : (G, v) \mapsto \nu_G(v)$$

for some encoding of the **vertex label** in $\nu_G(v) \in \Sigma$ in \mathbb{R}^d .

Message Passing Neural Networks (MPNNs)

- ▶ Liberally interpreted, Gilmer et al. proposed an **inductive** way of defining **vertex** and **graph embeddings**.
- ▶ Indeed, one has **initial vertex embeddings**

$$\xi : \mathcal{G} \rightarrow (\mathcal{V} \rightarrow \mathbb{R}^d) : (G, v) \mapsto \nu_G(v)$$

for some encoding of the **vertex label** in $\nu_G(v) \in \Sigma$ in \mathbb{R}^d .

- ▶ Then, let $\xi' : \mathcal{G} \rightarrow (\mathcal{V} \rightarrow \mathbb{R}^{d'})$ be an old vertex embedding. A **new vertex embedding** $\xi : \mathcal{G} \rightarrow (\mathcal{V} \rightarrow \mathbb{R}^d)$ can be obtained by:

Message Passing Neural Networks (MPNNs)

- ▶ Liberally interpreted, Gilmer et al. proposed an **inductive** way of defining **vertex** and **graph embeddings**.
- ▶ Indeed, one has **initial vertex embeddings**

$$\xi : \mathcal{G} \rightarrow (\mathcal{V} \rightarrow \mathbb{R}^d) : (G, v) \mapsto \nu_G(v)$$

for some encoding of the **vertex label** in $\nu_G(v) \in \Sigma$ in \mathbb{R}^d .

- ▶ Then, let $\xi' : \mathcal{G} \rightarrow (\mathcal{V} \rightarrow \mathbb{R}^{d'})$ be an old vertex embedding. A **new vertex embedding** $\xi : \mathcal{G} \rightarrow (\mathcal{V} \rightarrow \mathbb{R}^d)$ can be obtained by:

$$(G, v) \mapsto \xi(G, v) := \sum_{u \in N_G(v)} \xi'(G, u)$$

Message Passing Neural Networks (MPNNs)

- ▶ Liberally interpreted, Gilmer et al. proposed an **inductive** way of defining **vertex** and **graph embeddings**.
- ▶ Indeed, one has **initial vertex embeddings**

$$\xi : \mathcal{G} \rightarrow (\mathcal{V} \rightarrow \mathbb{R}^d) : (G, v) \mapsto \nu_G(v)$$

for some encoding of the **vertex label** in $\nu_G(v) \in \Sigma$ in \mathbb{R}^d .

- ▶ Then, let $\xi' : \mathcal{G} \rightarrow (\mathcal{V} \rightarrow \mathbb{R}^{d'})$ be an old vertex embedding. A **new vertex embedding** $\xi : \mathcal{G} \rightarrow (\mathcal{V} \rightarrow \mathbb{R}^d)$ can be obtained by:

$$(G, v) \mapsto \xi(G, v) := \underbrace{\text{Update}}_{\substack{\text{any function} \\ \mathbb{R}^{2d'} \rightarrow \mathbb{R}^d}} \left(\xi'(G, v), \sum_{u \in N_G(v)} \xi'(G, u) \right)$$

Message Passing Neural Networks (MPNNs)

- ▶ Finally, one can also construct **graph embeddings** $\xi : \mathcal{G} \rightarrow \mathbb{R}^d$ from a **vertex embedding** $\xi' : \mathcal{G} \rightarrow (\mathcal{V} \rightarrow \mathbb{R}^{d'})$ as follows:

$$G \mapsto \xi(G) := \sum_{v \in V_G} \xi'(G, v)$$

Message Passing Neural Networks (MPNNs)

- ▶ Finally, one can also construct **graph embeddings** $\xi : \mathcal{G} \rightarrow \mathbb{R}^d$ from a **vertex embedding** $\xi' : \mathcal{G} \rightarrow (\mathcal{V} \rightarrow \mathbb{R}^{d'})$ as follows:

$$G \mapsto \xi(G) := \underbrace{\mathbf{Readout}}_{\substack{\text{any function} \\ \mathbb{R}^{d'} \rightarrow \mathbb{R}^d}} \left(\sum_{v \in V_G} \xi'(G, v) \right)$$

Message Passing Neural Networks (MPNNs)

- ▶ Finally, one can also construct **graph embeddings** $\xi : \mathcal{G} \rightarrow \mathbb{R}^d$ from a **vertex embedding** $\xi' : \mathcal{G} \rightarrow (\mathcal{V} \rightarrow \mathbb{R}^{d'})$ as follows:

$$G \mapsto \xi(G) := \underbrace{\text{Readout}}_{\substack{\text{any function} \\ \mathbb{R}^{d'} \rightarrow \mathbb{R}^d}} \left(\sum_{v \in V_G} \xi'(G, v) \right)$$

Easy exercise: The GNNs 101 we have seen before are MPNNs.

Layer-based vector embedding computation

$$\mathbf{F}^{(0)} := \text{Initial vertex labels} \quad \mathbf{F}^{(t)} := \sigma \left(\mathbf{F}^{(t-1)} \mathbf{W}_1^{(t)} + \mathbf{A} \mathbf{F}^{(t-1)} \mathbf{W}_2^{(t)} + \mathbf{B}^{(t)} \right)$$

- σ is a non-linear activation function $\mathbb{R} \rightarrow \mathbb{R}$;
- \mathbf{A} is adjacency matrix in $\mathbb{R}^{n \times n}$ of a graph;
- $\mathbf{F}^{(t)}$ vertex embedding in $\mathbb{R}^{n \times d}$ computed in layer t ; and
- learnable weight matrices $\mathbf{W}_1^{(t)}, \mathbf{W}_2^{(t)} \in \mathbb{R}^{d \times d}$ and bias $\mathbf{B}^{(t)} \in \mathbb{R}^{n \times d}$.

$$\mathbf{F}_{i \rightarrow v}^{(t)} := \sigma \left(\mathbf{F}_{i \rightarrow v}^{(t-1)} \mathbf{W}_1^{(t)} + \sum_{u \in N_G(v)} \mathbf{F}_{i \rightarrow u}^{(t-1)} \mathbf{W}_2^{(t)} + \mathbf{b}^{(t)} \right) \quad \mathbf{F} := \sigma \left(\sum_{v \in V_G} \mathbf{F}^{(t)} \mathbf{W} + \mathbf{b} \right)$$

Let's put on our database glasses

Let's put on our database glasses (you can find those under your seats).

We can turn MPNNs into a **specification language** for vertex and graph embeddings.

- ▶ We fix the **output space** \mathbb{Y} to be $\cup_d \mathbb{R}^d$.
- ▶ We take **two variables** x_1 and x_2 .
- ▶ Inductively define **MPNN expressions** φ .

We can turn MPNNs into a **specification language** for vertex and graph embeddings.

- ▶ We fix the **output space** \mathbb{Y} to be $\cup_d \mathbb{R}^d$.
- ▶ We take **two variables** x_1 and x_2 .
- ▶ Inductively define **MPNN expressions** φ . Each expression φ comes with:
 - ▶ a **dimension** $\dim(\varphi) \in \mathbb{N}$ and
 - ▶ a set of **free variables** $\text{fv}(\varphi)$.

We can turn MPNNs into a **specification language** for vertex and graph embeddings.

- ▶ We fix the **output space** \mathbb{Y} to be $\cup_d \mathbb{R}^d$.
- ▶ We take **two variables** x_1 and x_2 .
- ▶ Inductively define **MPNN expressions** φ . Each expression φ comes with:
 - ▶ a **dimension** $\dim(\varphi) \in \mathbb{N}$ and
 - ▶ a set of **free variables** $\text{fv}(\varphi)$.
- ▶ **Semantics**: p -vertex embeddings $\xi_\varphi : \mathcal{G} \rightarrow (\mathcal{V}^p \rightarrow \mathbb{Y})$ (p determined by number free variables, output space \mathbb{R}^d by dimension of φ).

- ▶ Initially, **atomic** MPNN expressions are of the form

$$\varphi(x_i) := \text{lab}_j(x_i)$$

with free variable x_i , $j = 1, 2, 3, \dots$ and of dimension 1.

- ▶ Initially, **atomic** MPNN expressions are of the form

$$\varphi(x_i) := \text{lab}_j(x_i)$$

with free variable x_i , $j = 1, 2, 3, \dots$ and of dimension 1.

For example, let

$$\varphi(x_1) := \text{Lab}_2(x_1)$$

and G a graph with vertex labelling $\nu_G : V_G \rightarrow \mathbb{R}^3$. The corresponding semantics is the **vertex embedding**

$$\xi_\varphi : (G, \nu) \mapsto (\nu_G(v))_2 \in \mathbb{R}$$

i.e., the second component of $\nu_G(v) \in \mathbb{R}^3$.

We close under **application of functions** coming from some set Ω .

- ▶ Consider MPNN expressions $\varphi_1(x_i), \dots, \varphi_\ell(x_i)$ with free variable x_i and of dimension d_1, \dots, d_ℓ , respectively.
- ▶ Take a **function** $\mathbf{F} : \mathbb{R}^{d_1 + \dots + d_\ell} \rightarrow \mathbb{R}^d$ in Ω . Then,

$$\varphi(x_i) := \mathbf{F}(\varphi_1(x_i), \dots, \varphi_\ell(x_i))$$

is again an MPNN expression, of dimension d and free variable x_i .

We close under **application of functions** coming from some set Ω .

- ▶ Consider MPNN expressions $\varphi_1(x_i), \dots, \varphi_\ell(x_i)$ with free variable x_i and of dimension d_1, \dots, d_ℓ , respectively.
- ▶ Take a **function** $\mathbf{F} : \mathbb{R}^{d_1 + \dots + d_\ell} \rightarrow \mathbb{R}^d$ in Ω . Then,

$$\varphi(x_i) := \mathbf{F}(\varphi_1(x_i), \dots, \varphi_\ell(x_i))$$

is again an MPNN expression, of dimension d and free variable x_i .

Example:

$$\varphi(x_1) := \text{ReLU}(\varphi'(x_1)) \text{ with } \text{ReLU}(x) := \max\{0, x\} \text{ in } \Omega$$

with **semantics**

$$\xi_\varphi : (G, v) \mapsto \text{ReLU}(\xi_{\varphi'}(G, v))$$

We close under **restricted** application of aggregation functions from a set Θ .

- ▶ Let $\varphi_1(x_1)$ and $\varphi_2(x_2)$ expressions with free variables x_1 and x_2 of dimension d_1 and d_2 , respectively.
- ▶ Let θ be any **aggregate function** from **bags of elements** in $\mathbb{R}^{d_1+d_2}$ to \mathbb{R}^d .
- ▶ Then,

$$\varphi(x_1) := \text{agg}_{x_2}^{\theta} \left(\varphi_1(x_1), \varphi_2(x_2) \mid \underbrace{E(x_1, x_2)}_{\text{edge relation}} \right)$$

is an MPNN expression of dimension d and free variable x_1 (similarly with roles of x_1 and x_2 reversed.)

We close under **restricted** application of aggregation functions from a set Θ .

- ▶ Let $\varphi_1(x_1)$ and $\varphi_2(x_2)$ expressions with free variables x_1 and x_2 of dimension d_1 and d_2 , respectively.
- ▶ Let θ be any **aggregate function** from **bags of elements** in $\mathbb{R}^{d_1+d_2}$ to \mathbb{R}^d .
- ▶ Then,

$$\varphi(x_1) := \text{agg}_{x_2}^{\theta} \left(\varphi_1(x_1), \varphi_2(x_2) \mid \underbrace{E(x_1, x_2)}_{\text{edge relation}} \right)$$

is an MPNN expression of dimension d and free variable x_1 (similarly with roles of x_1 and x_2 reversed.)

Example:

$$\varphi(x_1) := \text{agg}_{x_2}^{\text{sum}} \left(\varphi_1(x_1), \varphi_2(x_2) \mid E(x_1, x_2) \right)$$

and corresponding embedding

$$\xi_{\varphi} : (G, v) \mapsto \sum_{(v, u) \in E_G} (\xi_{\varphi_1}(G, v), \xi_{\varphi_2}(G, u))$$

We can also express graph embeddings.

- ▶ Let $\varphi'(x_1)$ be an MPNN expression with free variable x_1 and of dimension d' .
- ▶ Let θ be an **aggregation function** from bags of elements in $\mathbb{R}^{d'}$ to \mathbb{R}^d .
- ▶ Then,

$$\varphi := \text{agg}_{x_1}^{\theta}(\varphi'(x_1))$$

is an MPNN expression of dimension d and no free variables.

We can also express graph embeddings.

- ▶ Let $\varphi'(x_1)$ be an MPNN expression with free variable x_1 and of dimension d' .
- ▶ Let θ be an **aggregation function** from bags of elements in $\mathbb{R}^{d'}$ to \mathbb{R}^d .
- ▶ Then,

$$\varphi := \text{agg}_{x_1}^{\theta}(\varphi'(x_1))$$

is an MPNN expression of dimension d and no free variables.

Example:

$$\varphi := \text{agg}_{x_1}^{\text{sum}}(\varphi_1(x_1))$$

and corresponding embedding

$$\xi_{\varphi} : G \mapsto \sum_{v \in V_G} (\xi_{\varphi_1}(G, v))$$

- ▶ We have thus defined a language $\text{MPNN}(\Omega, \Theta)$.
- ▶ A very limited fragment of **calculus with aggregates**.⁹
- ▶ It **differs** from classical MPNNs because in those one **restricts** how function application and aggregation interleave:

$$\varphi^{(t)}(x_1) := \mathbf{F}^{(t)}\left(\varphi^{(t-1)}(x_1), \text{agg}_{x_2}^{\theta^{(t)}}\left(\varphi^{(t-1)}(x_2) \mid E(x_1, x_2)\right)\right)$$

⁹  Hella, Libkin, Nurmonen and Wong, Logics with aggregate operators, JACM 2001.

- ▶ Existing architectures can be easily cast as $\text{MPNN}(\Omega, \Theta)$ expressions, due to more flexible definition when compared to classical MPNNs.

GraphSage GINs GCNs SGNs
GATs GatedGCNs extended GINs 2-IGNs ChebNet ...
Walk GNNs 2WL-GNNs ring-GNNs 1-Dropout GNNs
Id-aware GNNs CayleyNet 3-IGNs 2-FGNNs ...
kWL-GNNs k-FGNNs (k+1)-IGNs GSNs k-Dropout GNNs ...

What about separation power of $MPNN(\Omega, \Theta)$?

Color refinement: Iteratively computes a **coloring** of vertices of a graph:

1. **Initialization:** all vertices have their original colors (labels)
2. **Refinement Step:** two vertices v and w get different colors if there is a color c such that v and w have a different number of neighbors of color c .

This process terminates and a **graph** will get a color based on the **multiset of colors** of all its vertices.

ρ (color refinement) contains pairs of graphs/vertices with the same coloring.

Color refinement: Iteratively computes a **coloring** of vertices of a graph:

1. **Initialization:** all vertices have their original colors (labels)
2. **Refinement Step:** two vertices v and w get different colors if there is a color c such that v and w have a different number of neighbors of color c .

This process terminates and a **graph** will get a color based on the **multiset of colors** of all its vertices.

ρ (color refinement) contains pairs of graphs/vertices with the same coloring.

Theorem

For any Ω and Θ , $\rho_{0/1}(\text{color refinement}) \subseteq \rho_{0/1}(\text{MPNN}(\Omega, \Theta))$.

- ▶ For MPNNs, this was shown in the seminal papers^{10,11} and then expanded.¹²
- ▶ This can also be shown using the correspondence

$$\rho(\text{color refinement}) = \rho(\text{guarded } C_2)$$

and elimination of function and aggregation functions by detour to infinitary counterparts of guarded C_2 .¹³

¹⁰ [X] Xu, Hu, Leskovec, Jegelka. How Powerful are Graph Neural Networks? ICLR (2019)

¹¹ [X] Morris, Ritzert, Fey, Hamilton, Lenssen, Rattan, Grohe. Weisfeiler and Leman Go Neural: Higher-Order Graph Neural Networks. AAAI (2019)

¹² [X] G* and Reutter. Expressiveness and Approximation Properties of GNNs. ICLR 2022

¹³ [X] Hella, Libkin, Nurmonen and Wong, Logics with aggregate operators, JACM 2001.


Which functions are needed to match color refinement in separation power?

Theorem

if Ω contains concatenation, linear combinations and non-linear activation functions and Θ consists of summation, then

$$\rho_{0/1}(\text{color refinement}) = \rho_{0/1}(\text{MPNN}(\Omega, \Theta))$$

- ▶ Shown by explicit construction of GNNs 101.¹⁴

¹⁴  Morris, Ritzert, Fey, Hamilton, Lenssen, Rattan, Grohe. Weisfeiler and Leman Go Neural: Higher-Order Graph Neural Networks. AAAI (2019)

Theorem

For sufficiently rich sets Ω of functions and compact \mathcal{C} sets of \mathcal{G} , $\text{MPNN}(\Omega, \text{sum})$ can \mathcal{C} -approximate any embedding $\Psi : \mathcal{G} \rightarrow \mathbb{R}$ satisfying $\rho(\text{color refinement}) \subseteq \rho(\{\Psi\})$.

- ▶ For example: Ω is rich enough when it is **mlp-closed**. That is, for any $q \in \mathbb{N}$, any **multilayered perceptron**¹⁵ $\text{mlp} : \mathbb{R}^q \rightarrow \mathbb{R}$ and any functions f_1, \dots, f_q already in Ω ,

$$\text{mlp}(f_1, \dots, f_q)$$

is also in \mathcal{F} .

¹⁵MLPs: Layered architectures $\mathbf{F}^{(t)} := \sigma(\mathbf{W}^{(t)}\mathbf{F}^{(t-1)} + \mathbf{b}^{(t)})$, $\mathbf{F}^{(0)} := \mathbf{x}$ with learnable weight matrices $\mathbf{W}^{(t)}$, bias vectors $\mathbf{b}^{(t)}$, and activation functions σ .

Can we get rid of compact domain?

Assuming graphs with discrete labels.¹⁶

Theorem

MPNN(Ω, Θ) can express any unary query expressible in graded modal logic.

GNNs 101 already suffice for this.

Theorem

If a first-order logic unary query is expressible in MPNN(Ω, Θ), then it is a query expressible in graded modal logic.

¹⁶  Barceló, Kostylev, Monet, Pérez, Reutter and Silva. The Logical Expressiveness of Graph Neural Networks. ICLR 2020


- ▶ Can any MPNN(Ω, Θ) be converted into a “normal form” MPNN(Ω, Θ)

$$\varphi^{(t)}(x_1) := \mathbf{F}^{(t)}\left(\varphi^{(t-1)}(x_1), \text{agg}_{x_2}^{\theta^{(t)}}\left(\varphi^{(t-1)}(x_2) \mid E(x_1, x_2)\right)\right)$$

- ▶ Important for implementation purposes!
- ▶ Partial results when Ω consists of linear combinations and activation functions σ , and Θ is summation.¹⁷

Theorem

- ▶ Every MPNN(Ω, sum) is equivalent to a normal form MPNN if $\sigma = \text{ReLU}$.
- ▶ On compact domain \mathcal{C} , normal form ReLU MPNNs \mathcal{C} -approximate embeddings in MPNN(Ω, sum).

¹⁷  G*, Steegmans and Van den Bussche: On the Expressive Power of Message-Passing Neural Networks as Global Feature Map Transformers. FoIKS 2022

- ▶ Many other embedding methods exist which are not in $\text{MPNN}(\Omega, \Theta)$.

Only “simple” GNNs are MPNNs :-)

	GraphSage	GINs	GCNs	SGNs	
GATs	GatedGCNs	extended GINs	2-IGNs	ChebNet	...
Walk GNNs	2WL-GNNs	ring-GNNs	1-Dropout GNNs		
	Id-aware GNNs	CayleyNet	3-IGNs	2-FGNNs	...
kWL-GNNs	k-FGNNs	(k+1)-IGNs	GSNs	k-Dropout GNNs	...

How to analyze all other embedding methods?

4. Embedding Language #2:

GEL(Ω, Θ)

and its finite variable fragments

We expand the language of $\text{MPNN}(\Omega, \Theta)$: $\text{GEL}(\Omega, \Theta)$

1. **More** variables x_1, x_2, \dots ;
2. **More** atomic MPNNs;
3. **More** general function and aggregation applications.

We expand the language of $\text{MPNN}(\Omega, \Theta)$: $\text{GEL}(\Omega, \Theta)$

1. **More** variables x_1, x_2, \dots ;
2. **More** atomic MPNNs;
3. **More** general function and aggregation applications.

Label: $\text{Lab}_j(x_i)$ of dimension 1, free variable x_i ;

Edge: $E(x_i, x_j)$ of dimension 1, free variables x_i and x_j ;

Equality: $\mathbf{1}[x_i \text{ op } x_j]$ with $\text{op} \in \{=, \neq\}$, of dimension 1 and free variables x_i and x_j .

Label: $\text{Lab}_j(x_i)$ of dimension 1, free variable x_i ;

Edge: $E(x_i, x_j)$ of dimension 1, free variables x_i and x_j ;

Equality: $\mathbf{1}[x_i \text{ op } x_j]$ with $\text{op} \in \{=, \neq\}$, of dimension 1 and free variables x_i and x_j .

Example:

$$\varphi(x_1, x_2) := E(x_1, x_2) \text{ and } \psi(x_1, x_2) := \mathbf{1}[x_1 \text{ op } x_2]$$

and corresponding 2-vertex embeddings

$$\xi_\varphi : (G, v, w) \mapsto \begin{cases} 1 & \text{if } (v, w) \in E_G \\ 0 & \text{otherwise.} \end{cases} \text{ and } \xi_\psi : (G, v, w) \mapsto \begin{cases} 1 & \text{if } v \text{ op } w \\ 0 & \text{otherwise.} \end{cases}$$

We close GEL under function applications with function in Ω .

- ▶ As before, Ω still set of functions $\mathbb{R}^{d'} \rightarrow \mathbb{R}^d$.
- ▶ Consider expressions $\varphi_1(\mathbf{x}_1), \dots, \varphi_\ell(\mathbf{x}_\ell)$ with free variables \mathbf{x}_i and of dimension d_1, \dots, d_ℓ , respectively.
- ▶ Take a function $\mathbf{F} : \mathbb{R}^{d_1 + \dots + d_\ell} \rightarrow \mathbb{R}^d$ in Ω . Then,

$$\varphi(\mathbf{x}) := \mathbf{F}(\varphi_1(\mathbf{x}_1), \dots, \varphi_\ell(\mathbf{x}_\ell))$$

is again an MPNN expression, of dimension d and free variables $\mathbf{x} := \mathbf{x}_1 \cup \dots \cup \mathbf{x}_\ell$.

We close GEL under function applications with function in Ω .

- ▶ As before, Ω still set of functions $\mathbb{R}^{d'} \rightarrow \mathbb{R}^d$.
- ▶ Consider expressions $\varphi_1(\mathbf{x}_1), \dots, \varphi_\ell(\mathbf{x}_\ell)$ with free variables \mathbf{x}_i and of dimension d_1, \dots, d_ℓ , respectively.
- ▶ Take a function $\mathbf{F} : \mathbb{R}^{d_1 + \dots + d_\ell} \rightarrow \mathbb{R}^d$ in Ω . Then,

$$\varphi(\mathbf{x}) := \mathbf{F}(\varphi_1(\mathbf{x}_1), \dots, \varphi_\ell(\mathbf{x}_\ell))$$

is again an MPNN expression, of dimension d and free variables $\mathbf{x} := \mathbf{x}_1 \cup \dots \cup \mathbf{x}_\ell$.

Example:

$$\varphi(x_1, x_2, x_3) := f_x\left(E(x_1, x_2), f_x\left(E(x_1, x_2), E(x_2, x_3)\right)\right)$$

with $f_x : \mathbb{R}^2 \rightarrow \mathbb{R} : (a, b) \rightarrow a \times b$. Then, we obtain a 3-vertex embedding

$$\xi_\varphi : (G, u, v, w) \mapsto \begin{cases} 1 & (u, v), (u, w), \text{ and } (v, w) \in E_G \\ 0 & \text{otherwise.} \end{cases}$$

We close GEL under aggregation with aggregations in Θ .

- ▶ Let $\varphi_1(\mathbf{x}, \mathbf{y})$ and $\varphi_2(\mathbf{x}, \mathbf{y})$ expressions with free variables (\mathbf{x}, \mathbf{y}) and of dimension d_1 and d_2 , respectively.
- ▶ Let θ be any aggregate function from bags of elements in \mathbb{R}^{d_1} to \mathbb{R}^d .
- ▶ Then,

$$\varphi(\mathbf{x}) := \text{agg}_{\mathbf{y}}^{\theta}(\varphi_1(\mathbf{x}, \mathbf{y}) \mid \varphi_2(\mathbf{x}, \mathbf{y}))$$

is an expression of dimension d and free variables \mathbf{x} .

We close GEL under aggregation with aggregations in Θ .

- ▶ Let $\varphi_1(\mathbf{x}, \mathbf{y})$ and $\varphi_2(\mathbf{x}, \mathbf{y})$ expressions with free variables (\mathbf{x}, \mathbf{y}) and of dimension d_1 and d_2 , respectively.
- ▶ Let θ be any aggregate function from bags of elements in \mathbb{R}^{d_1} to \mathbb{R}^d .
- ▶ Then,

$$\varphi(\mathbf{x}) := \text{agg}_{\mathbf{y}}^{\theta}(\varphi_1(\mathbf{x}, \mathbf{y}) \mid \varphi_2(\mathbf{x}, \mathbf{y}))$$

is an expression of dimension d and free variables \mathbf{x} .

Semantics

$$\xi_{\varphi} : (G, \mathbf{v}) \mapsto \theta\left(\left\{\left\{\xi_{\varphi_1}(G, \mathbf{v}, \mathbf{w}) \mid \mathbf{w} \in V_G^p \text{ s.t. } \xi_{\varphi_2}(G, \mathbf{v}, \mathbf{w}) \neq \mathbf{0}\right\}\right\}\right).$$

with $p := |\mathbf{y}|$.

Important special fragments:

- ▶ $GEL_k(\Omega, \Theta)$: **only k variables** x_1, \dots, x_k may be used;
- ▶ $GGEL_2(\Omega, \Theta)$: **guarded fragment** of GEL_2 in which aggregation and function application are restricted = $MPNN(\Omega, \Theta)$

MPNN(Ω, Θ) GEL₂(Ω, Θ) GEL₃(Ω, Θ) GEL_k(Ω, Θ)

GraphSage GINs GCNs SGNs

GATs GatedGCNs extended GINs 2-IGNs ChebNet ...

Walk GNNs 2WL-GNNs ring-GNNs 1-Dropout GNNs

Id-aware GNNs CayleyNet 3-IGNs 2-FGNNs ...

kWL-GNNs k-FGNNs (k+1)-IGNs GSNs k-Dropout GNNs ...

For many of these GNNs, their **layer definitions** translate **naturally** into expressions in our language.

What can we say about separation power of $\text{GEL}_k(\Omega, \Theta)$?

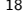
The k -dimensional Weisfeiler-Leman algorithm:^{18,19} Iteratively computes a coloring of k -tuples of vertices of a graph

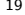
Intuitively, it can be seen as color refinement on a k -fold product of a graph

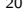
Again, has been subject to many theoretical studies and is used in graph isomorphism algorithms

$\rho(k\text{-WL})$ contains pairs of k -tuples of vertices with the same k -WL coloring

It is known:²⁰ $\rho(\text{CR}) \supseteq \rho(1\text{-WL}) \not\supseteq \rho(2\text{-WL}) \not\supseteq \rho(3\text{-WL}) \not\supseteq \dots \not\supseteq \rho(\text{graph iso})$.

¹⁸  Grohe. The logic of graph neural networks LICS, 1–17 (2021)

¹⁹  Morris, Lipman, Maron, Rieck, Kriege, Grohe, Fey, Borgwardt. Weisfeiler and Leman go Machine Learning: The Story So Far. CoRR abs/2112.09992 (2021)

²⁰  Cai, Fürer, Immerman: An optimal lower bound on the number of variables for graph identification. Comb. 12(4):389-410 (1992)

Theorem

For any Ω and Θ , $\rho(k\text{-WL}) \subseteq \rho(GEL_{k+1}(\Omega, \Theta))$


Proofs²¹ rely on connections to logics and techniques from database theory^{22,23}

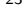
Theorem

if Ω contains concatenation, linear combinations and non-linear activation functions and Θ consists of summation, then $\rho(k\text{-WL}) = \rho(GEL_{k+1}(\Omega, \Theta))$

Approximation properties can be derived as well.

²¹  G., Reutter, Expressiveness and approximation properties of GNNs, ICLR (2022)

²²  Hella, Libkin, Nurmonen, Wong: Logics with Aggregates, JACM, 48(4): 880-907 (2001)

²³  Cai, Fürer, Immerman: An optimal lower bound on the number of variables for graph identification. Comb. 12(4):389-410 (1992)

color refinement 1-WL 2-WL k-WL

GraphSage GINs GCNs SGNs

GATs GatedGCNs extended GINs 2-IGNs ChebNet ...

Walk GNNs 2WL-GNNs ring-GNNs 1-Dropout GNNs

Id-aware GNNs CayleyNet 3-IGNs 2-FGNNs ...

kWL-GNNs k-FGNNs (k+1)-IGNs GSNs k-Dropout GNNs ...


6. What's Next?

1. Impact of different aggregation functions on expressive power.

- ▶ In initial work is investigated when summation MPNNs can be approximated by mean or max MPNNs, and vice versa.²⁴

2. Quantitative approximation results.

- ▶ What is complexity of embeddings needed to approximate within ϵ ?


²⁴  Rosenbluth, Tönshoff and Grohe: Some Might Say All You Need Is Sum, arxiv 2023.

1. Impact of different aggregation functions on expressive power.

- ▶ In initial work is investigated when summation MPNNs can be approximated by mean or max MPNNs, and vice versa.²⁴

2. Quantitative approximation results.

- ▶ What is complexity of embeddings needed to approximate within ϵ ?

²⁴  Rosenbluth, Tönshoff and Grohe: Some Might Say All You Need Is Sum, arxiv 2023.

Technique presented only gives **upper bounds**. Lower bounds, still case by case analysis.

3. Can we find “**reductions**” between embedding methods that **preserve expressive power**.

- ▶ This may help to show lower bounds by simulating “hard” embedding methods.

4. Fining the **minimal k** in $\text{GEL}_k(\Omega, \Theta)$ needed for your method.

- ▶ The lower k the better the upper bound.
- ▶ (Semantic) treewidth notion for GEL expressions?²⁵
- ▶ Related to work on FAQ-AIs when functions and aggregations can be seen as semiring operators.²⁶
- ▶ We need to accommodate for non-linear activation functions.

²⁵ G* and Reutter. Expressiveness and Approximation Properties of GNNs. ICLR 2022

²⁶ Khamis, Ngo and Rudra: "FAQ: Questions Asked Frequently", PODS 2016

Technique presented only gives **upper bounds**. Lower bounds, still case by case analysis.

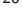
3. Can we find “**reductions**” between embedding methods that **preserve expressive power**.

- ▶ This may help to show lower bounds by simulating “hard” embedding methods.

4. Finding the **minimal k** in $\text{GEL}_k(\Omega, \Theta)$ needed for your method.

- ▶ The lower k the better the upper bound.
- ▶ (Semantic) treewidth notion for GEL expressions?²⁵
- ▶ Related to work on FAQ-AIs when functions and aggregations can be seen as semiring operators.²⁶
- ▶ We need to accommodate for non-linear activation functions.


²⁵  G* and Reutter. Expressiveness and Approximation Properties of GNNs. ICLR 2022

²⁶  Khamis, Ngo and Rudra: "FAQ: Questions Asked Frequently", PODS 2016

5. Other hierarchies than Weisfeiler-Leman hierarchy.

- ▶ Some embedding methods can be cast in $\text{GEL}_k(\Omega, \Theta)$ but are not as expressive as $(k-1)$ -WL.
- ▶ E.g., Reconstruction GNNs, ESANs, ID-Aware GNNs, Nested GNNs, Dropout-GNNs.
- ▶ By imposing further restrictions on expressions in $\text{GEL}_k(\Omega, \Theta)$ a more fine-grained hierarchy can be obtained. ^{27, 28}
- ▶ More work needed is to understand corresponding separation power. For example, hom count characterizations.


²⁷  Qian, Rattan, G*, Morris and Niepert. Ordered Subgraph Aggregation Networks. NeurIPS 2022


²⁸  Rattan and Seppelt: Weisfeiler-Leman and Graph Spectra. SODA 2023


6. We can see embedding methods as **views**.

- ▶ **Query rewriting** using such views?
- ▶ **View embedding**: First embed graph using **complex fixed embedding**, followed by **simple learnable embedding** of the view.²⁹
- ▶ Other?

- ▶ **Generalization** properties (VC dimension, Rademacher complexity)
- ▶ **Semi-ring** valued embeddings and learning?
- ▶ **Zero-one laws** of embeddings ³⁰
- ▶ **Continuous WL** ³¹
- ▶ More connections with **logic** and **descriptive complexity**. ³²

³⁰  Adam-Day, Iliant and Ceylan: Zero-One Laws of Graph Neural Networks, arxiv 2023


³¹  Böker, Levie, Huang, Villar and Morris: Fine-grained Expressivity of Graph Neural Networks, arxiv 2023

³²  Grohe. The Descriptive Complexity of Graph Neural Networks, arxiv 2023

And of course, as also mentioned in Jure's keynote:

Relational embeddings.

- ▶ Initial work by considering multi-relation graphs and analyzing power.³³

³³  Barceló, Galkin, Morris, Orth: Weisfeiler and Leman Go Relational Learning on Graphs (LoG) 2022.

- ▶ The $\text{GEL}(\Omega, \Theta)$ language and use for graph embeddings is heavily influenced by earlier work on the **expressive power of linear algebra and matrix query languages** together with Robert Brijder, Thomas Muñoz, Cristian Riveros, Jan Van den Bussche, and Domagoj Vrgoč.
- ▶ **Transferral to ML**: Pablo Barceló, Martin Grohe, Christopher Morris, Gaurav Rattan, Juan Reutter, Jasper Steegmans and Jan Van den Bussche.

- ▶ There is interest in ML community for these kind of theoretical analyses (but preferably accompanied with some experiments).
- ▶ Great opportunity for our community to contribute.
- ▶ So far, these papers are in ML conference. Would be great to also have some at PODS or ICDT!
- ▶ Maybe you got some inspiration for doing so :-)