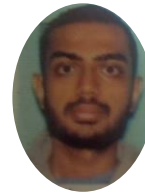


# Graph Representation Learning: Where Probability Theory, Data Mining, and Neural Networks Meet

Bruno Ribeiro  
Assistant Professor  
Department of Computer Science  
Purdue University



Joint work with R. Murphy\*, B. Srinivasan\*, V. Rao

GrAPL Workshop @ IPDPS  
May 20<sup>th</sup>, 2019

Sponsors:



Army Research Lab  
Network Science  
CTA



# Overview

---

- ▶ What is the most powerful<sup>+</sup> graph model / representation?
- ▶ How can we make model learning<sup>\$</sup> tractable<sup>\*</sup>?
  - How can we make model learning<sup>\$</sup> scalable?

<sup>+</sup> powerful → expressive

<sup>\*</sup> tractable → works on small graphs

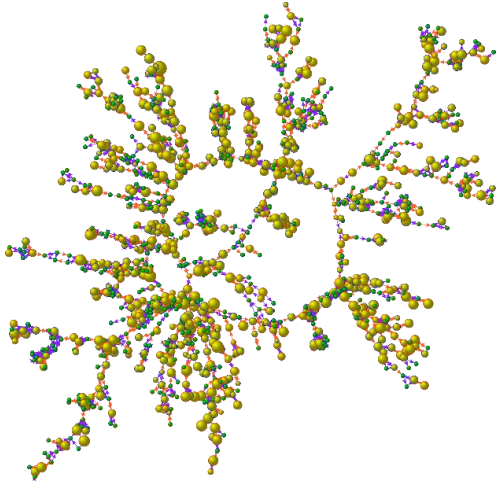
<sup>\$</sup> learning and inference

# What is a Graph\*?

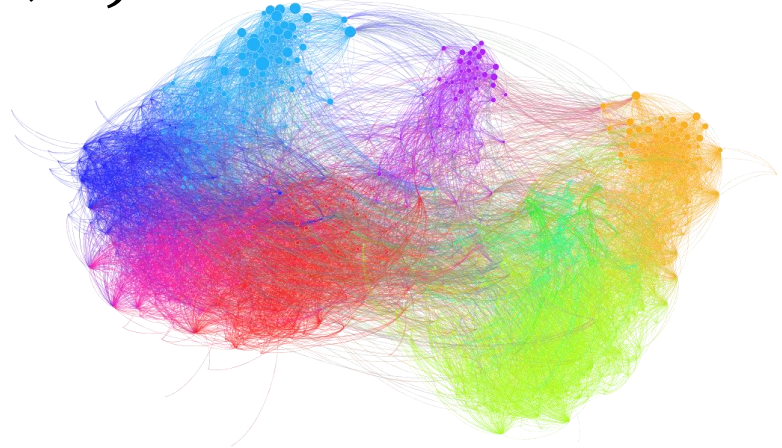
\* static

# Examples of (Static) Graphs

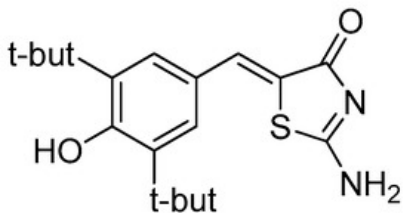
$$G = (V, E)$$



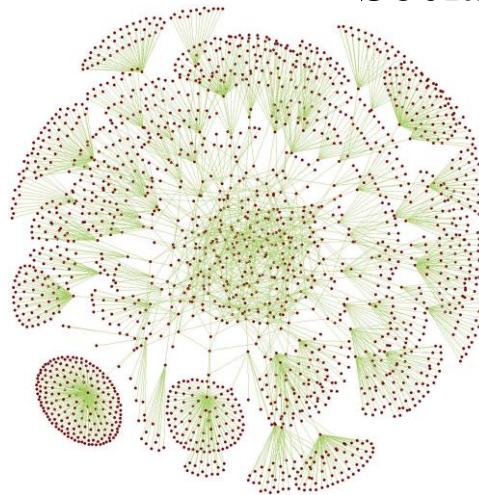
Biological Graphs



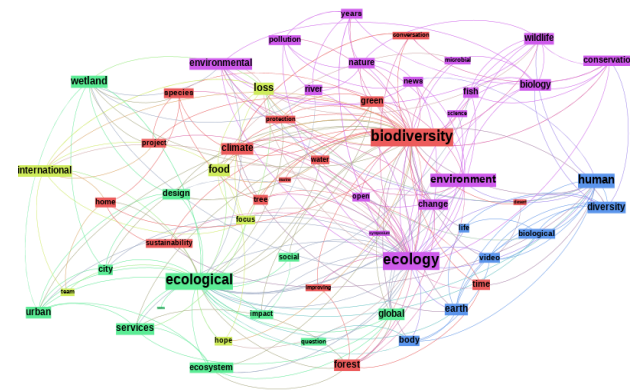
Social Graphs



Molecules

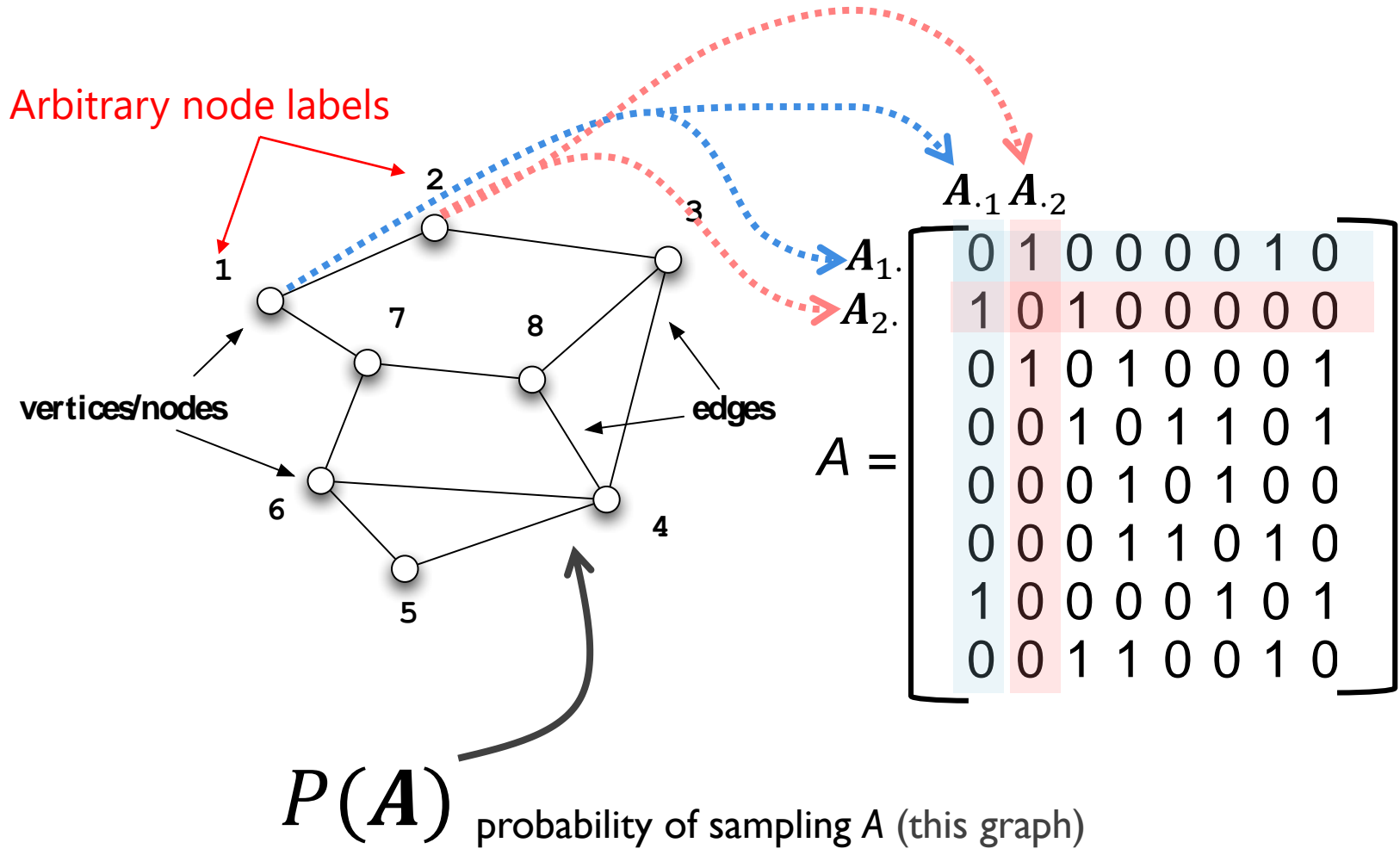


The Web



Ecological Graphs


# Statistical View of Graphs



Statistical definitions: from sequences to graphs....

# Definition: Joint Probability

- Consider a **sequence** of  $n$  random variables:

$$X_1, \dots, X_n \quad \text{with } X_i \in \Omega \quad \text{e.g. } \Omega = \left\{ \begin{array}{c} \text{countable} \\ \text{pile of words} \end{array} \right\}$$


with joint probability distribution

$$P(X_1, \dots, X_n)$$

- Sequence example: "The quick brown fox jumped over the lazy dog"

$$P(X_1 = \text{the}, X_2 = \text{quick}, \dots, X_9 = \text{dog})$$

- The joint probability is just a function

$$P: \Omega^n \rightarrow [0,1] \quad (\text{w/ normalization})$$

- $P$  takes an ordered sequence and outputs a value between zero and one (w/ normalization)

# Probabilities over Unordered Sequences (Multisets)

- Consider a **set** of  $n$  random variables (*representing a multiset*):

$$\begin{array}{ccc} X_2 & X_1 & X_4 \\ & X_3 & X_5 \end{array} \quad \text{with } X_i \in \Omega$$

how should we define their *joint probability distribution*?

**Recall:** Probability function  $P: \Omega^n \rightarrow [0,1]$  is order-dependent

- Definition:* For multisets the probability function  $P$  is such that

$$P(X_1, \dots, X_n) = P(X_{\pi(1)}, \dots, X_{\pi(n)})$$

is true for any permutation  $\pi$  of  $(1, \dots, n)$

Useful references:

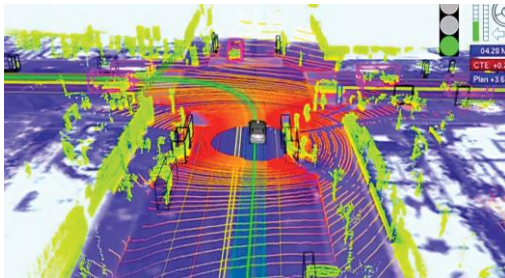
(Diaconis, *Synthese* 1977). Finite forms of de Finetti's theorem on exchangeability

(Murphy et al., ICLR 2019) Janossy Pooling: Learning Deep Permutation-Invariant Functions for Variable-Size Inputs



# Examples of Sets

- ▶ Point clouds



Lidar maps

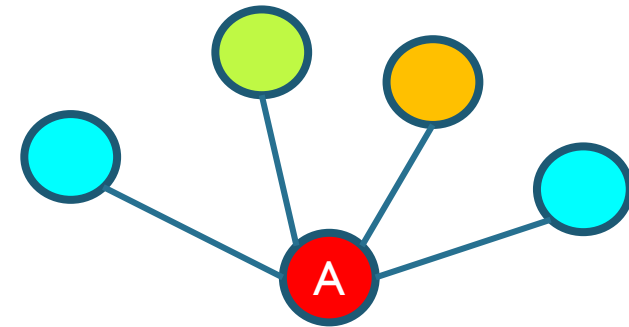
- ▶ Bag of words



- ▶ Our friends



- ▶ Neighbors of a node



extension: set-of-sets (Meng et al., KDD 2019)

# Probabilities over Array-ordered Sequences (Graphs)

- ▶ Consider an **array** of  $n^2$  random variables:

$$\begin{array}{ccccccc} & X_{21} & X_{11} & X_{22} & & & \\ & & & & \dots & & \\ & X_{12} & X_{nn} & & & & \\ & & & & & & X_{ij} \in \Omega \end{array}$$

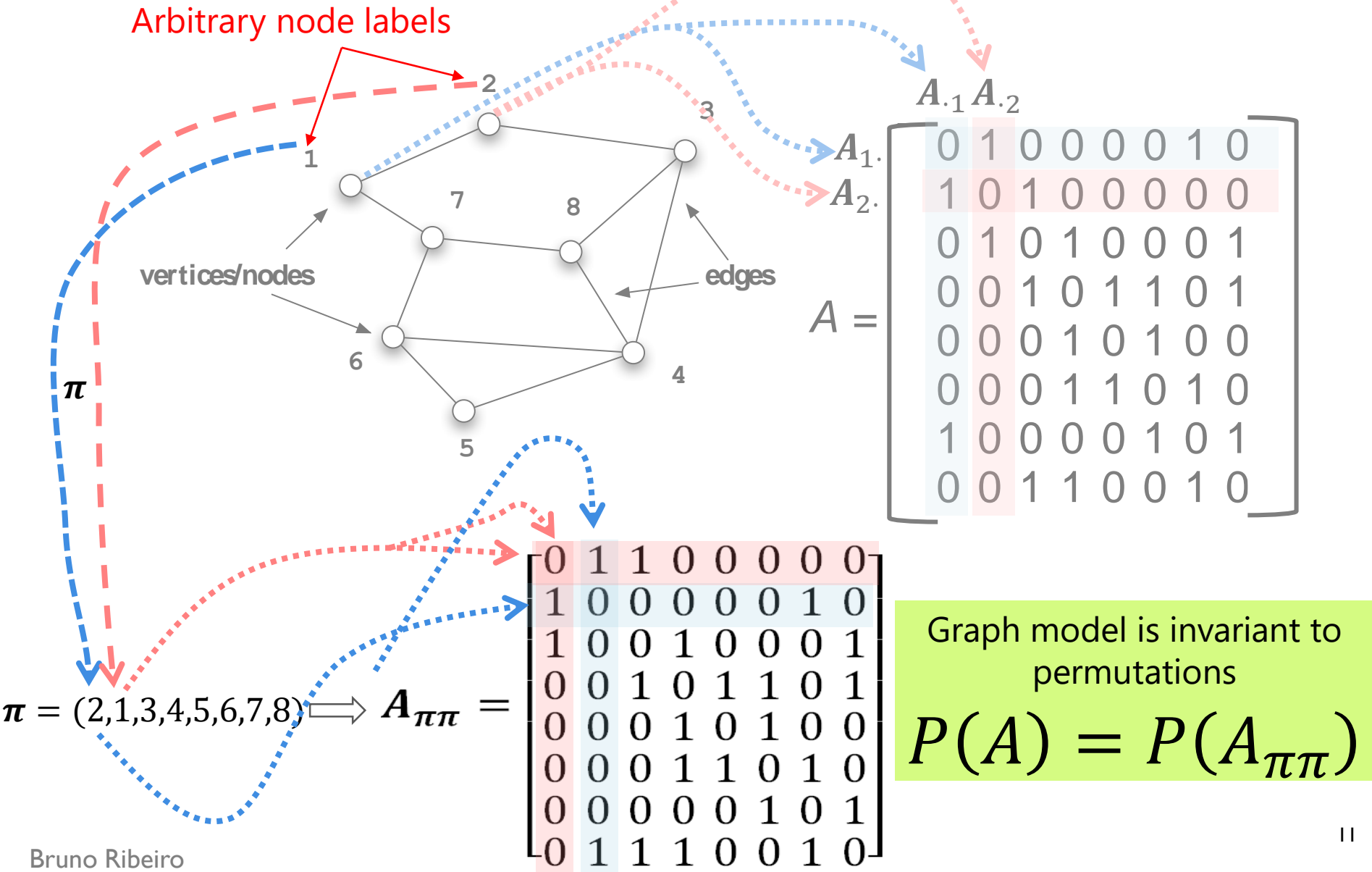
and  $P: \Omega^{n \times n} \rightarrow [0,1]$  such that

$$P(X_{11}, X_{12}, X_{21}, \dots, X_{nn}) = P(X_{\pi(1)\pi(1)}, X_{\pi(1)\pi(2)}, X_{\pi(2)\pi(1)}, \dots, X_{\pi(n)\pi(n)})$$

for any permutation  $\pi$  of  $(1, \dots, n)$

- ▶ Then,  $P$  is a model of a **graph** with  $n$  vertices, where  $X_{ij} \in \Omega$  are edge attributes (e.g., weights)
  - For each graph,  $P$  assigns a probability
  - Trivial to add node attributes to the definition
- ▶ If  $\Omega = \{0,1\}$  then  $P$  is a probability distribution over adjacency matrices
  - Most statistical graph models can be represented this way

# Statistical Model Invariant to Permutations of $A$



# Consequences of Invariances...

# A Historical Perspective

- ▶ Invariances have deep implications in nature
  - Noether's (first) theorem (1918):  
**invariances  $\Rightarrow$  laws of conservation**
  - e.g.:
    - time and space translation invariance  $\Rightarrow$  energy conservation
- ▶ The study of probabilistic invariances (symmetries) has a long history
  - Laplace's "rule of succession" dates to 1774 (Kallenberg, 2005)
  - Maxwell's work in statistical mechanics (1875) (Kallenberg, 2005)
  - Permutation invariance for **infinite sets**:
    - de Finetti's theorem (de Finetti, 1930)
    - Special case of the ergodic decomposition theorem, related to integral decompositions (see Orbanz and Roy (2015) for a good overview)
  - Kallenberg (2005) & (2007): de-facto references on probabilistic invariances



In 1981...

Aldous, D. J. Representations for partially exchangeable arrays of random variables. J. Multivar. Anal., 1981.

# Aldous-Hoover **Representation** Theorem (1979-1983)

- Consider an **infinite** set of random variables:

$$\begin{matrix} X_{21} & & X_{11} & & X_{22} & & \dots & & X_{ij} \\ & & X_{12} & & & & & & \end{matrix} \in \Omega$$

such that

$$P(X_{11}, X_{12}, \dots) = P(X_{\pi(1)\pi(1)}, X_{\pi(1)\pi(2)}, \dots)$$

is true for any permutation  $\pi$  of the positive integers

- Then,

$$P(X_{11}, X_{12}, \dots) \propto \int_{U_1 \in [0,1]} \cdots \int_{U_\infty \in [0,1]} \prod_{ij} P(X_{ij} | U_i, U_j)$$

is a mixture model of uniform distributions over  $U_i, U_j, \dots \sim \text{Uniform}(0,1)$

**(Aldous-Hoover representation is sufficient only for infinite graphs)**

# Connection Between Representation and Probability Theory



# Representations (*Sufficient Statistics*)

*Relationship between deterministic functions and probability distributions*

## ► **Noise outsourcing:**

- Tool from measure theory
- Any conditional probability  $P(Y|X)$  can be represent as  

$$Y = g(X, \epsilon), \quad \epsilon \sim \text{Uniform}(0,1)$$
 where  $g$  is a deterministic function
- The randomness is entirely outsourced to  $\epsilon$

## ► **Representation $s(X)$ :**

- $s(X)$ : deterministic function, makes  $Y$  independent of  $X$  given  $s(X)$
- Then,  $\exists g'$  such that  

$$(Y, X) = (g'(s(X), \epsilon), X), \quad \epsilon \sim \text{Uniform}(0,1)$$

We call  $s(X)$  a representation of  $X$

*Representations are generalizations of “embeddings”*

## Example Aldous-Hoover-type Representation in Graph Analysis....

# A Graph **Model** (based on Aldous-Hoover Representation)

*Gaussian Linear Model:* (each node  $i$  represented by a random vector)

Node  $i$  vector  $\mathbf{U}_i \sim \text{Normal}(\mathbf{0}, \sigma_U^2 \mathbf{I})$

Adjacency matrix:  $\mathbf{A}_{ij} \sim \text{Normal}(\mathbf{U}_i^T \mathbf{U}_j, \sigma^2)$

**Q:** For a given  $\mathbf{A}$ , what is the most likely  $\mathbf{U}$ ?

**Answer:**  $\mathbf{U}^* = \operatorname{argmax}_{\mathbf{U}} P(\mathbf{A}|\mathbf{U})$  , a.k.a. maximum likelihood

*Equivalent optimization:* Minimizing Negative Log-Likelihood:

$$\mathbf{U}^* = \operatorname{argmin}_{\mathbf{U}} \left\| \mathbf{A} - \mathbf{U}\mathbf{U}^T \right\|_2^2 + \frac{\sigma^2}{\sigma_U^2} \|\mathbf{U}\|_2^2$$



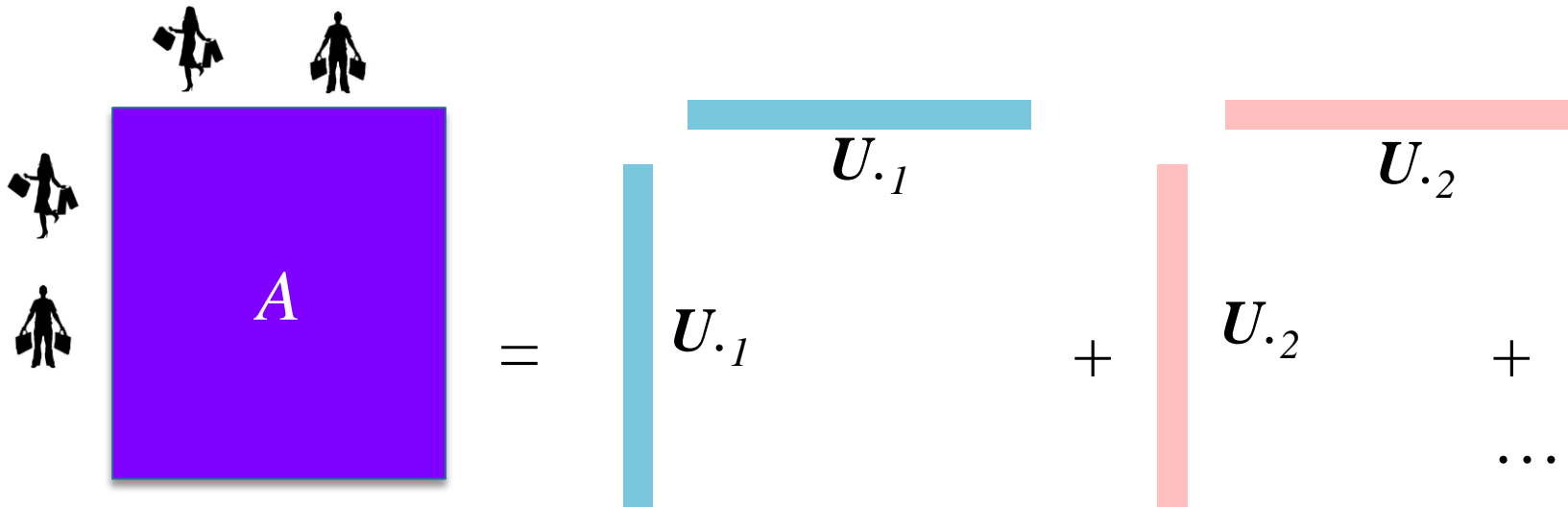
That will turn out to be the same

# Graph **Embedding** (Representation)

- ▶ Embedding of adjacency matrix  $A$

$$A \approx UU^T$$

$U_{\cdot i} = i$ -th column vector of  $U$



# Factorization for Low-rank **Representation** of a Graph

- ▶ Matrix factorization can be used to compute a **low-rank representation of A**
- ▶ A *reconstruction* problem:

Find

$$\mathbf{A} \approx \mathbf{U}\mathbf{U}^T$$

by optimizing

$$\min_{\mathbf{U}} \|\mathbf{A} - \mathbf{U}\mathbf{U}^T\|_2^2 + \lambda \|\mathbf{U}\|_2^2$$

Sum squared error

L2 regularization

Regularization strength

where  $\mathbf{U}$  has  $k$  columns\*

\*sometimes we will force orthogonal columns in  $\mathbf{U}$

Matrix factorization as an  
Aldous-Hoover representation

(Aldous-Hoover representation is sufficient only for infinite graphs)

⇒ Matrix factorization not sufficient for finite graphs

# Search for More Expressive Representations of Finite Graphs



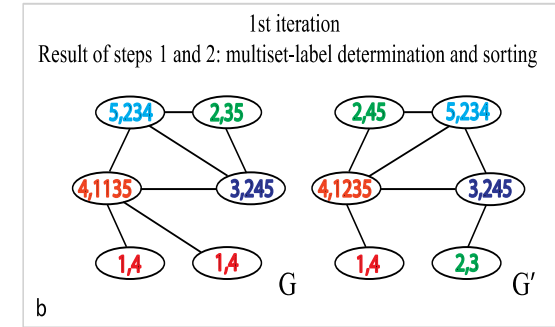
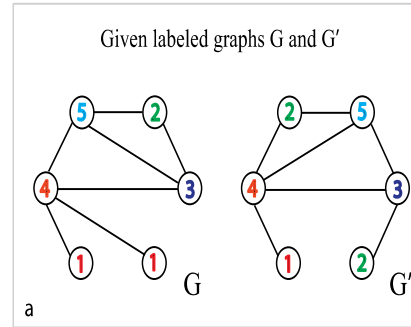
# Graph Neural Networks

(representations of **finite and infinite** graphs)

**DETOUR:** Isomorphism test for attributed graphs

# Weisfeiler-Lehmann Algorithm

- ▶ Recursive algorithm to determine if two graphs are isomorphic
  - Valid isomorphism test for most graphs (Babai and Kucera, 1979)
  - Cai et al., 1992 shows examples that cannot be distinguished by it
  - Belongs to class of color refinement algorithms that iteratively update vertex "colors" (hash values) until it has converged to unique assignments of hashes to vertices
  - Final hash values encode the structural roles of vertices inside a graph
  - Often fails for graphs with a high degree of symmetry, e.g. chains, complete graphs, tori and stars



Shervashidze et al. 2011

## Initialize:

$h_v$  is the attribute vector of vertex  $v \in G$   
(if no attribute, assign 1)  
 $k = 0$

## function WL-fingerprints( $G$ ):

**while** vertex attributes change **do**:

$k \leftarrow k + 1$

**for all** vertices  $v \in G$  **do**

$h_{k,v} \leftarrow \text{hash}(h_{k-1,v}, \{h_{k-1,u} : \forall u \in \text{Neighbors}(v)\})$

**Return**  $\{h_{k,v} : \forall v \in G\}$

neighbors of node  $v$

# Graph Representation Idea

- ▶ The hardest task for graph representation is:
  - Give different *tags* to different graphs
    - Isomorphic graphs should have the same *tag*
  - **Task:** Given adjacency matrix  $\mathbf{A}$ , predict *tag*
    - *Goal:* Find a representation  $s(\mathbf{A})$  such that
$$P(\text{tag}|\mathbf{A}) = g(s(\mathbf{A}), \epsilon)$$
- Then,  $s(\mathbf{A})$  must give:
  - same representation to isomorphic graphs
  - different representations to non-isomorphic graphs

# Back to Graph Neural Networks

# Graph Neural Networks (GNNs)

## Main idea:

Graph Neural Networks: Use the WL algorithm to *compute representations* that are *related* to a **task**

**Initialize**  $h_{0,v}$  = node  $v$  attribute

**function**  $\vec{f}(A, W_1, \dots, W_K, b_1, \dots, b_K)$ :

**while**  $k < K$  **do**: # K layers

$k \leftarrow k + 1$

**for all** vertices  $v \in V$  **do**

$$h_{k,v} = \sigma(W_k(h_{k-1,v}, \underbrace{A_v \cdot h}_{\text{could be another permutation-invariant function (see Murphy et al. ICLR 2019)}}) + b_k)$$

**return**  $\{h_{K,v} : \forall v \in V\}$

Example supervised task: predict label  $y_i$  of graph  $G_i$  represented by  $A_i$

Optimization for loss  $L$ : Let  $\theta = (W_1, \dots, W_K, b_1, \dots, b_K, W_{\text{agg}}, b_{\text{agg}})$

$$\theta^* = \operatorname{argmax}_{\theta} \sum_{i \in \text{Data}} L(y_i, W_{\text{agg}} \operatorname{Pooling}(\vec{f}(A_i, W_1, \dots, W_K, b_1, \dots, b_K)) + b_{\text{agg}})$$

permutation-invariant function  
(see Murphy et al. ICLR 2019)

# Graph Neural Network (GNN) Expressive Power

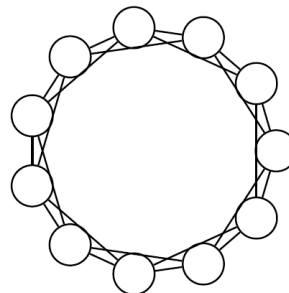
**GNN** representations can be **as expressive as** the **Weisfeller-Lehman** (WL) isomorphism test  
(Xu et al., ICLR 2019)

By construction, GNN representation  $\vec{f}$  is guaranteed permutation invariance (equivariance)

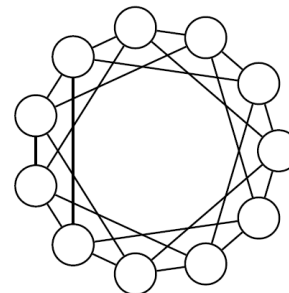
$$\vec{f}(\mathbf{A}) = \vec{f}(\mathbf{A}_{\pi\pi})$$

But WL test can sometimes fail ....

$$M = 11, L = 2$$



$$M = 11, L = 3$$



E.g. in a family of  
circulant graphs:

Figure 1: Two non-isomorphic graphs deemed isomorphic by the WL test.

Are there even more expressive graph representations?



# Universal Graph Approximator?

---

- ▶ Multilayer perceptron (MLP) is universal function approximator

(Hornik et al. 1989)

- What about using  $\vec{f}_{\text{MLP}}(\text{vec}(\mathbf{A}))$ ?

- **No!** Permutation-sensitive\*

$$\vec{f}_{\text{MLP}}(\text{vec}(\mathbf{A})) \neq \vec{f}_{\text{MLP}}(\text{vec}(\mathbf{A}_{\pi\pi}))$$

for some permutation  $\pi$

\* unless neuron weights nearly all the same (Maron et al., 2018)

**Back** to the **Definition** of a Graph Model

- ▶ Extension: A graph model  **$P$**  is an **array** of  $n^2$  random variables,  $n > 1$ ,

$$\begin{array}{ccccccc} X_{21} & X_{11} & X_{22} & & & & \\ & X_{12} & X_{nn} & \dots & X_{ij} & \in \Omega \end{array}$$

and  $P: \Omega^{\cup \times} \rightarrow [0,1]$ , where  $\Omega^{\cup \times} \equiv \bigcup_{i=2}^{\infty} \Omega^{i \times i}$ , such that

$$P(X_{11}, X_{12}, X_{21}, \dots, X_{nn}) = P(X_{\pi(1)\pi(1)}, X_{\pi(1)\pi(2)}, X_{\pi(2)\pi(1)}, \dots, X_{\pi(n)\pi(n)})$$

for any value of  $n$  and any permutation  $\pi$  of  $(1, \dots, n)$

- ▶ **(Murphy et al., ICML 2019) insight:**

$P$  is average of an unconstrained probability function  $\vec{P}$  applied over the Abelian group defined by the permutation operator  $\pi$

- Average is invariant to the group action of a *permutation*  $\pi$  (see Bloem-Reddy & Teh, 2019)
- *Works for variable size graphs*

# Relational Pooling (RP) (Murphy et al., ICML 2019)

- ▶  $\mathbf{A}$  is a tensor encoding: adjacency matrix & edge attributes
- ▶  $\mathbf{X}^{(v)}$  encodes node attributes
- ▶  $\Pi$  is the set of all permutation of  $(1, \dots, |V|)$ , where  $|V|$  is number of vertices
- ▶  $\vec{f}$  is any permutation-sensitive function

$$\bar{\bar{f}}(\mathbf{A}) = E_{\pi} \left[ \vec{f}(\mathbf{A}_{\pi\pi}, \mathbf{X}_{\pi}^{(v)}) \right] = \frac{1}{|V|!} \sum_{\pi \in \Pi} \vec{f}(\mathbf{A}_{\pi\pi}, \mathbf{X}_{\pi}^{(v)})$$

average over  $\pi \sim \text{Uniform}(\Pi)$

- ▶ **Theorem 2.1:** Necessary and sufficient representation of finite graphs
  - (details)  $\vec{f}$  is an universal approximator (MLP, RNNs), then  $\bar{\bar{f}}(\mathbf{A})$  is the most expressive representation of  $\mathbf{A}$

# Approximating Intractable Sum (Murphy et al., ICML 2019)

$$\bar{\bar{f}}(\mathbf{A}) \propto \sum_{\pi \in \Pi} \vec{f}(\mathbf{A}_{\pi\pi}, \mathbf{X}_{\pi}^{(v)})$$

1. Canonical orientation (some order of the vertices), so that  $\text{canonical}(\mathbf{A}) = \text{canonical}(\mathbf{A}_{\pi\pi})$
2. **k**-ary dependencies:
  - Nodes k-by-k independent in  $\vec{f}$
  - $\vec{f}$  considers only the first  $k$  nodes of any permutation  $\pi$
3. Stochastic optimization (proposes  $\pi$ -SGD)

# 1. Tractability through Canonical Orientations

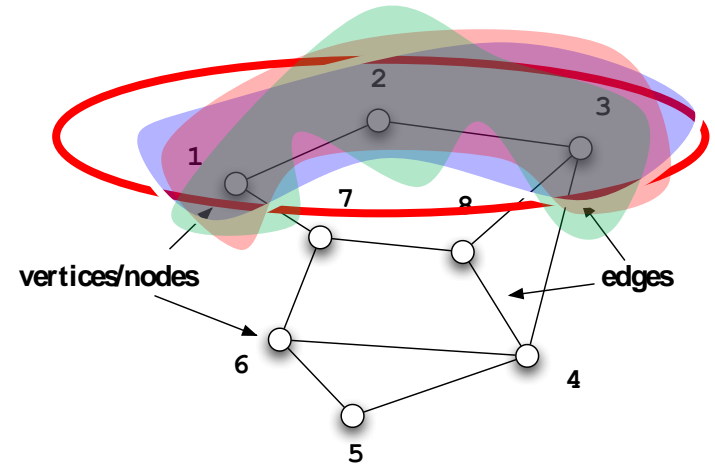
---

$$\bar{\bar{f}}(\mathbf{A}) \propto \sum_{\pi \in \Pi} \vec{f}(\mathbf{A}_{\pi\pi}, \mathbf{X}_{\pi}^{(v)})$$

- ▶ Order nodes with a **sort** function
  - E.g.: order nodes by PageRank
- ▶ Arrange  $\mathbf{A}$  with  $\text{sort}(\mathbf{A})$  (*assuming no ties*)
  - Note that  $\text{sort}(\mathbf{A}) = \text{sort}(\mathbf{A}_{\pi\pi})$  for any permutation  $\pi$

## 2. Tractability through k-ary Dependencies

$$\bar{f}(A) \propto \sum_{\pi \in \Pi} \vec{f}(A_{\pi\pi}, X_{\pi}^{(v)})$$



**k-ary dependencies:**

- Nodes k-by-k independent in  $\vec{f}$
- $\vec{f}$  considers only the first  $k$  nodes of any permutation  $\pi$ :  $\binom{n}{k}$  permutations

$$\pi = (1, 2, 3, 4, 5, 6, 7, 8)$$

$$\pi = (2, 1, 3, 4, 5, 6, 7, 8)$$

$$\pi = (2, 1, 3, 4, 5, 6, 8, 7)$$

$$\sum_{\pi \in \Pi} \vec{f}(A_{\pi\pi}) = \vec{f}(\cdot) + \vec{f}(\cdot) + \dots + \vec{f}(\cdot) + \dots$$

$$A_{\pi\pi} = \begin{bmatrix} 0 & 1 & 0 & 0 & 0 & 0 & 1 & 0 \\ 1 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 1 & 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 & 1 & 1 & 0 & 1 \\ 0 & 0 & 0 & 1 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 1 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 & 1 \\ 0 & 1 & 1 & 1 & 0 & 0 & 1 & 0 \end{bmatrix}$$

$$A_{\pi\pi} = \begin{bmatrix} 0 & 1 & 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \\ 1 & 0 & 0 & 1 & 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 & 1 & 1 & 0 & 1 \\ 0 & 0 & 0 & 1 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 1 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 & 1 \\ 0 & 1 & 1 & 1 & 0 & 0 & 1 & 0 \end{bmatrix}$$

$$A_{\pi\pi} = \begin{bmatrix} 0 & 1 & 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 1 & 0 & 0 & 1 & 0 \\ 0 & 0 & 1 & 0 & 1 & 1 & 1 & 0 \\ 0 & 0 & 0 & 1 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 1 & 0 & 0 & 1 \\ 0 & 1 & 1 & 1 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 & 1 & 1 & 0 \end{bmatrix}$$

### 3. Tractability through Stochastic Optimization ( $\pi$ -SGD)

- ▶ SGD: standard Stochastic Gradient Descent
  1. SGD will sample a batch of  $n$  training examples
  2. Compute gradients (backpropagation using chain rule)
  3. Update model following negative gradient (one gradient descent step)
  4. GOTO 1:
  
- ▶  $\pi$ -SGD (as fast as SGD per gradient step)
  1. Sample a batch of training examples
  2. For each example  $\mathbf{x}^{(j)}$  in the batch
    - Sample one permutation  $\pi^{(j)}$
  3. Perform a forward pass over the examples with the single sampled permutation
  4. Compute gradients (backpropagation using chain rule)
  5. Update model following negative gradient (one gradient descent step)
  6. GOTO 1:

### 3. $\pi$ -SGD guarantees

---

- ▶ **Proposition 2.1** (Murphy et al., ICML 2019)
  - $\pi$ -SGD behaves just like SGD
  - If loss is MSE, cross-entropy, negative log-likelihood, then  $\pi$ -SGD is minimizing an upper bound of the loss
  - However, the solution  $\pi$ -SGD converges to is not the solution of SGD.
    - But still a valid graph representation



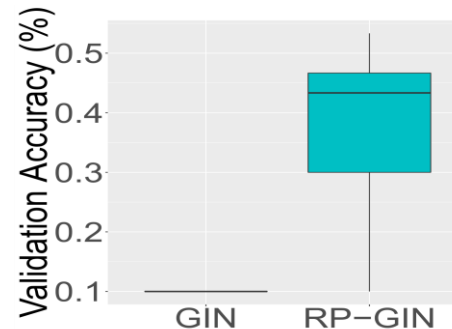
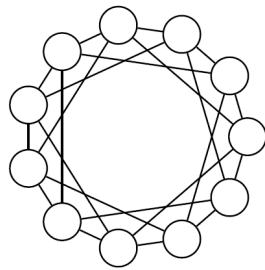
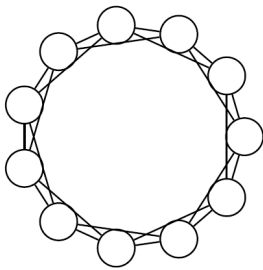
# Mixing GNNs with Relational Pooling (Murphy et al. ICML 2019)

- ▶ Consider a GNN  $\vec{f}$ , e.g., GIN of (Xu et al., ICLR 2019)
  - By definition  $\vec{f}$  is insensitive to permutations
- ▶ Let's make  $\vec{f}$  sensitive to permutations by adding node id (label) as unique node feature
  - And use RP to make the entire representation insensitive to permutations (learnt approximately via  $\pi$ -SGD)

## Task: Classify circulant graphs

$M = 11, L = 2$

$M = 11, L = 3$



## Task: Molecular classification

Model	<i>HIV</i>	<i>MUV</i>	<i>Tox21</i>
Molecule GCN	81.2 (1.4)	79.8 (2.5)	79.4 (1.0)
RP Mol GCN	<b>83.2</b> (1.3)	79.4 (2.5)	79.9 (0.6)

# RP gives New Class of Graph Representations

- ▶  $\vec{f}$  can be a logistic model (logistic regression)
- ▶  $\vec{f}$  can be a Recurrent Neural Network (RNN)
- ▶  $\vec{f}$  can be a Convolutional Neural Network (CNN)
  - Treat A as image

$$A = \begin{array}{c} \begin{array}{cccccccc} \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \end{array} \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \end{array}$$


- ▶ These are all valid graph representations in Relational Pooling (RP)

# Take-home

Relational Pooling (**RP**) framework gives a new class of graph representations and models

- Until now  $\vec{f}$  has been hand-designing to be permutation-invariant
- (Murphy et al ICML 2019) **RP**  $\vec{f}$  can be permutation-sensitive, allows more expressive models
  - Trade-off: can only be learnt approximately

$$\text{RP: } \bar{\bar{f}}(\mathbf{A}) \propto \sum_{\pi \in \Pi} \vec{f}(\mathbf{A}_{\pi\pi}, \mathbf{X}_{\pi}^{(v)})$$

Thank You!  
 @brunofmr  
 ribeiro@cs.purdue.edu

“Any”  $\vec{f}$  made permutation-invariant by **RP**

learns  $\bar{\bar{f}}(\mathbf{A})$  approximately via  $\pi$ -SGD

$\bar{\bar{f}}(\mathbf{A})$  inference with Monte Carlo  
 = approximately permutation-invariant

Permutation-invariant  
 $\vec{f}$  without permutation  
 averaging

learns exact  $\vec{f}$

$\vec{f}$  is always perm-invariant

# References

1. Murphy, R. L., Srinivasan, B., Rao, V., and Ribeiro, B., Janossy pooling: Learning deep permutationinvariant functions for variable-size inputs. ICLR 2019
2. Murphy, R.L., Srinivasan, B., Rao, V., Ribeiro, B., Relational Pooling for Graph Representations, ICML 2019
3. Meng, C., Yang, J., Ribeiro, B., Neville, J., HATS: A Hierarchical Sequence-Attention Framework for Inductive Set-of-Sets Embeddings. KDD 2019
4. de Finetti, B.. Fuzione caratteristica di un fenomeno aleatorio. Mem. R. Acc. Lincei, 1930
5. Aldous, D. J. Representations for partially exchangeable arrays of random variables. J. Multivar. Anal., 1981
6. Diaconis, P. and Janson, S. Graph limits and exchangeable random graphs. Rend. di Mat. e delle sue Appl. Ser. VII, 28:33–61, 2008
7. Kallenberg, O. (2005). Probabilistic Symmetries and Invariance Principles. Springer.
8. Kallenberg, O. (2017). Random Measures, Theory and Applications. Springer International Publishing.
9. Diaconis P. Finite forms of de Finetti's theorem on exchangeability. Synthese. 1977
10. Orbanz, P. and Roy, D. M. Bayesian models of graphs, arrays and other exchangeable random structures. IEEE transactions on pattern analysis and machine intelligence, 37(2):437–461, 2015.
11. Bloem-Reddy B, Teh YW. Probabilistic symmetry and invariant neural networks, arXiv:1901.06082. 2019.
12. Robbins, H. and Monro, S. A stochastic approximation method. The annals of mathematical statistics, pp. 400–407, 1951.
13. Hornik, K., Stinchcombe, M., and White, H. Multilayer feedforward networks are universal approximators. Neural networks, 2(5):359–366, 1989.
14. Cai, J.-Y., Furer, M., and Immerman, N. An optimal lower bound on the number of variables for graph identification. Combinatorica, 12(4):389–410, 1992
15. Shervashidze, N., Schweitzer, P., Leeuwen, E. J. V., Mehlhorn, K., & Borgwardt, K. M., Weisfeiler-lehman graph kernels. *Journal of Machine Learning Research*, 2011
16. Maron, H., Ben-Hamu, H., Shamir, N., and Lipman, Y. Invariant and equivariant graph networks. arXiv preprint arXiv:1812.09902, 201
- Maron, H., Fetaya, E., Segol, N., and Lipman, Y. On the universality of invariant networks. arXiv preprint arXiv:1901.09342, 2019.