# Graph and Subgraph Isomorphism Using GNNs

An overview of the essential concepts in Stanford CS224W (Lectures 9 and 12) with only oversimplified examples

## Ng Yen Kaow

# Graph Isomorphism

- Complexity of graph isomorphism is unknown
- Weisfeiler-Lehman graph kernel traditionally used to obtain graph-level embedding
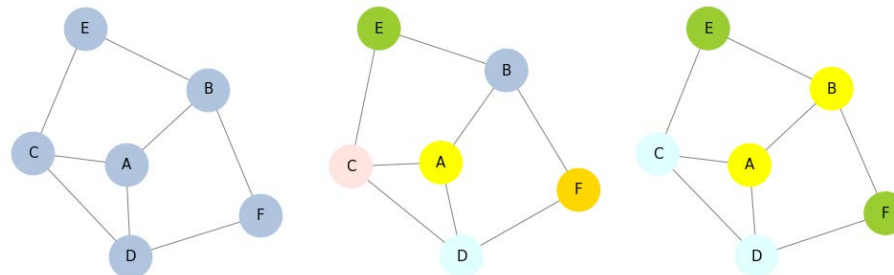
WL color-refinement algorithm
1. Assign initial color $c^{(0)}(v)$ to each node $v$
2. Iteratively refine node colors by
$$c^{(k+1)}(v) = \text{HASH}\left(c^{(k)}(v), \left\{c^{(k)}(u)\right\}_{u \in N(v)}\right)$$

- HASH function maps input to distinct values (colors)
- After $K$ steps, $c^{(K)}(v)$ summarizes the structure of the $K$-hop neighborhood

- To run as GNN, need to implement HASH as AGG
  - Need notion of distinguishing node embeddings

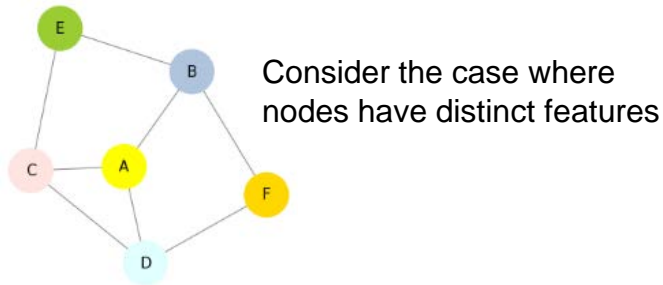# Distinguishing node embeddings

- ☐ Consider the task of keeping embeddings feature distinguishable

- ☐ Two factors affect embedding

  - ■ (Initial) feature
    - ☐ Feature representation will affect whether features can be converted into each other
      - ■ If $\boxed{1\ 0}$, $\boxed{0\ 1}$ and $\boxed{1\ 1}$ represent three distinct features, then the sum of $\boxed{1\ 0}$ and $\boxed{0\ 1}$ would become $\boxed{1\ 1}$, the third feature
    - ☐ To simplify this discussion assume that each distinct feature corresponds to one distinct dimension in the feature vector

  - ■ Neighborhood structure

- ☐ Nodes with the same (initial) **feature** and **neighborhood structure** should be assigned the same embedding, and vice versa
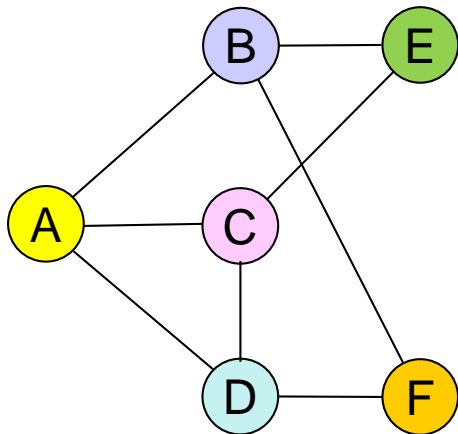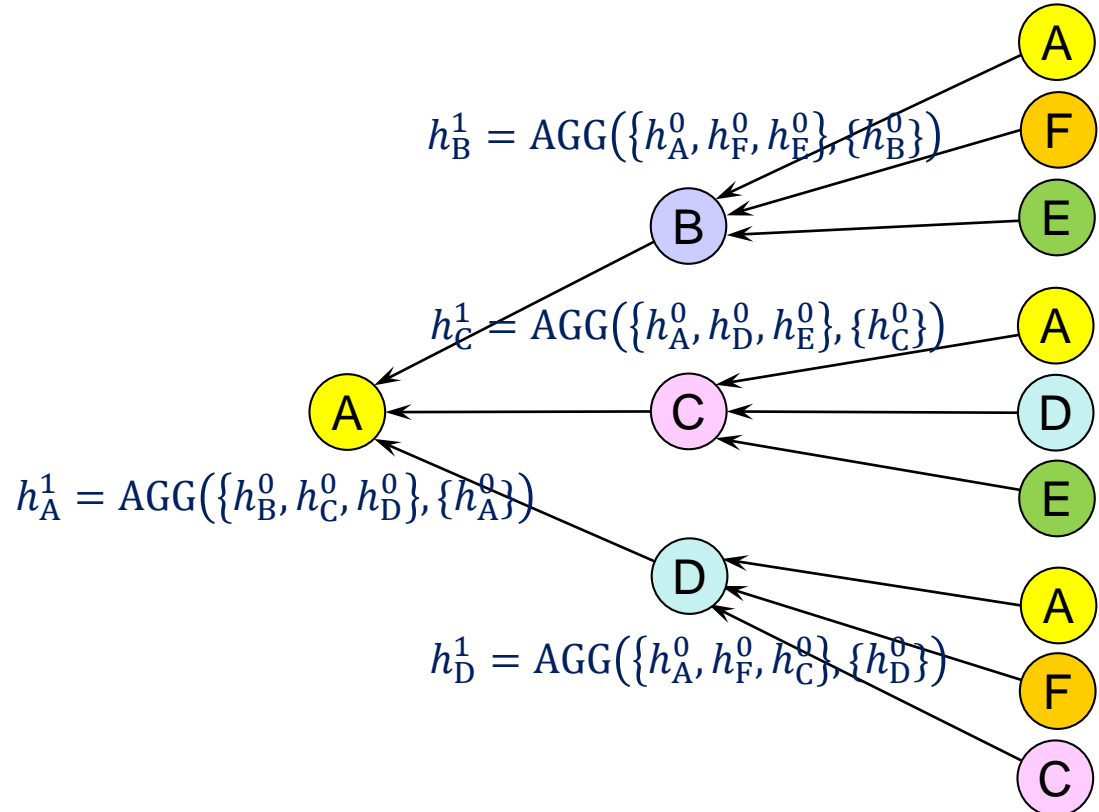
# Computing node embeddings

□ Recall that given a GCN of 2 layers, the embedding of A is computed through the computation graph as follows



Consider the case where nodes have distinct features

Computation graph of A's embedding

Rearranged w.r.t. distance from A

$$h_B^1 = AGG(\{h_A^0, h_F^0, h_E^0\}, \{h_B^0\})$$

$$h_C^1 = AGG(\{h_A^0, h_D^0, h_E^0\}, \{h_C^0\})$$

$$h_A^1 = AGG(\{h_B^0, h_C^0, h_D^0\}, \{h_A^0\})$$

$$h_D^1 = AGG(\{h_A^0, h_F^0, h_C^0\}, \{h_D^0\})$$

# Distinguishability under **sum**

☐ Let $\text{AGG}=$**sum**, then for the GCN of 2 layers, the embeddings are as follows

▪ Let $h_A^0=$ | 1 | 0 | 0 | 0 | 0 | 0 |, $h_B^0=$ | 0 | 1 | 0 | 0 | 0 | 0 |, $h_C^0=$ | 0 | 0 | 1 | 0 | 0 | 0 |,

$h_D^0=$ | 0 | 0 | 0 | 1 | 0 | 0 |, $h_D^0=$ | 0 | 0 | 0 | 0 | 1 | 0 |, $h_D^0=$ | 0 | 0 | 0 | 0 | 0 | 1 |, and let $\text{AGG}$ be **sum**. Then

☐ $h_B^1 = \text{AGG}\left(\{h_A^0, h_F^0, h_E^0\}, \{h_B^0\}\right) =$ | 1 | 1 | 0 | 0 | 1 | 1 |

☐ $h_C^1 = \text{AGG}\left(\{h_A^0, h_D^0, h_E^0\}, \{h_C^0\}\right) =$ | 1 | 0 | 1 | 1 | 1 | 0 |

☐ $h_D^1 = \text{AGG}\left(\{h_A^0, h_F^0, h_C^0\}, \{h_D^0\}\right) =$ | 1 | 0 | 1 | 1 | 0 | 1 |

☐ $h_A^1 = \text{AGG}\left(\{h_B^0, h_C^0, h_D^0\}, \{h_A^0\}\right) =$ | 1 | 1 | 1 | 1 | 0 | 0 |

▪ Finally the embedding of A is

☐ $h_A^2 = \text{AGG}\left(\{h_B^1, h_C^1, h_D^1\}, \{h_A^1\}\right) =$ | 4 | 2 | 3 | 3 | 2 | 2 |

▪ Similarly, $h_B^2 =$ | 2 | 4 | 2 | 2 | 2 | 2 |

$h_C^2 =$ | 3 | 2 | 4 | 3 | 2 | 1 |

$h_D^2 =$ | 3 | 2 | 3 | 4 | 1 | 2 | — Distinct embeddings

$h_E^2 =$ | 2 | 2 | 2 | 1 | 3 | 1 |

$h_F^2 =$ | 2 | 2 | 1 | 2 | 1 | 3 |

# Distinguishability under **sum**

- Let AGG=**sum**, then for the GCN of 2 layers, the embeddings are as follows

Summarized



E 0 0 0 0 1 0   B 0 1 0 0 0 0
C 0 0 1 0 0 0   A 1 0 0 0 0 0   F 0 0 0 0 0 1
D 0 0 0 1 0 0

E 0 1 1 0 1 0   B 1 1 0 0 1 1
C 1 0 1 1 1 0   A 1 1 1 1 0 0   F 0 1 0 1 0 1
D 1 0 1 1 0 1

E 2 2 2 1 3 1   B 2 4 2 2 2 2
C 3 2 4 3 2 1   A 4 2 3 3 2 2   F 2 2 1 2 1 3
D 3 2 3 4 1 2

- By induction on the respective distinct feature dimension, the embeddings will be distinct for subsequent iterations

- **If every node has distinct feature, then they have distinct embeddings under sum** regardless of neighborhood structure or iterations (with the exception of the graph of only two nodes A—C )

# Distinguishability under **sum**

- **If some nodes have the same features?**

  - Let $h_A^0 = h_B^0 = \boxed{1\ 0\ 0}$, $h_C^0 = h_D^0 = \boxed{0\ 1\ 0}$, $h_E^0 = h_F^0 = \boxed{0\ 0\ 1}$, and let $\mathrm{AGG}$ be **sum**. Then, as can be seen from the following example



- Two nodes with the **same feature** will have the same embedding under **sum** if they have the same **neighborhood structure**
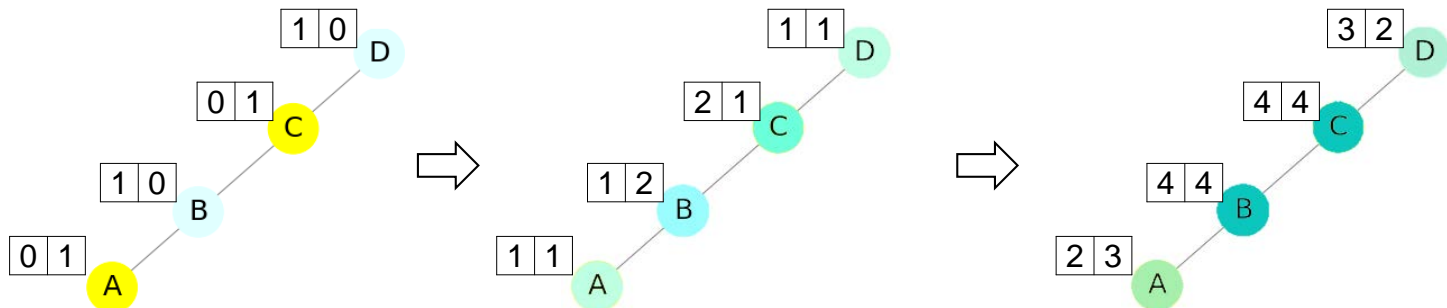
  - However different features and neighborhood structure cannot guarantee distinct embeddings
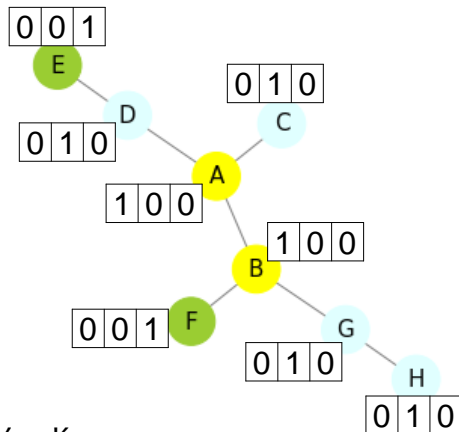
# Distinguishability under **sum**

- ☐ **If some nodes have the same features?**
  - ■ Different features and neighborhood structure cannot guarantee distinct embeddings for various reasons
    - ☐ Complementary neighborhood structure
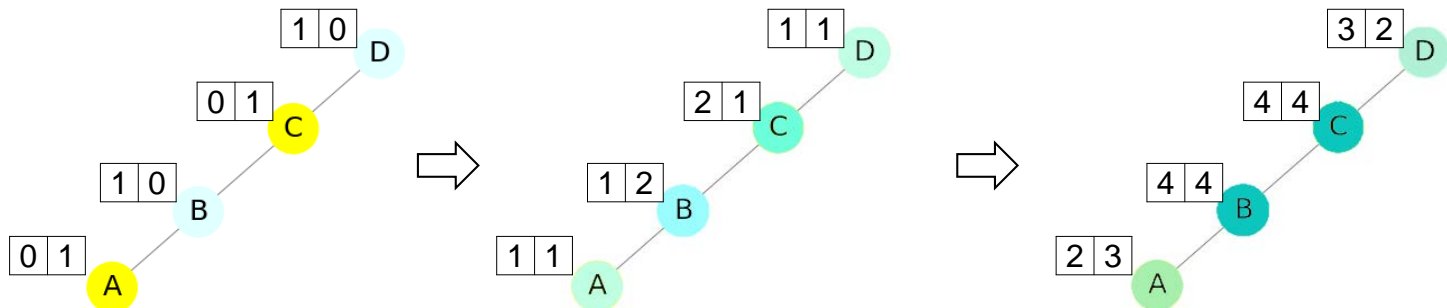


    - ☐ Hard-to-predict cases



Will A and B ever become the same again after the first step?
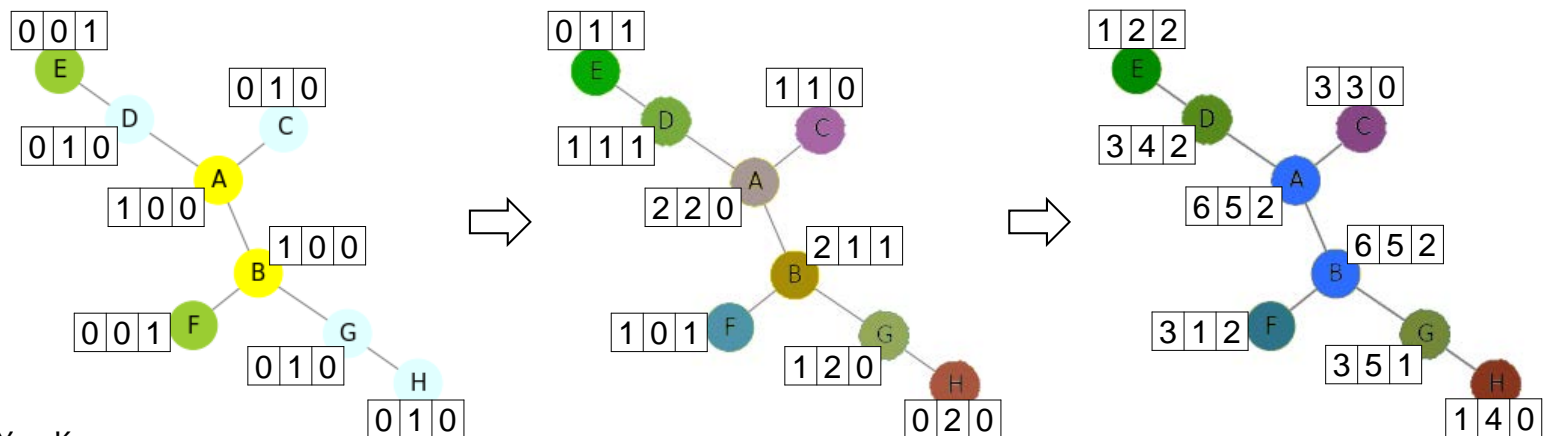
# Distinguishability under **sum**

☐ **If some nodes have the same features?**

■ Different features and neighborhood structure cannot guarantee distinct embeddings for various reasons
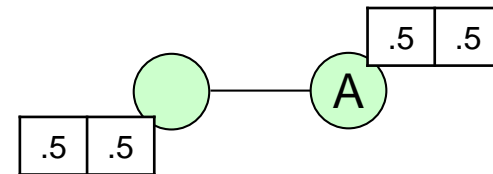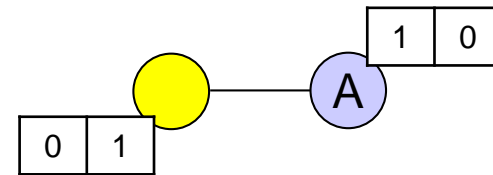
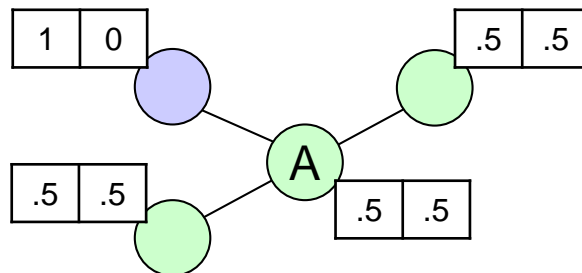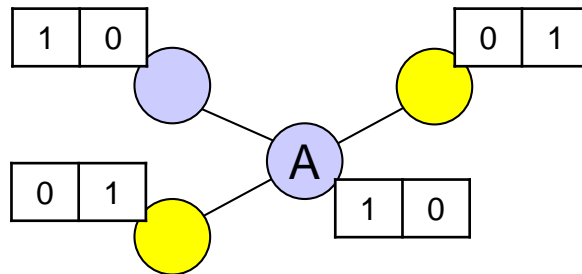☐ Complementary neighborhood structure
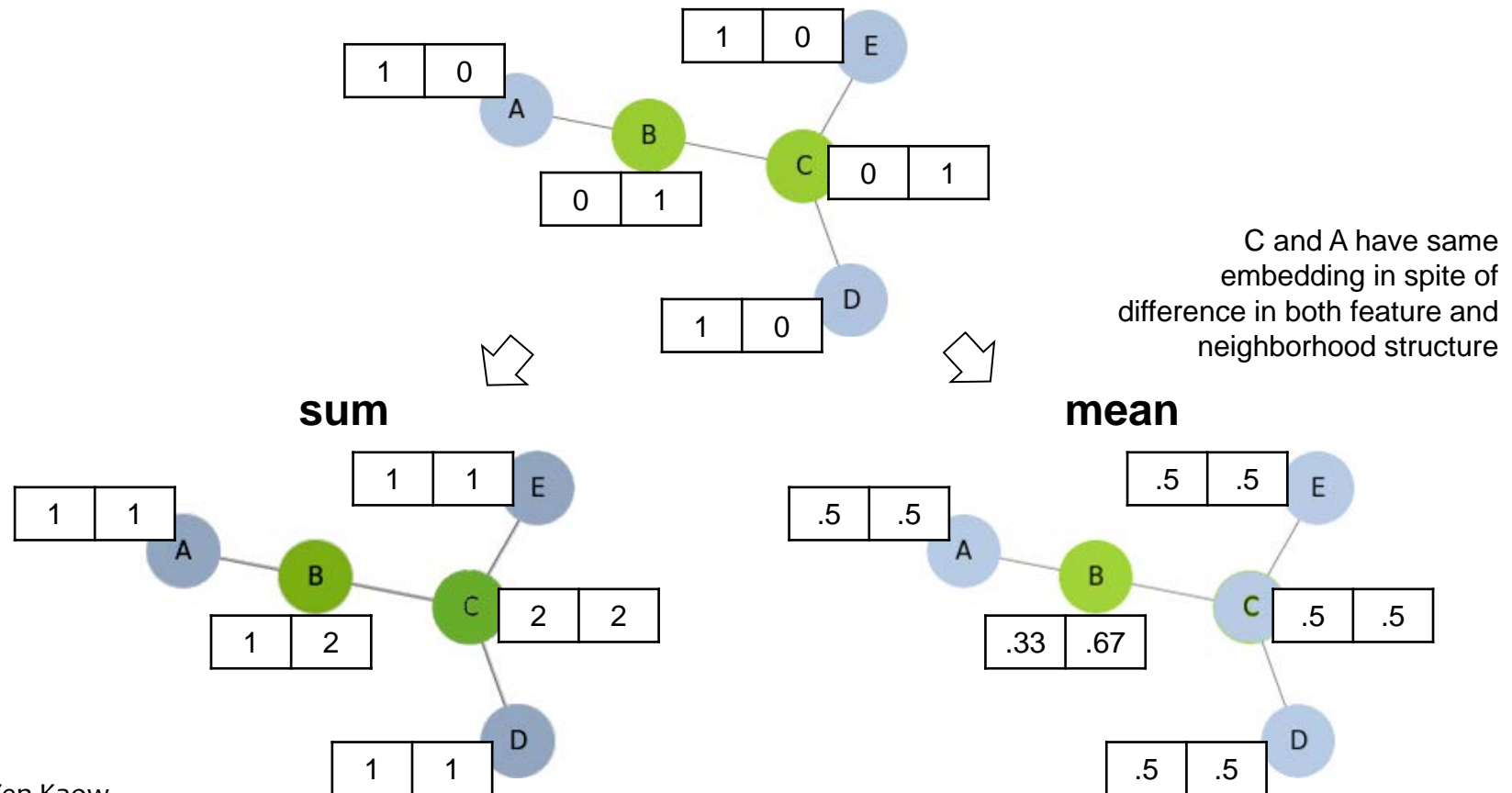


☐ Hard-to-predict cases

# Distinguishability under **mean**

□ Using **mean** as AGG results in even less desirable behavior

■ For instance, node A in both graphs below would give the same embedding under **mean** with one iteration

# Distinguishability under **mean**

- ☐ Using **mean** as AGG results in even less desirable behavior
  - ■ The following example shows **sum** to result in more consistent embeddings than **mean**



C and A have same embedding in spite of difference in both feature and neighborhood structure

**sum**

**mean**

# Injective function for isomorphism

- ☐ An **injective** function will output distinguishable embeddings for nodes of distinct feature and neighborhood structure
  - ■ **sum**, **mean**, and **max** are not injective

- ☐ **Theorem** (Xu *et al.* 2019). Any **injective** $\mathrm{AGG}$ function can be expressed as $\Phi(\sum_{x \in S} f(x))$ for some non-linear $\Phi$ and linear $f$

- ☐ Since MLP is able to approximate any function, we can learn $\Phi$ and $f$ with non-linear $\mathrm{MLP}_\Phi$ and linear $\mathrm{MLP}_f$

$$\mathrm{AGG} = \mathrm{MLP}_\Phi\left(\sum_{x \in S} \mathrm{MLP}_f(x)\right)$$

$\Rightarrow$ **Graph Isomorphism Network (GIN)**

# Subgraph Isomorphism

- Subgraph isomorphism is NP-complete

  ⇒ Compare **neighborhood around each node**
- The $k$-**hop neighborhood around node** $u \in Q$ is
  1. All the nodes within $k$ hops from $u$, and
  2. All the edges in between those nodes
     - Such a neighborhood is a subgraph of $Q$
       - However, not every subgraph of $Q$ is a neighborhood of some node $u \in Q$
         - At most $k|V|$ neighborhoods for each $k$
- If $P$ is a subgraph of $Q$, then every neighborhood of $P$ is a subgraph of some neighborhood of $Q$
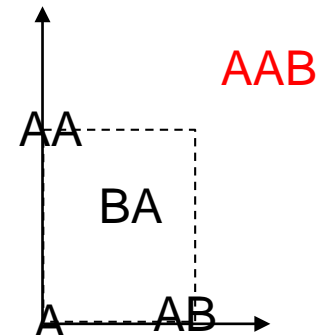
# Order embedding space

□ Idea: We want an $d$-dimensional embedding space $z$ such that for every neighborhoods $p \in P$ and $q \in Q$

$$q \subseteq p \Leftrightarrow \forall_{i=1}^{d} z_q[i] \leq z_p[i]$$

□ With embedding space, we can test subgraph isomorphism through the following

> For each neighborhood $p \in P$ and $q \in Q$
>   If $(\exists i)\left[z_q[i] \leq z_p[i]\right]$, return false
> Return true

# Order embedding space

- ☐ Idea: We want an $d$-dimensional embedding space $z$ such that for every neighborhoods $p \in P$ and $q \in Q$

$$q \subseteq p \Leftrightarrow \forall_{i=1}^{d} z_q[i] \leq z_p[i]$$

- ☐ Whether such an embedding space exist depends on $D$ and the class of graphs

  - ■ Even for substring relation ($\preccurlyeq$) with only alphabet {A, B}, a 2D embedding space is insufficient for

  $$q \preccurlyeq p \Leftrightarrow \forall_{i=1}^{2} z_q[i] \leq z_p[i]$$

    - • AAB needs to cover the dotted box (since it includes both AA and AB)
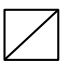    - • On the other hand, BA cannot be brought out of the dotted box (otherwise it would cover AA or AB)
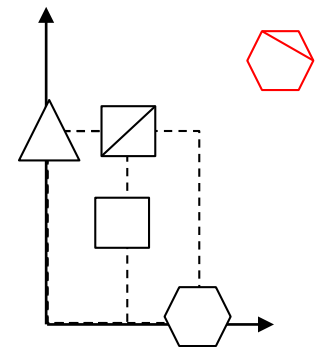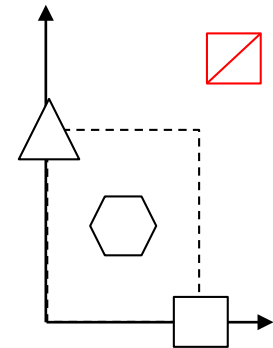
# Order embedding space

- 2D embedding space example for graph
    - ▨ needs to cover the dotted box (since it must cover both triangle and square) but that would cover ⬡
        - On the other hand, ⬡ cannot be brought out of the dotted box (otherwise it would cover △ or ☐)
    - Toggling hexagon and square will allow ▨ to be placed, but now ⬡ cannot be placed
- Assume that $d$ is sufficiently large for reasonable embeddings

# Training order embedding space

- ☐ Use a **node embedding space of** $k$-hop
- ☐ Denote the embedding of a node $u$ as $z_u$
- ☐ Use the loss function
$$\text{loss}(u, v) = \sum_{i=1}^{d} \max(0, z_u[i] - z_v[i])^2$$
It is clear that
  - ■ $\text{loss}(u, v) = 0$ when $\forall_{i=1}^{d}(z_u[i] \leq z_v[i])$
  - ■ $\text{loss}(u, v) > 0$ otherwise

- ☐ Generate random pair of graphs $P$, $Q$ and train GNN such that embeddings of $u_P \in P$ and $u_Q \in Q$ has
  - ■ $\text{loss}(u_Q, u_P) = 0$ when $u_Q \subseteq u_P$, and
  - ■ $\text{loss}(u_Q, u_P) > 0$ otherwise