

Ch. CS224W: Graph ML



§01. Introduction

- ① Traditional Methods:
Graphlets, Graph Kernels
- ② Node Embedding:
DeepWalk, Node2Vec
- ③ Graph Neural Networks:
GNN, GraphSAGE, GAT, GIN
- ④ graphs & reasoning:
TransE, RotaE
- ⑤ Deep Generative models:
GraphIMN

▲ 图形器学习, 图数据挖掘, 图神经网络 (GNN), 图卷积神经网络 (GCN), 深度图谱, 知识计算, 对话问答, 推荐系统, 社交计算, 媒体传播, 网络分析, 搜索引擎。

< Why Graphs? >

- ① 表达关系和表示语义。
Graphs are a general language for describing and analyzing entities with interactions.
- 实体图 (entity graph) — 非线性
Graph = (V, E) (实体图 + 图像 (img))
- event graphs / computer networks / disease pathways / particle networks / knowledge graphs / regulatory networks / scene graphs / 3D shapes ...
Königsberger Brückenproblem 七桥问题 (欧拉 - 4点)

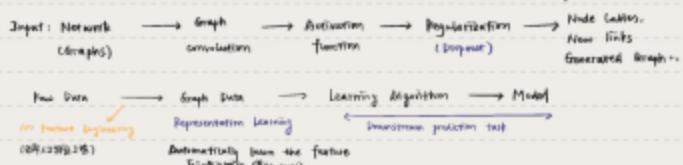
< Why not just use linear algebra? >

↳ relational graph.

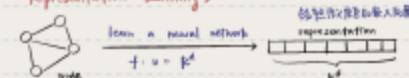
- modern deep learning toolbox is designed for simple sequences & grids.

arbitrary size and complex topological structure. (任意尺寸) — no spatial (localizing)
No fixed node ordering or reference points (没有固定的节点序数和参考点)
Often dynamic and have multimodal features. (动态变化 + 多模态特征)

predictions.



< Representation Learning? >



↳ 给定一个输入向量。

- map nodes to d -dimensional embeddings such that similar nodes in the network are embedded together

- PyG (PyTorch Geometric): the ultimate library for Graph Neural Networks

Graph Gym: Platform for designing Graph Neural Networks,

↳ modularized GIN implementation, simple hyperparameter tuning ...

< Different types of Tasks? >

Graph-level prediction { Node level (单个节点)

Graph generation { Community (子图) level (社群)

{ Edge level (连接 - 单条边) level (连接)

- ① Node classification: 预测类，预测一个节点的属性。

- ② Link prediction: 预测预测；预测两个节点之间是否存在连接。

- ③ Graph classification: 图分类：识别不同的图。

④ Clustering 集羣: Detect if nodes form a community.

⑤ Graph generation: Drug discovery

Graph evolution: Physical simulation

§ 02. Graph Representation (Network)

Objects, nodes, vertices (N)

Interactions, links, edges (E)

System: network, graph $G = (N, E)$

How to build a graph? What are edges? 如何设计图的节点连接
what are nodes? 和对称性相关的问题

<图的种类>

undirected: 无向图 —— symmetrical, reciprocal

directed: 有向图 —— arcs

• Homogeneous Graph: 单质图. → node & edge 都有相同属性,

$G = (V, E, \rho, T)$ Nodes with node types: $\forall i \in V$

Edges with relation types: $(v_i, v_j) \in E$

Node type: $T(v)$

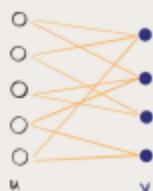
Relation type: $\rho(v, v')$

• Bipartite Graph: 二部图.

• A graph whose nodes can be divided into two disjoint sets U & V .

such that every link connects a node in U to one in V .

U and V are independent sets.



<Node Degrees> 顶点度数

* Node degree: k_i : the number of edges adjacent to node i .

* Avg. degree: $\bar{k} = \frac{1}{N} \sum_i k_i = \frac{2E}{N}$.

* in / out-degree: $\overline{k^+} = \frac{\sum_i k_i^+}{N}$



<Adjacency Matrix> 邻接矩阵.

$A_{ij} = 1$ If there is a link from node $i \rightarrow j$.

$A_{ij} = 0$ otherwise.

$$A = \begin{bmatrix} 0 & 1 & 0 & 1 \\ 1 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 \\ 1 & 1 & 1 & 0 \end{bmatrix} \xrightarrow{\text{node } 2 \rightarrow 4} \begin{bmatrix} 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 1 & 1 & 0 \end{bmatrix} \xrightarrow{\text{node } 4 \rightarrow 1} \begin{bmatrix} 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 1 & 1 & 0 \end{bmatrix} \xrightarrow{\text{node } 1 \rightarrow 2} \begin{bmatrix} 0 & 1 & 0 & 1 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 1 & 1 & 0 \end{bmatrix} \xrightarrow{\text{node } 2 \rightarrow 3} \begin{bmatrix} 0 & 0 & 1 & 1 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 1 & 1 & 0 \end{bmatrix}$$

$\xrightarrow{\text{node } 3 \rightarrow 2}$

$\xrightarrow{\text{node } 2 \rightarrow 3} \begin{bmatrix} 0 & 0 & 1 & 1 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 1 & 1 & 0 \end{bmatrix}$

\rightarrow not symmetric

* Networks are Sparse Graphs (稀疏图)

$$\hookrightarrow E \ll E_{\text{max}}, k = N-1.$$

Edge list: 各边连在连接的 node pair.

adjacency list: 各个 node 相关的边.

* unweighted: 无权重

$$\left\{ \begin{array}{l} E = \frac{1}{2} \sum A_{ij} \\ \bar{k} = \frac{2E}{N} \end{array} \right.$$



* self-edges (自环 - loop)

$$E = \frac{1}{2} \sum_{i,j} A_{ij} + \sum_{i,j} A_{ii}$$

* weighted: 有带权

$$\left\{ \begin{array}{l} E = \frac{1}{2} \sum \text{non-zero}(A_{ij}) \\ \bar{k} = \frac{2E}{N} \end{array} \right.$$

* multigraph

$$E = \frac{1}{2} \sum \text{non-zero}(A_{ij}) \\ \bar{k} = \frac{2E}{N}$$



< Connectivity > 连通性.

任意两个结点之间都有路径.

* connected graph: Any two vertices can be joined by a path.

A disconnected graph is made up by two or more connected components.

* The adjacency matrix of a network with several components can be written in a block-diagonal form, so that non-zero elements are confined to squares, with all other elements being zero.

Strongly, 强连通
弱连通.

Strongly: has a path from each node to every other node and vice versa. (A=B, B=A)
weakly: is connected if we disregard the edge directions.

* Strongly connected components (SCCs): 强连通分量.

↳ can be identified, but not every node is part of a non-trivial strongly connected component.

§ 0.3. Traditional Methods for ML in Graphs

{ Design features for nodes/links/graphs

{ Extract features for all training data

{ Build ML models based on these features

{ 基于特征的模型

* Traditional ML pipeline uses hand-designed features on undirected graphs.

↳ Goal: Make predictions for a set of objects

{ Features: 2-dimensional vectors

{ Objects: Nodes, edges, sets of nodes, entire graphs.

{ Objective function: What task are we aiming to solve?

< Node-Level > 节点级别, $f: V \longrightarrow k$

* Characterize the structure and position of a node in the network.

→ ① Node degree: 节点度数 → treats all neighboring nodes equally.

→ ② Degree Centrality: 度中心性. [Adjacency Matrix] \times [Column Vector] = [Degree Centrality].
(Eigenvector / Betweenness / Closeness ...)

without capturing important

node importance

< Node Centrality > 重要度 / 权重 → node importance in a graph

- Eigenvector centrality: 特征向量。

A node v is important if surrounded by important neighboring nodes. $v \in N(u)$.

- sum of the centrality of neighboring nodes:

$$\text{recursive defn} \quad C_v = \frac{1}{\lambda} \sum_{u \in N(v)} C_u \quad \text{归一化常数, } \lambda \text{ largest eigenvalue of } A.$$

$$\lambda C = AC$$

$\left\{ \begin{array}{l} A: \text{adjacency} \\ C: \text{centrality vector} \rightarrow \text{eigenvector of } A. \\ \lambda: \text{eigenvalue} \end{array} \right.$

By Perron - Frobenius Theorem: largest eigenvalue λ_{\max} is always positive and unique.

- Betweenness centrality

A node is important if it lies on so many shortest paths between other nodes.

$$C_B = \frac{1}{|V|} \sum_{\text{shortest path from } u \neq v \text{ through } v} \frac{1}{|\text{shortest path}|} \quad \text{路径数} \times \text{路径数} \times \text{路径数}$$

- Closeness centrality (距离度)

A node is important if it has small shortest path lengths to all other nodes.

$$C_C = \frac{1}{2|V|} \sum_{u \neq v} \frac{1}{\text{shortest path length from } u \neq v} \quad \text{路径数} \times \text{路径数} \times \text{路径数}$$

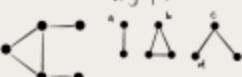
< Clustering Coefficient > 集合系数

Counts the triangles
in the ego-network
(周围子图)

↳ measures how connected v's neighboring nodes are.

$$C_L = \frac{\text{# edges among neighboring nodes}}{\binom{|N(v)|}{2}} \quad \text{路径数} \times \text{路径数} \rightarrow \text{路径数}$$

counts pre-specified
subgraphs



• Importance-based

node degree → number of nodes
node centrality measure → importance.

Structure-based:

node degree
Wasserman coefficient → measured
regularity metric is measures

< Graphlet >

- Graphlet Degree Vector (GDV): A count vector of graphlets count_i at a given node.

↳ 路径子图统计

small subgraphs that describe the structure of node v's network neighborhood.

provides a measure of a node's local network topology. 本地网络拓扑 (局部网络拓扑).

- Induced Subgraph (ISG)

$\cong \text{ISG}_v$

由图中所有与该节点相关的边和顶点组成的子图，称为邻接子图。

formed from a subset of vertices and all of the edges connecting the vertices in that subset.

isomorphic (同构) \Leftrightarrow 两个子图具有相同的结构。

↳ 同构子图的集合。

Two graphs which contain the same number of nodes connected in the same way.

Part 2: Link Prediction Task and Features

基础特征三类
历史 (H)

预测未来 (F)
综合 (HF).

* The task is to predict new links based on the existing links. 通过已知的特征和历史

node pairs (with no existing links) are ranked, and top K node pairs are predicted.

直接预测未来链接，即历史预测 ✓

记下所有历史的预测链接数 × → 52% 的精确度是相当令人惊讶的。

① Links missing or random:

↳ Remove a random set of links and then aim to predict them.

② Links over time:

Given $G(t_0, t_1)$ a graph defined by edges up to time t_0 , output a ranked list L of edges,

(not in $G(t_0, t_1)$) that are predicted to appear in $G(t_0, t_1)$.

* Evaluation: $|L| = |E_{test}|$: # new edges that appear during the test period $[t_0, t_1]$

Take top n elements of L and count correct edges. 比单纯随机抽样，它需要更多的计算工作。

* Methodology: 基础特征模型 → 特征集 → G(t_0, t_1)

① For each pair of nodes (x, y) computes score $C(x, y)$

② Sort pairs (x, y) by the decreasing score $C(x, y)$

③ Predict top n pairs as new links. MAP@B.

④ See which of these links actually appear in $G(t_0, t_1)$.

* shortest-path distance between two nodes. 最短路径距离，邻居数。

↳ This does not capture the degree of neighborhood overlap.

* Captures # neighboring nodes shared between 2 nodes $x \neq y$:

Common neighbors: $\frac{|N(x) \cap N(y)|}{|N(x)|}$ → BRUTSFAR

Second's neighbors: $\frac{|N(x) \cap N(y)|}{|N(x) \cup N(y)|}$ → BRUN

Balance-factor: $\frac{1}{2(|N(x) \cap N(y)| / |N(x)|)}$ → BRUNBALANCE

* Limitations: Metric is always 0 if the two nodes do not have any neighbors in common.

$N(x) \cap N(y) = \emptyset \rightarrow |N(x) \cap N(y)| = 0 \rightarrow$ 基本上忽略边缘。

* Note today: count the number of walks of all lengths between a given pair of nodes.

↳ $\sum_k N(x) \times N(y) \times k! \rightarrow$ 太慢了！

power of Adjacency matrix: $A_{xy} = 1$ if $y \in N(x)$, 否则为0，否则为0。

① Let $P_x^{(k)} =$ # walks of length k between x and y .

② $P_x^{(k)} = A^k$

③ $P_x^{(1)} =$ # walks of length 1 (direct neighborhood) between x and $y = A_{xy}$.

④ $P_x^{(0)}$, sum up these #walks among x 's neighbors.

$$P_x^{(0)} = \sum_k P_x^{(k)} = P_x^{(1)} + \sum_k P_x^{(k)} = A_{xx}$$

邻居数 = 走路数之和。

< Edge-based Features >

Distance-based feature

Local neighborhood overlap (LOLO)

Global neighborhood overlap (GNO)

$$\text{S}_{\text{WLO}} = \sum_{k=1}^{\infty} p^k A_{\text{WLO}}^k$$

weight: discount factor

$$S = \sum_{k=1}^{\infty} p^k N^k = (1 - p\alpha)^{-1} - I$$

↓ by geometric series of matrices.

几何级数。

PART 3. Graph-Level &

Graph kernels

全局层面的特征工程

② 特例性问题判断 2/2

是否同构 (验证 + 构造)

k=3:



Limitation: expensive!

时间复杂度: n^k (多项式级数)

worst-case: 2图同构测试。

subgraph isomorphism test
(NP-hard)

Computationally efficient.

Linear to #edges.

Characterize the structure of an entire graph. 提取出能很好地反映全局的特征。

<Kernel Method> 将非线性空间的非线性映射问题, 转化为高维空间的线性映射问题。

↳ design kernels instead of feature vectors.

* Graph kernels, $\{ \text{Graphlet kernel} \}$ $\rightarrow \text{Bag-of-...}$

measure similarity between 2 graphs. Weisfeiler-Leman kernel.

* Goal: Design graph feature vector $\mathbf{p}(G)$ 难度大。Key Idea: Bag-of-Words (BoW) \rightarrow word counts (将图中所有词频统计出来)

Naïve extension to graph: regard nodes as words. 语义相似度, 热带相似度。

↳ 先看每一个节点的，再看整体的。

Bag of Node Degrees \rightarrow 只算Node degree 1%。 像节点，不看连通结构。

<Graphlet Features>

Let $G = (g_1, g_2, \dots, g_m)$ be a list of graphlets of size k .

(k=3, 4 graphlets k=4, 11 graphlets)

Count the number of different graphlets in a graph.

+ (1) 统计 Graphlet 1 的频率, 然后统计 Graphlet 2 的频率。

Given 2 graphs G & G' : graphlet kernel is computed as:

$$K(G, G') = f_G^T f_{G'}^T \quad (\text{即 } G \text{ 和 } G' \text{ Graphlet Count Vector 相乘})$$

B: If G and G' have different sizes. \rightarrow Show the value.

S: normalize each feature vector. (归一化)

$$h_G = \frac{f_G}{\sum_i f_G} \quad K(G, G') = h_G^T \cdot h_{G'}^T$$

<Weisfeiler-Leman Kernel> — Colour Refinement

↳ Use neighbourhood structure to iteratively enrich node vocabulary.

Given: A graph G with a set of nodes V .① Assign an initial colour $C^{(0)}(v)$ to each node v .② Iteratively refine node colours. Ig: \rightarrow 逐层，用图结构。

$$C^{(t+1)}(v) = \text{HOST}\left(\{C^{(t)}(v), \{C^{(t)}(u)\}_{u \in N(v)}\}\right).$$

↳ Maps different inputs to different colors.

③ After k steps of colour refinement. $C^{(k)}(v)$ summarizes the structure of k -hop neighbourhood.