

第 1 章 图机器学习概论

引言

现实世界中充满了各种各样的系统，各个系统由大量类型各异、彼此交互的组件构成，例如生态系统、社交网络和计算机网络等，在这些系统中相互作用的组件可以抽象为图结构。得益于图结构数据对复杂交互关系的强大建模能力，在当今数据密集的研究领域和商业应用中，图论和图机器学习技术扮演着至关重要的角色。本章旨在揭示图的基础知识及机器学习在图中的应用，特别是如何利用图结构对复杂关系和多维数据进行建模和分析。本章首先介绍图的基本概念，包括定义、表示和类型，为理解图数据结构奠定基础；随后探讨图机器学习的核心概念和应用领域，通过具体任务展示其实际用途和挑战；最后介绍图机器学习的发展历程，从早期的理论研究到现代的图神经网络，展示了图机器学习的演变过程并对其未来进行了展望。

本章学习目标

- (1) 理解图的基本概念、定义及其重要属性，以及这些属性如何影响图的表示和处理；
- (2) 掌握不同的图表示方法，并理解各自的使用场景与优缺点；
- (3) 区分图的不同类型，并了解它们在实际应用中的差异；
- (4) 理解图机器学习的基本概念和技术，包括其原理及其在多种图分析任务中的应用；
- (5) 探索图机器学习的历史发展，以及当前的研究热点和未来可能的发展方向。

1.1 图基础知识

在现实生活中，图数据无处不在。作为强大的工具，它们在各个领域中用于表示复杂的关系。在化学和材料科学中，图用于表示分子的结构关系，其中原子被视为节点，化学键作为边，形成化合物的图表示^[1]。在社会科学中，图可以用来表示个体之间的交互关系，用节点代表人，边就可以表示人际关系^[2]，例如友谊、同学关系或师生关系。此外，推荐系统中的图同样重要，可以用节点表示用户或商品，边表示购买或点击行为，从而揭示用户的消费模式^[3]。因为图数据具有强大的表达能力，利用机器学习来分析图受到越来越多的关注。

1.1.1 图的定义和表示

1. 图的定义

图可以表示为 $G = (V, E)$ ，其中 $V = \{v_1, \dots, v_{|V|}\}$ 表示节点集合， $E = \{e_1, \dots, e_{|E|}\}$ 表示边集合。其中： $|V|$ 是图中节点的数量， $|E|$ 是图中边的数量。边也可以用其两端的节点进行表示，例如连接节点 v_1 和 v_2 的边也可以表示为 (v_1, v_2) 。在许多情况下，人们只关心简单无向图，如图 1-1 所示，即每对节点之间最多只有一条边，没有节点与自身相连的边，并且这些边都是无向的，即 $(u, v) \in E \Leftrightarrow (v, u) \in E$ 。

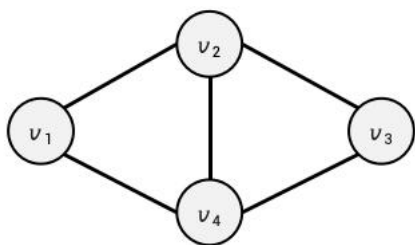


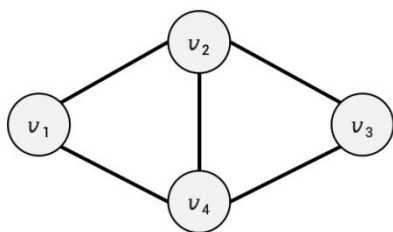
图 1-1 简单无向图

2. 图的表示

1) 邻接矩阵表示法

如图 1-2 所示，给定图 $G = (V, E)$ ，可以使用邻接矩阵 $A \in \{0,1\}^{|V| \times |V|}$ 表示边的分布。邻接矩阵的第 (i, j) 项表示为 A_{ij} ，其表示节点 v_i 与 v_j 之间的连接性。如果 v_i 与 v_j 之间存在一条边，则 $A_{ij} = 1$ ，否则 $A_{ij} = 0$ 。

特别地，在有向图中，边是从一个节点指向另一个节点的，而在无向图中，同一条边的两个节点的指向顺序没有区别。在无向图中，若节点 v_i 与节点 v_j 相邻，则 v_j 与 v_i 相邻，因此对所有图中的 v_i 和 v_j 有 $A_{ij} = A_{ji}$ ，对应的邻接矩阵是对称的。注意，除非特别说明，讨论将会限制在无向图中。使用邻接矩阵，可以轻松地计算出一个节点的度，即该节点与其他节点的连接数。



(1) 图

节点	v_1	v_2	v_3	v_4
v_1	0	1	1	0
v_2	1	0	1	1
v_3	0	1	1	0
v_4	1	1	1	0

(2) 邻接矩阵

图 1-2 图和邻接矩阵

邻接矩阵是一种直观、简单且易于理解的图表示方法。其可以方便地检查任意一对节点之间是否存在边，便于找出任一节点的所有邻接点，以及便于计算任一节点的度。然而，邻接矩阵也存在一些缺点。首先，增加或删除节点时需要调整矩阵的行列，操作不便；其次，对于稀疏图，邻接矩阵会浪费大量空间，会有许多无效元素。

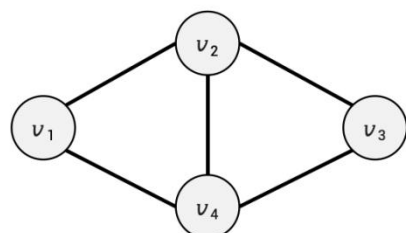
邻居图： G 中节点 v_i 的邻居集合表示为 $N(v_i)$ ，它包含所有与 v_i 直接相连的节点。

度：节点 v_i 的度数可以表示为 $d_i = \sum_{j=1}^N A_{ij}$ 。节点 v_i 的度数等于 $N(v_i)$ 中节点的数量，即 $d_i = |N(v_i)|$ 。对角度数矩阵亦称作对角度矩阵，可以表示为 $D = \text{diag}(d_1, d_2, \dots, d_n)$ 。

2) 邻接表表示法

如图 1-3 所示，给定图 $G = (V, E)$ ，其邻接表是一个包含 $|V|$ 个列表的数组，每个列表对应于图 G 的一个节点。对于每个节点 $v \in V$ ，图 G 的邻接表中 v 的列表包含所有与 v 相邻的节点。

如果图 G 是无向的，那么对于每条边 $(v, w) \in E$ ，节点 w 将出现在节点 v 的邻接表中，同时节点 v 也会出现在节点 w 的邻接表中。如果图 G 是有向的，那么对于每条边 $(v, w) \in E$ ，节点 w 将出现在节点 v 的邻接表中，但节点 v 不一定会出现在节点 w 的邻接表中，除非存在一条反向的边 (w, v) 。



(1) 图

节点	邻接点
v_1	v_2, v_4
v_2	v_1, v_3, v_4
v_3	v_2, v_4
v_4	v_1, v_2, v_3

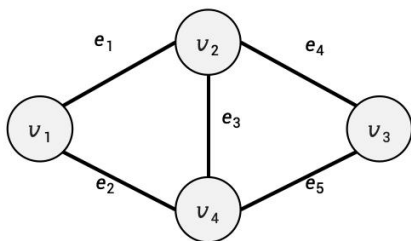
(2) 邻接表

图 1-3 图和邻接表

邻接表特别适用于存储稀疏图，这种图中的边相较于节点组合数量要少得多。相比于邻接矩阵，邻接表在存储空间上更为高效，并且能够更快速地找到与某个节点相连的所有邻接节点，邻接表在遍历邻接节点时效率较高。然而，邻接表的缺点是实现相对复杂，且在需要频繁查询两个节点之间是否存在边时效率较低。

3) 关联矩阵表示法

如图 1-4 所示，关联矩阵是一个维度为 $|V| \times |E|$ 的矩形矩阵。矩阵中的每个元素 a_{ij} 表示第 i 个节点和第 j 条边之间的连接关系：如果第 i 个节点与第 j 条边关联（即该边连接到该节点），则 $a_{ij} = 1$ 。如果第 i 个节点与第 j 条边不关联，则 $a_{ij} = 0$ 。



(1) 图

节点 \ 边	e_1	e_2	e_3	e_4	e_5
v_1	1	1	0	0	0
v_2	1	0	1	1	0
v_3	0	0	0	1	1
v_4	0	1	1	0	1

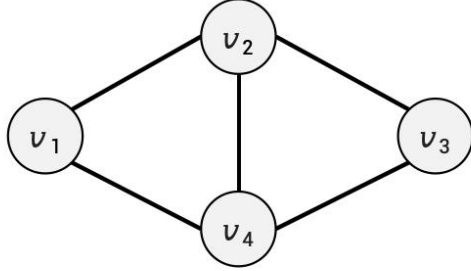
(2) 关联矩阵

图 1-4 图和关联矩阵

关联矩阵适用于稀疏图，可以清晰地表示每个节点与边的关系，适合于处理边的属性和权重，还可以方便地表示节点和边之间的关系，适用于求解边的度数等问题。但对于边数较多的图，关联矩阵的列数会显著增加，导致内存消耗大。

4) 拉普拉斯矩阵表示法

如图 1-5 所示，对于一个以 A 为邻接矩阵的图 G ，其拉普拉斯矩阵定义为 $L = D - A$ ，其中 $D = \text{diag}(d(v_1), \dots, d(v_N))$ 是对角度矩阵。



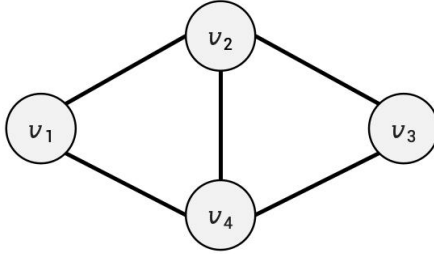
(1) 图

$$\begin{bmatrix} 2 & -1 & -1 & 0 \\ -1 & 3 & -1 & -1 \\ -1 & -1 & 3 & -1 \\ 0 & -1 & -1 & 2 \end{bmatrix}$$

(2) 拉普拉斯矩阵

图 1-5 图和拉普拉斯矩阵

标准化拉普拉斯矩阵：如图 1-6 所示，对于一个给定的以 A 为邻接矩阵的图 G ，其标准化拉普拉斯矩阵记作 \tilde{L} ，定义为 $\tilde{L} = D^{-\frac{1}{2}}(D - A)D^{-\frac{1}{2}} = I - D^{-\frac{1}{2}}AD^{-\frac{1}{2}}$ 。



(1) 图

$$\begin{bmatrix} 1 & -\frac{1}{\sqrt{6}} & -\frac{1}{\sqrt{6}} & 0 \\ -\frac{1}{\sqrt{6}} & 1 & -\frac{1}{3} & -\frac{1}{2\sqrt{3}} \\ -\frac{1}{\sqrt{6}} & -\frac{1}{3} & 1 & -\frac{1}{2\sqrt{3}} \\ 0 & -\frac{1}{2\sqrt{3}} & -\frac{1}{2\sqrt{3}} & 1 \end{bmatrix}$$

(2) 标准化拉普拉斯矩阵

图 1-6 图和标准化拉普拉斯矩阵

谱图理论是图论与线性代数相结合的产物，它通过分析图的某些矩阵的特征值与特征向量来研究图的性质。拉普拉斯矩阵是谱图理论中的核心概念，是研究图的结构特性和节点间关系的重要工具。它可以通过考虑节点的度和邻接性来捕获图的几何和拓扑特性。拉普拉斯的特征值和特征向量在谱图论中有着广泛的应用。标准化的拉普拉斯矩阵则是为了消除不同节点度对分析结果的影响。这保持了拉普拉斯矩阵的对称性和半正定性，因此在运行谱聚类算法时更有用，有助于图信号处理和机器学习工作更好地利用图结构信息。拉普拉斯矩阵和标准化拉普拉斯矩阵在机器学习与深度学习中有着重重要的应用，例如：流形学习数据降维算法中的拉普拉斯特征映射^[4]、局部保持投影^[5]，无监督学习中的谱聚类算法^[6]，半监督学习中基于图的算法^[7]，以及目前炙手可热的图神经网络^[8]等。

1.1.2 图的类型

在现实生活中，复杂系统无处不在。然而，这些系统中庞大而多样的节点和关系使得使用简单的图模型无法完全捕捉现实世界的复杂性。以社交网络为例，用户与用户之间的关系可能是单向的。在这种情况下，可以使用有向图来清楚地描绘这种单向的关系。在交通网络中，不仅有不同的道路或路线连接在一起，而且还与距离、时间或成本等其他因素相关。带权图引入了边权重，可以更好地捕捉这些信息。在电商网络中，节点的类型以及它们之间的交互关系往往是多样的。例如节点可能表示商品，用户、商家、品类等，边可能表示它们之间的购买、浏览、评价等关系。异质图在捕捉这种多样性方面更有效。同时，许多网络中的节点和边可以存储一些相关的属性信息。例如，在社交网络中，每个节点（用户）具有年龄、性别、职业等属性，而边则可能表示如朋友关系或关注关系等，属性图可以更好地分析这些属性，有助于识别特定群体的社交行为模式。除此之外，许多现实世界的网络都涉及组件之间不断发展的交互。动态图可以反映某个特定快照时间的网络结构，同时捕捉该结构随时间的变化。例如社交网络中的联系人圈可能会随着时间的推移而扩大或缩小，这种情况下，动态图就能够反映出网络在某一时间点的状态，同时还能捕捉到网络结构随时间的变化。这些不同类型的图提供了强大的工具，使人们能够更深入地理解和分析现实世界中的复杂系统。

1. 无向 vs. 有向图

$G = (V, E)$ 由一个非空的节点集 V 和一组有向边集 E 组成。 E 中的每条边 e 是由一个有序顶点对 $(u, v) \in V$ 指定的。无向图被视为有向图的一种特殊情况，其中如果两个节点之间连接，则存在一对方向相反的边。当且仅当邻接矩阵是对称的时，图才是无向的。

无向图中的边没有方向，如果存在一条边 (u, v) ，则 u 和 v 是相邻的，可以从 u 到 v 或从 v 到 u 。每个节点有一个度数，等于与其相连的边的数量。如果存在一条从节点 u 到节点 v 的路径，则必定存在一条从节点 v 到节点 u 的路径。无向图常用于表示双向关系，如社交网络。有向图中的边则有方向。如果存在一条边 (u, v) ，则可以从 u 到 v ，但不一定能从 v 到 u 。每个节点有入度和出度两个度数，入度是指向该节点的边的数量，出度是从该节点出发的边的数量。有向图常用于表示单向关系，如网页之间的超链接关系、推荐系统。

如图 1-7 所示，在现实生活中，节点之间关系并不总是简单的双方面的互动。例如在社交网络中，用户关系可能是单一方向的，一个用户可能关注另一个用户，但反之则不一定发生。在这种情况下，用无向图来建模并不完全合适，在无向图中，所有节点之间的关系会被看成双向的，这是适用于需要对等关系的情况，如朋友关系。但现实世界中很多交互关系是有明确方向的，这体现出有向图的重要作用。

有向图则允许在节点之间引入边的方向性，从而更精确地捕捉和表达各个对象之间的不同关系。在有向图中，边的方向性给每种关系赋予更加丰富的含义。例如在一个社交网络中，边的方向可以表示关注这种关系。如果有一条边从节点 A 指向节点 B，那么它传达的信息是 A 关注了 B，同时也表明 B 并不一定也关注了 A。

在有向图中，边可以分为入边和出边，分别表示其他节点对某个节点的影响以及某个节点对其他节点的影响。这种信息的引入使得有向图在分析复杂系统时能够提供更多维度的洞察力。通过对有向图的深入分析，研究者能够识别出网络中的关键节点、信息传播路径等，从而为推荐系统、社交网络分析等领域提供有力的支持。

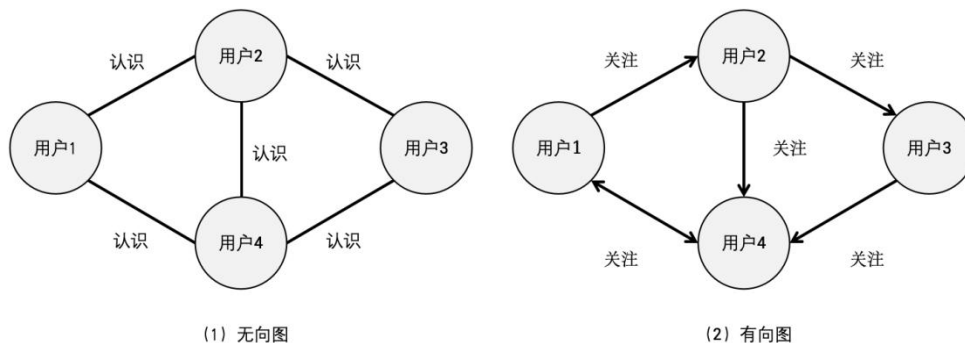


图 1-7 无向图和有向图的对比

2. 无权 vs. 带权图

一个带权图是一个图 $G = (V, E)$ ，以及一个权重函数 $w: E \rightarrow \mathbf{Z}$ 。即：为每条边 $e = (u, v) \in E$ 分配一个整数权重，一般用 $w(e) = w(u, v)$ 来表示。

无权图在边上没有权重，即边的权重统一为 1。这意味着所有边都具有同等的重要性或成本。带权图可以在边上赋予权重，这些权重可以代表成本、距离、时间或任何其他表示两个节点之间关系强度的数值度量。

如图 1-8 所示，在现实生活中，节点之间的关系并不总是具有相同的强度或重要性。例如，在交通网络中，不同道路之间的通行时间或距离可能各不相同。在这种情况下，使用无权图来建模可能并不完全合适。在无权图中，边没有权重，即每个边都是等价的。这使得无权图适合处理仅需要识别节点之间是否存在边而不必识别该边的权重或强度的情况。例如，在表示朋友之间关系的社交网络中，如果两个人是朋友，那么他们之间存在一条边，并且没有与这条边相关的其他特征。此时无权图便是一个良好的建模工具，无权图大大简化了图的复杂性，在处理大型图网络时更为高效。

带权图则允许在边上赋予权重，这些权重可以代表距离、成本、时间或其他度量标准，从而更精确地模拟现实世界中的各类关系。许多实际情况可以通过加权图更准确地建模，例如，道路网络中城市之间的交通流量或通信网络的带宽限制。通过对带权图的分析，可以更好地揭示出隐藏在网络中的重要模式和规律，从而为交通管理、物流优化、经济活动分析等领域提供强有力的支持。

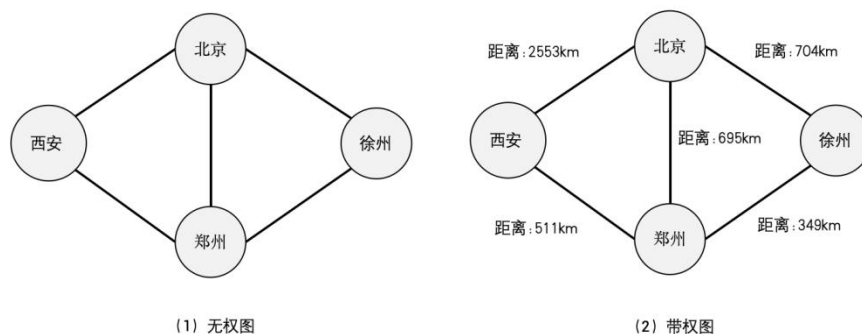


图 1-8 无权图和带权图的对比

3. 同质 vs. 异质图

一个异质图 G 由一组节点 $V = \{v_1, \dots, v_N\}$ 和一组边 $E = \{e_1, \dots, e_M\}$ 组成。每个节点 v 和边 e 都与其类型映射函数 $\phi_v: V \rightarrow T_v$ 和 $\phi_e: E \rightarrow T_e$ 相关联，其中 $|T_v| + |T_e| > 2$ 。其中 $|T_v|$ 表示图中节点的种类数， $|T_e|$ 表示图中边的种类数。

同质图中所有节点和边都属于同一类型，结构相对简单，适合使用统一的分析方法。相比之下，异质图包含多种不同类型的节点和边，能够更准确地捕捉复杂关系，适合用于建模多层次、多维度的系统。

如图 1-9 所示，在现实生活中，各种复杂系统通常由多种不同类型、相互作用的组件组成，例如生物系统、社交网络和计算机系统等。在这种情况下，使用同质图来建模可能并不完全合适。在同质图中，所有的节点都是同一类型，且节点之间的边也是相同的类型。同质图建模方法往往忽略了实际交互系统中对象及其关系的异质性，只捕捉了部分信息，导致不可逆的信息损失。

异质图则包含多种类型的节点和边，各类型的节点和边具有各自独特的属性和关系类型。例如，在一个电子商务平台上，用户、商品、商家等可以作为不同类型的节点，而用户购买商品、商品属于商家等则代表不同类型的边。这种多样性使得异质图能够更全面地表示现实世界中的复杂关系，实现对现实世界更完整自然的抽象。

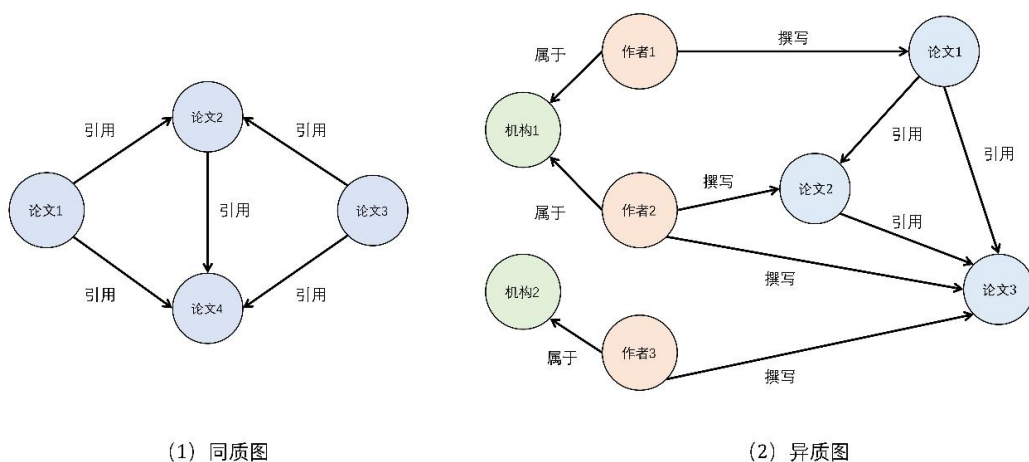


图 1-9 同质图和异质图的对比

4. 无属性 vs. 属性图

一个属性图 $G = (V, E, F)$ 由节点集 V 、边集 E 以及一组映射 $F = \{f_1, \dots, f_N\}$ 组成，使得对于 $i \in [1, \dots, N]$ ， $f_i: V \rightarrow \text{dom}(a_i)$ 将属性 a_i 的值 $f_i(v)$ 分配给节点 v ，其中 $\text{dom}(a_i)$ 是属性 a_i 的定义域。

无属性图（又称为纯图）是一种简单的图结构，仅关注节点之间的连接关系，不包含任何额外的属性信息，而属性图则在节点和边上附加了丰富的属性信息。这种结构不仅能表示连接关系，还能提供更多的语义信息，适用于更复杂的分析。

如图 1-10 所示，在现实生活中，对象之间的关系往往不仅仅是简单的连接，而是包含了丰富的属性信息。例如在蛋白质-蛋白质相互作用网络中，蛋白质可以作为节点，它们之间的相互作用可以作为边。在这个图中，节点需要带有如名称、类型、功能等属性，边需要带有如激活、抑制等相互作用类型的属性。在这种复杂的网络中，用无属性图来建模并不完全合适，无属性图中所有节点和边不包含任何额外的属性信息。无属性图可能无法充分表达节点和边的丰富属性信息，这便体现出属性图的重要。

属性图则可以在节点和边上附加丰富的属性信息，从而更精确地表达各个对象之间的不同关系。例如在电子商务平台中，用户、商品、商家等都可以作为节点，用户购买商品、商品属于商家等行为可以作为边。在这个属性图中，用户节点可以带有年龄、性别、购物历史等属性，商品节点可以带有价格、品牌、评价等属性，边可以带有购买时间、数量等属性。这种丰富的属性信息使属性图能更好地捕捉现实世界的信息，更准确地模拟现实世界中的复杂系统。

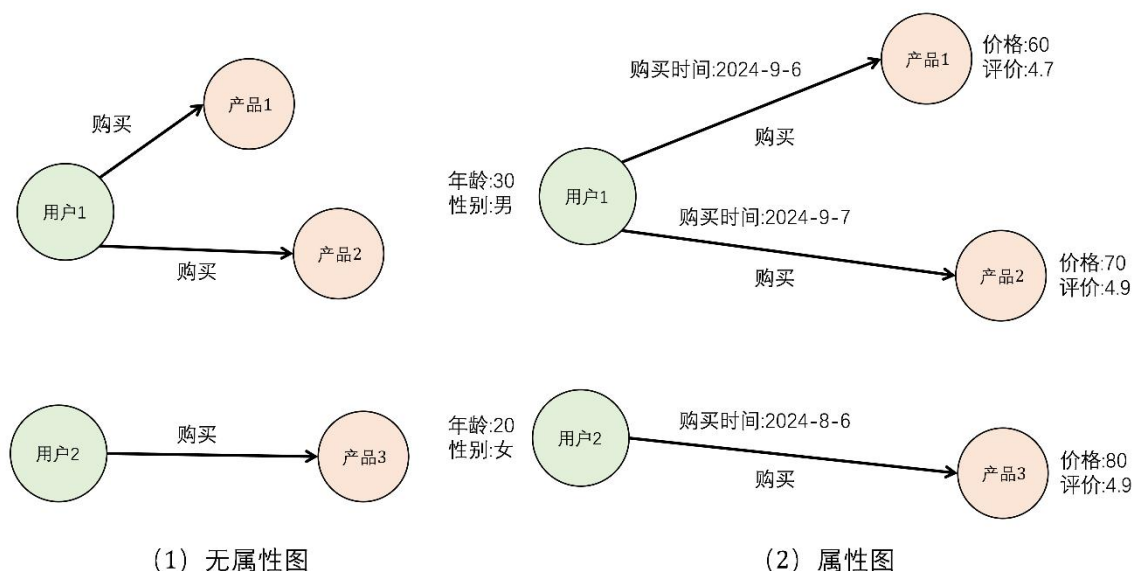


图 1-10 无属性图和属性图的对比

5. 静态 vs. 动态图

动态图 $G = (V, E, T)$ 由节点集 V 、边集 E 和时间映射 T 组成。具体来说，每个节点或每条边都与时间戳信息相关联，这些时间戳指示它们出现的时间。时间映射 $T: (V, E) \rightarrow (V', E')$ ，表示节点集 V 和边集 E 及其属性随时间的变化关系，其中 V' 是映射后的节点集， E' 是映射后的边集。

动态图是一种节点、边本身以及他们各自属性可以随着时间的变化进行增加、删除或修改的图，这意味着节点和边的关系随着时间的变化是可以变化的，可以更准确地反映现实世界中的复杂系统。而静态图在节点和边上都没有时间属性，即图的节点和边都是固定的，不会随时间变化。动态图适合用于描述一种随时间变化的关系，而静态图则适合用于描述一种稳定的关系。

如图 1-11 所示，在现实世界中，许多复杂系统是动态变化的，它们的节点和边以及各种属性都会随着时间的变化而变化。例如，在社交网络中，会有新的用户加入和现有的用户退出，用户之间的关系会建立和断裂，用户的属性例如爱好等也会随着时间的变化而变化。在这种情况下，静态图的建模可能并不完全合适，可能会忽视这些重要的动态变化，从而影响对系统的理解和预测的效果。

动态图的节点、边和属性则可以随着时间的变化进行增加、删除或修改。这使得研究者可以更好地理解系统的演化过程，预测系统的未来状态。例如，当一个新用户注册社交平台时，可以将其作为一个新的节点加到图中。当这个用户关注其他用户时，可以将关注关系表示为边添加到图中。当这个用户发布新的推文，或者对其他用户的推文进行回复时，可以将这些行为表示为节点的属性，并更新这些属性。随着时间的推移，社交网络的具体结构和其中的属性会不断发生着变化。基于静态图去解决问题，虽然处理起来更简单，但不符合实际，预测效果也会大打折扣。当使用动态图来建模时，可以更好地理解和预测网络的演化过程，例如可以预测哪些用户可能会获得大家的关注，哪些话题可能会成为热门话题等。

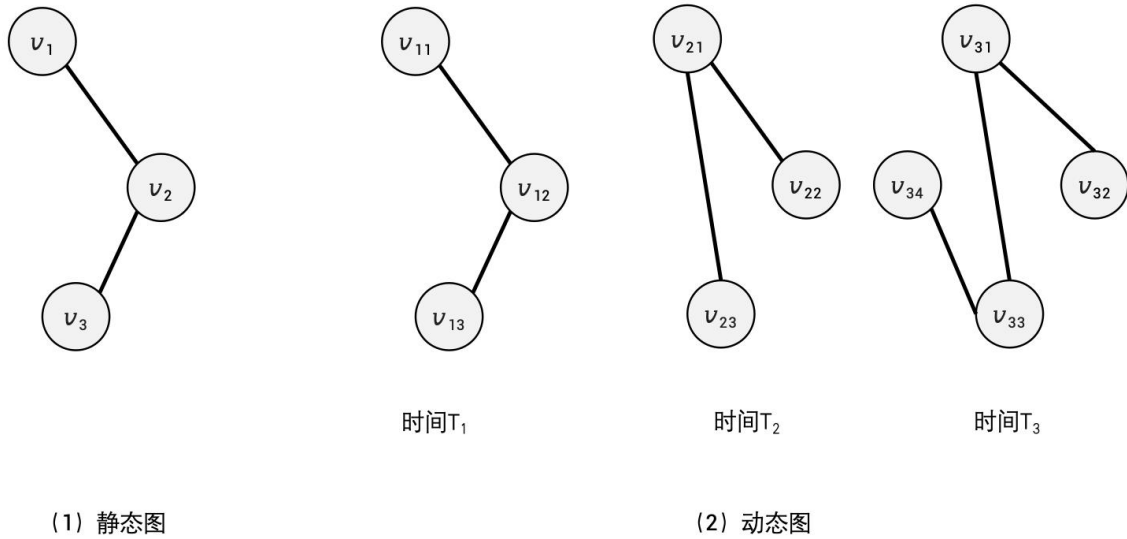


图 1-11 静态图和动态图的对比

1.2 图机器学习

1.2.1 基本概念

近年来，机器学习在各类任务建模方面取得了显著成功。与传统计算机程序中的手工编写规则的计算方式不同，机器学习是一种数据驱动的计算方法，它通过让计算机从大量数据中自动学习模式和规律来解决实际问题^[9]。例如，在计算机视觉（Computer Vision, CV）领域，基于卷积神经网络（Convolutional Neural Networks, CNN）的模型能够从大量图像中学习并提取特征，从而完成图像分类、目标检测和图像分割等任务^[10]。在自然语言处理（Natural Language Processing, NLP）领域，基于循环神经网络（Recurrent Neural Networks, RNN）的模型能够从大量语料中学习、理解和生成文本，实现语言翻译和对话问答等功能^[11]。

机器学习是一门跨学科的领域，结合了计算机科学和概率统计学的基本原理。其核心思想在于利用计算机技术对现有数据进行分析，从而构建概率统计模型，以实现未知数据的预测和分析。这一过程通常涉及三个主要研究对象：数据、模型和算法。具体而言，机器学习的基本假设是观测到的数据呈现一定的统计规律，这些规律可以推广至未知的同类数据。因此，收集并构建适合特定任务的数据集至关重要，因为它们直接影响后续过程中模型的预测表现；不同的任务和数据具有各自特点，因此需要根据数据的性质和具体任务的需求，对模型的函数空间进行合理的假设；为了从所选择的函数空间中获得任务最优模型，还需要设计有效的算法以支持模型的求解。机器学习算法种类繁多，为了便于研究，通常会根据任务的性质和数据标注的情况进行分类：首先，基于数据样本的任务标注情况，算法可分为监督学习、半监督学习和无监督学习。监督学习中所有训练数据都得到了标注；半监督学习中只有部分训练数据得到了标注；无监督学习则需要在无标注数据上进行学习。其次，根据任务标签是否连续，机器学习算法可以分为分类算法和回归算法。其中，分类算法用于将输入数据分配至预定义类别，适用于离散标签的任务；而回归算法则用于预测连续数值输出，通常用于解决连续值预测问题。

尽管机器学习在处理图像和文本等数据方面取得了惊人的效果，其有效性在很大程度上依赖于针对数据特性所设计的归纳偏置（Inductive Bias）^[12]。例如，基于图像数据的平移等变性和局部相似性假设，设计出了基于局部卷积的模型；基于文字数据的序列依赖假设，设计出了基于循环递归的模型。然而传统机器学习模型往往缺乏针对图结构数据的特殊设计，因此，图机器学习的概念应运而生。图机器学习是将机器学习方法应用于图结构数据的技术。其从应用

场景中抽象构建出图结构数据，通过针对性的模型设计来处理图中节点、边及其之间的复杂交互关系，从而捕捉图结构数据中的潜在模式来进行预测。

在数据、模型和算法三个层面上，图机器学习和传统的机器学习相比有所区别。在数据层面，图机器学习主要关注于拓扑结构数据中的对象，例如图中的节点、节点间的连边或整个图。与文本（序列结构）和图像（网格结构）等规则数据结构相比，图数据的结构显得更加不规则，例如图中不同节点的邻居数量往往也不同。这种不规则的结构特性使得图成为一种更通用的数据结构，即可以将文本和图像视为其特例：例如可以通过将像素点建模为节点，将像素的颜色和位置作为节点属性，并在相邻像素之间建立连边，从而将图像表示成图；将词建模为节点，并基于词的先后顺序或语法关系等建立连边，从而将文本序列建模成图^[13]。在模型层面，不同于传统机器学习中的样本独立性假设，图机器学习研究的对象往往具有关联关系，例如图中的节点通过连边与其他节点建立起不同的依赖关系，使得图机器学习需要对模型的函数空间引入新的合理假设。在算法层面，由于图上的标注更加困难，因此图机器学习算法更关注无监督、半监督特别是小样本的学习，例如节点分类任务往往要求模型在仅有极少标注节点的图上进行训练，并对剩下大部分节点给出预测结果，社区发现任务则要求模型能在无标注的图上识别出其中的社区结构。

作为机器学习中的一个新兴方向，图机器学习面临着诸多挑战，可以将这些挑战主要归纳于模型设计和系统优化两个层面：

1. 模型设计

传统的机器学习模型往往采用独立性假设，即假设样本点之间是相互独立的，然而图数据中的样本（节点）则会与相邻样本产生交互，从而产生关联性假设，要求模型能够充分挖掘利用节点之间的交互信息。而不规则的结构使得这种交互难以用统一的方式进行建模，因此如何针对不同类型图数据进行模型设计，如何有效地在节点间进行信息交互充满挑战。

2. 系统优化

真实世界中的图数据通常具有超大的规模和复杂的连接关系，这对图的存储提出了高要求。同时由于结构的不规则性，传统的并行计算和分布式算法的设计与实现难以直接应用到图机器学习中，因此如何设计有效的分布式采样算法，如何高效地进行并行化计算是充满挑战的问题。

1.2.2 任务与应用

机器学习是面向实际应用的学科，旨在利用真实世界中的数据构建模型以解决特定的任务。而不同的任务场景往往具有不同的特性，因此这些具体的任务将会被分门别类，例如：根据任务标签是否连续，机器学习可以被分为回归任务和分类任务；根据训练数据的标注情况，机器学习可以被分为无监督任务，半监督任务和监督任务。这种任务分类方法帮助研究者在设计和选择模型时，更加针对性地考虑所面临的问题情境，从而提高模型的有效性。

图机器学习也不例外，但在图数据的背景下，传统的任务分类方法未必是区分图相关任务最重要的依据。事实上，在处理图相关任务时，研究者常会按照标签所属对象类别将这些任务划分为节点级任务，边级任务和图级任务。在本节中，将简要概述图机器学习中最主要的任务类别，并介绍这些任务在真实世界中的实际应用。

1. 节点级任务

现实场景中常常会面临这样的问题：在社交网络中预测用户的兴趣类别，在推荐系统中预测用户对商品的喜好程度，或者在金融网络中检查出现异常交易的账户，这些任务往往需要模型对图中节点进行预测，所以被统一归纳为节点级任务。其中，节点分类任务是最常见的节点级任务，也是当前图相关研究中的热点问题，其目标是基于图中部分节点的标注信息来预测其余节点的类别。另外，在某些任务场景中，模型被要求在不给定任何数据标注的情况下自主发掘节点的潜在类别，在这一类别任务中社区发现和异常检测是两个主要的研究热点。

在解决节点级任务时，数据集通常为整张图，如社交网络，金融网络，电商网络等，将图数据表示为 $G = (V, E)$ ，图中的节点 $v \in V$ 被定义为样本点，图中的边 $e \in E$ 提供样本点之间的关联关系。在模型进行推断时，会将图数据作为整体进行输入，来预测目标节点的标签信息。

1) 节点分类

该类任务旨在利用图中少量节点的标签信息，推断剩余未标注节点的标签^[4]，这类似于一般的监督学习范式。但不同点在于，图中的样本（节点）往往是相互依赖的，这打破了一般监督分类任务中对样本数据独立同分布（Independent and Identically Distributed, i.i.d）的假设。如图 1-12 所示，为了建模这种样本（节点）间的依赖关系，图数据在训练时会被被视为一个整体，这使得模型也能从验证集和测试集的无标注样本（节点）上进行学习，这种训练数据被部分标注的学习策略也被称为半监督学习。在某些场景中，任务需要对样本之间的显式关联进行建模，节点分类能够有效帮助模型利用样本间的关联关系进行预测。

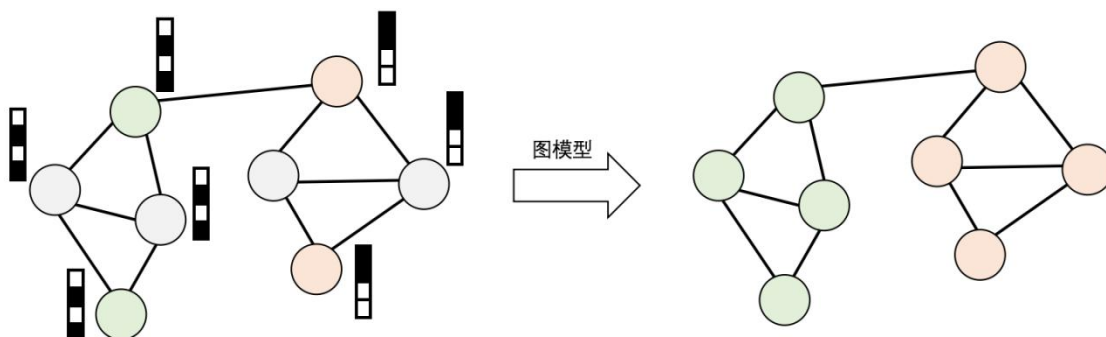


图 1-12 节点分类示例图

节点分类定义：给定图 $G = (V, E)$ ，其中节点集合被划分为训练集 V_{train} 和测试集 V_{test} ，即 $V = V_{\text{train}} \cup V_{\text{test}}$ 。在训练过程中，只有 V_{train} 中的节点有标签，而 V_{test} 中的节点没有标签。即存在映射 $g: V_{\text{train}} \rightarrow Y$ ，其中 $V_{\text{train}} \subset V$ 表示带有标签的训练集节点， Y 表示标签空间。节点分类任务的目的是学习一个映射函数 $f: V \rightarrow Y$ ，使得该函数能够将测试集 V_{test} 中节点的标签准确预测出来。

为了更好的利用无标注的节点信息，并建模节点之间的依赖关系，研究者提出了一些新的假设来指导节点分类模型的设计，例如：

- 同配性假设（Homophily Hypothesis）认为图数据中具有相似属性的节点更倾向于相互连接，这种同配现象在社交网络，推荐系统等中十分常见，比如人们趋向于与自己有共同兴趣爱好的人结交朋友，具有相似购买经历的用户更可能喜好相同的商品；
- 异配性假设（Heterophily Hypothesis）认为图数据中节点偏好于连接与自己有不同属性的节点，这与同配性假设正好相对，这种异配现象在一些具有互补性质的网络中较为常见，例如在蛋白质交互网络中，不同蛋白间具有互补的交互模式，因而具有交互连边的蛋白往往并不相似；
- 结构等价性（Structural Equivalence）认为具有相似局部邻域结构的节点，标签也会相似，在这种情况下，即使两个节点之间没有直接连接，如果它们与相同的节点相连，那么它们可能会拥有相似的标签或属性。

节点分类有非常丰富的应用场景，例如：在社交网络中可以根据用户的社交行为来预测用户对物品的喜好；在知识图谱中利用节点分类帮助识别实体的类型；在电商网络中基于用户与用户，用户与商品之间的交互，得到用户的商品偏好，实现个性化推荐等。

2) 社区发现

现实世界的图中经常会出现这样的连接模式，一些节点彼此之间连接较为紧密，而与外部其他节点之间的连接则较为稀疏，这些彼此紧密连接的节点集合被定义为社区，换句话说，社

区就是内部节点间连接紧密而与外部节点连接稀疏的节点集合。如图 1-13 所示，社区发现是专门发掘识别图中的社区结构的任务^[15]。与其他无监督节点级任务相比，社区发现通常更加关注图的拓扑结构，而不一定依赖于节点的具体特征。在某些场景中，社区内的节点通常会具有相似的功能特性，因而对于社区结构的挖掘有利于研究网络中的群体行为、传播模式等。

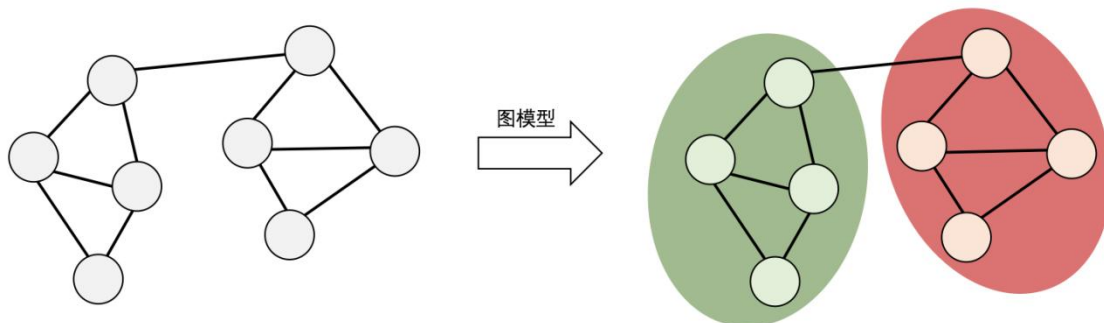


图 1-13 社区发现示例图

社区发现定义：给定一个图 $G = (V, E)$ 。社区发现的目标是在无标注的场景下，学习到映射函数 $f: V \rightarrow C$ ，其中 C 为社区标签空间。图中的节点 $v \in V$ 被映射到若干个社区 $\{C_1, C_2, \dots, C_k\}$ 中，且每个社区 $C_i \subseteq V$ 内的节点之间具有较高的连通性，而社区 C_i 与社区 C_j ($i \neq j$) 之间的连通性相对较弱。

由于社区是基于节点集合中的连接密度进行定义的，而边连接密度缺乏严格的度量标准，因此研究者提出了一些启发式的度量指标用于指导社区发现模型的设计，例如：

模块度 (Modularity) 是一种无标签下衡量社区划分质量的指标。它比较了社区内部的边数与在随机情况下的期望值，模块度越大表示社区划分质量越好。

标准化互信息 (Normalized Mutual Information) 是一种有标签下衡量社区划分质量的指标。它常用于衡量模型发现的社区和真实社区划分之间的相似程度，标准化互信息越大表示模型发现的社区与真实社区划分越相似。

社区发现具有许多实际应用场景，例如：在基因交互网络中识别功能模块，以揭示基因的生物学作用；在金融交易网络中揭示欺诈用户群体，以有效预防和检测金融犯罪；在社交网络分析中帮助理解用户关系，优化内容推荐；以及在推荐系统中细分用户群体，提供个性化的产品推荐。

3) 异常检测

异常检测是十分常见的机器学习任务，旨在从给定样本集合中找出与其他样本在某些方面显著不同的样本。传统的异常检测认为异常样本与正常样本在特征分布上有明显差异，因此会基于特征来进行推断。如图 1-14 所示，图上异常节点的检测则需要额外考虑节点间的异常连接模式，因此会基于特征和结构进行综合推断，这使得图上的异常检测更加具有挑战性^[16]。异常节点的检测可以帮助发现潜在的问题、异常活动或恶意行为，从而提高系统的安全性和可靠性。

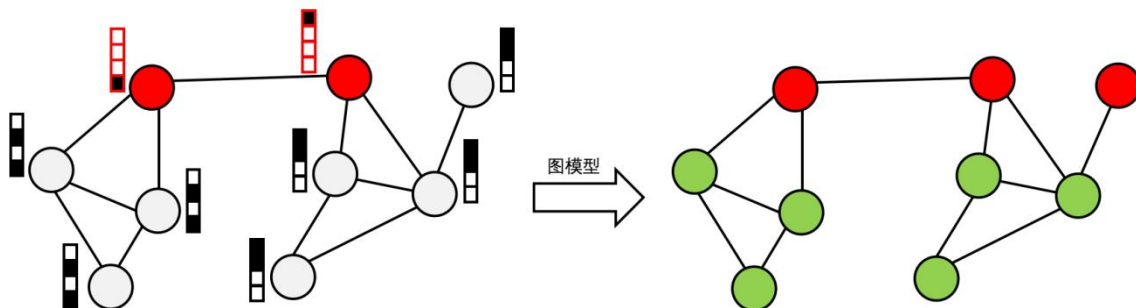


图 1-14 异常检测示例图

异常检测定义：给定一个图 $G = (V, E)$ ，其中每个节点 $v \in V$ 具有各自特征向量 X_v 。异常检测的目标是在无标注场景下，学习到映射 $f: (V \times X) \rightarrow Y$ ，其中 X 表示节点的特征空间， Y 表示异常程度空间。任务要求模型能够将与大部分节点具有明显特征分布差异或者结构差异的节点映射到更高的异常程度。

节点的异常模式往往具有多样性，研究者提出了不同假设来帮助更好的设计异常检测模型，例如：

异常样本在特征空间中往往较为孤立，因此可以通过启发式的异常度量指标来指导模型进行推断，例如局部离群因子(Local Outlier Factor, LOF)用于衡量样本在其邻域内的密度相对于邻域内其他节点的密度，如果样本的密度显著低于其邻域的密度，则被认为是异常。

相比正常样本，异常样本的分布模式难以被模型捕捉，因此可以通过重构的方式来进行推断，例如将重构误差较大的样本视为异常。

异常检测主要被应用在安全敏感的领域，例如在金融交易网络中检测交易模式异常的账户，以识别潜在的欺诈行为；在计算机网络中检测流量异常的 IP 地址，以防止网络攻击；在制造业中，通过设备间的交互检测出异常设备，以预测故障并减少停机时间。

2. 边级任务

现实场景中，除了需要对图中的节点进行预测，有时具体应用会更加关注节点之间是否存在交互关系，以及这种交互关系的类别，例如在社交网络中预测两个交互过的用户之间的关系类别（熟人，朋友，恋人等），在推荐系统中预测用户和商品之间是否有交互产生（点击，购买等）。这些任务往往需要模型预测给定的两个节点之间是否应该产生连边，以及连边的类型是怎么样的，前者通常在还未产生连边的节点之间进行预测，因此被定义为链接预测任务，而后者则通常会对已经产生的边进行预测，因此被定义为边分类。

类似于节点级任务，边级任务的数据集通常为整张图，其中链接预测会将图中所有的节点对 $\{v_i, v_j\}_{v_i, v_j \in V}$ 视为样本，而边分类则会将图中已经存在的边视为样本。模型在进行推断时，会将图数据作为整体进行输入，来预测节点对或边的标签信息。

1) 链接预测

链接预测是最常见的边级任务，其目标是预测图中节点对之间是否会形成边，因此可以被看做一个二元分类任务^[17]。由于图中已经存在显式连接的节点对信息，因此通常会将这些已经存在连边的节点对当做正样本加入到训练集中，而为了得到训练集中的负样本，一种常见的方式是从图中未建立连边的节点对中进行随机采样。然而如图 1-15 所示，不同于一般的二分类任务，链接预测除了要考虑节点对的特征，还需要利用节点已有的连接模式来进行预测，这使得模型设计时需要考虑更多的信息。另外，链接预测任务关注的是图中的潜在连接，因此其任务复杂度可能高达 $O(|V|^2)$ ，因此链接预测往往也更加需要关注模型的计算开销。

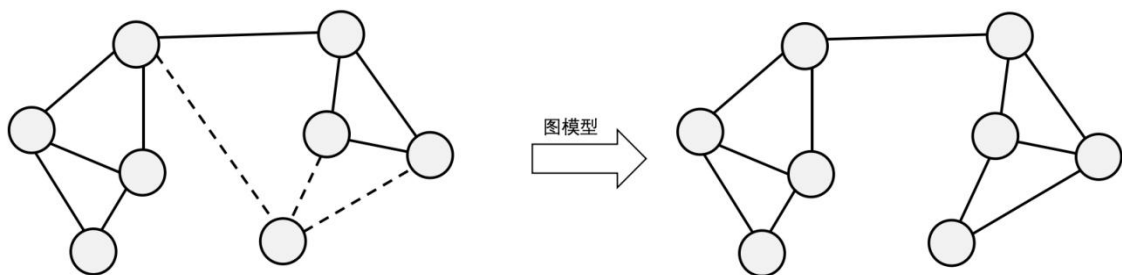


图 1-15 链接预测示例图

链接预测定义：给定一个图 $G = (V, E)$ ，链接预测任务的目标是基于图中已有的连接模式，即已经存在连边的节点对 $(v_i, v_j) \in E$ ，以及还未建立连边的节点对 $(v_i, v_j) \in E$ ，来学习到一个映射函数 $f: V \times V \rightarrow Y$ ，其中 $Y = \{0, 1\}$ 表示给定节点对之间是否会建立连边。

链接预测具有十分丰富实际应用场景，例如：在推荐系统中，通过链接预测可以发现用户可能感兴趣的商品或内容，在药物发现网络中，链接预测可以用于发现可能的药物-受体交互等。

2) 边分类

边分类任务的模型设计往往类似于链接预测模型，但也在某些方面会有所区别。包括采取的学习策略、关注的信息等。如图 1-16 所示，边分类任务通常采用半监督的学习策略，即在训练阶段会给定部分边的标签信息，因此同样会将边集合 E 划分为三个部分，即训练集 E_{train} ，验证集 E_{valid} 和测试集 E_{test} ；边分类任务关注的是已经存在的边而非潜在的边，使得任务复杂度从 $O(V^2)$ 下降到 $O(|E|)$ （通常情况下 $|E| \ll |V|^2$ ），这放宽了模型设计复杂度的要求；图数据中已存在边通常具有边属性，除了利用边两端节点对的特征外，还应该结合边自身的特征进行推断。

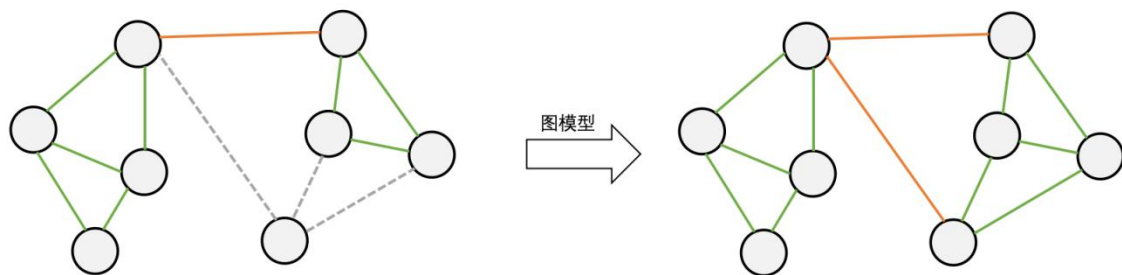


图 1-16 边分类示例图

边分类定义：给定一个图 $G = (V, E)$ ，其中边集合被划分为训练集 E_{train} 和测试集 E_{test} ，即 $E = E_{\text{train}} \cup E_{\text{test}}$ 。在训练过程中，只有 E_{train} 中的边有标签，而 E_{test} 中的边没有标签。即存在映射 $g: E_{\text{train}} \rightarrow Y$ ，其中 $E_{\text{train}} \subset E$ 表示带有标签的训练集， Y 表示标签空间。边分类任务的目的是学习一个映射函数 $f: E \rightarrow Y$ ，使得该函数能够将测试集 E_{test} 中边的标签准确预测出来。

边分类具有许多实际应用场景，例如：在金融交易网络，通过对交易边进行分类，可以检测异常交易或识别可能的欺诈行为；在交通物流网络中，通过对运输路径进行预测可以优化运输路线。

3. 图级任务

现实场景中，许多任务会将图本身作为预测的对象，而非图中的节点或者边。例如：在生物制药领域，常常会将原子表示为节点，原子间的化学键表示为边，使用整张图来建模药物分子，并依此预测分子的理化性质等；在人体姿态识别任务中，会将人体重要的关节视为节

点，关节间连接表示为边，使用整张图来建模人体骨架，并依此预测相应的姿态类别等。这些任务往往要求模型能将图视为一个整体进行推断，因此被归纳为图级任务。其中图分类，图匹配和图生成是图级别任务中最重要的子任务。

图级任务的建模方式与节点级和边级任务有所不同，图级任务的数据集往往不止一张图，将其表示为图集合的形式，即 $\mathcal{G} = \{G_1 \dots G_{|\mathcal{G}|}\}$ ，其中每张图 G_i 都代表一个独立的样本，图中的节点和边被视为样本特征的一部分。模型在进行推断时，通常会提取出图的整体特征进行预测。

1) 图分类

图分类是最常见的图级任务，其目标是预测给定图的标签类别。如图 1-17 所示，图分类任务往往采用监督学习的训练范式。具体来说，数据集集中的图会被分成三个部分，即训练集 $\mathcal{G}_{\text{train}}$ ，验证集 $\mathcal{G}_{\text{valid}}$ 和测试集 $\mathcal{G}_{\text{test}}$ ，只有训练集中的图样本 $G \in \mathcal{G}_{\text{train}}$ 得到了正确标注且会参与模型训练。

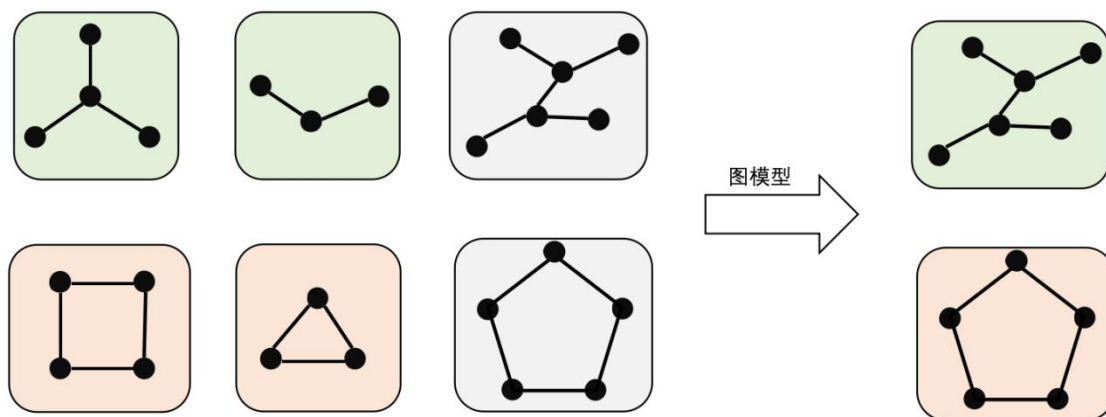


图 1-17 图分类示例图

图分类任务的定义：给定图集合 $\mathcal{G} = \{G_1 \dots G_{|\mathcal{G}|}\}$ ，其中图集合被划分为训练集 $\mathcal{G}_{\text{train}}$ 和测试集 $\mathcal{G}_{\text{test}}$ ，即 $\mathcal{G} = \mathcal{G}_{\text{train}} \cup \mathcal{G}_{\text{test}}$ 。在训练过程中，只有 $\mathcal{G}_{\text{train}}$ 中的图样本有标签，而 $\mathcal{G}_{\text{test}}$ 中的图样本没有标签。即存在映射 $g: \mathcal{G}_{\text{train}} \rightarrow Y$ ，其中 $\mathcal{G}_{\text{train}} \subset \mathcal{G}$ 表示带有标签的训练集， Y 表示标签空间。图分类任务的目的是学习一个映射函数 $f: \mathcal{G} \rightarrow Y$ ，使得该函数能够将测试集 $\mathcal{G}_{\text{test}}$ 中图的标签准确预测出来。

图分类任务依赖融合后的图表征，融合图中的节点和边信息的这一过程会丢失掉大量差异化信息，而图分类任务又需要模型对图样本间差异进行有效捕捉，这使得图分类任务尤其看中模型的表达能力。一种常见且有效的图模型表达能力的测试方法被称为多阶 Weisfeiler-Lehman 同构测试，也即 k-WL 测试，它是基于经典的图同构算法 WL 测试的改进版本。它将原本的 WL 测试从 1 阶邻域拓展到了 k 阶邻域，通过迭代节点特征更新，判断两个图是否为同构。假设如果两个图在 k-WL 测试中能够被区分开，那么它们也应该能够被图模型所区分。能够通过更高阶 WL-test 的模型通常被认为具有更强大的表达能力。

图分类有十分丰富的实际应用场景，如：在生物制药领域，通过湿实验来判断药物分子或蛋白质功能往往开销会十分昂贵，因此利用图机器学习模型来进行初步的性质预测能有效降低研发成本；在计算机视觉领域，通过捕捉不同图像中不同对象之间的交互关系能够更有效理解图像所包含的真实含义。

2) 图生成

在某些场景中，任务往往会要求模型基于需求生成对应的图数据，例如：在生物制药领域，基于已有的分子来生成具有相似结构和性质的分子，如图 1-18 所示，图模型需要基于给定的图数据集生成新的图数据。这需要模型能够理解和学习现有的图数据分布，以及生成新的图样

本。这对于研究图数据中潜在的图结构关系，理解现有数据中的模式、关联和隐藏的信息具有重要的意义^[18]。

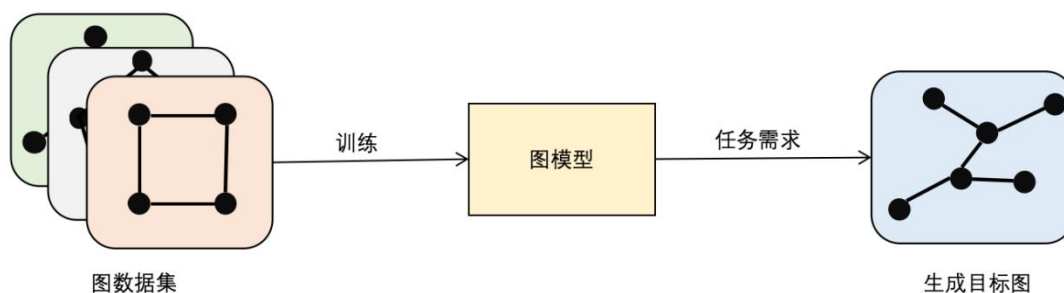


图 1-18 图生成示例图

图生成任务的定义：给定图集合 $\mathcal{G} = \{G_1 \dots G_{|\mathcal{G}|}\}$ 。图生成任务的目标是学习到这些图的分布 $P(\mathcal{G})$ ，然后通过采样的方式生成新的图样本 $G \sim P(\mathcal{G})$ 。也即任务希望模型学习到一个映射 $f: \mathcal{P} \rightarrow P(\mathcal{G})$ ，其中 \mathcal{P} 表示某种特定的分布（如高斯分布，均匀分布等）， $P(\mathcal{G})$ 表示给定图集合的真实分布，从而将随机分布的采样过程映射为图样本的采样过程。

图生成有许多实际应用场景，例如：在药物发现任务中，研究者希望模型能够基于给定的药物分子来生成结构或者性质相似的图结构，以此来帮助发现具有相似药理属性的药物。

1.3 图机器学习的发展历程

虽然图机器学习相较于传统机器学习起步较晚，但人们对图及其算法的研究却有着悠久的历史。从最开始使用图论算法进行图基本理论研究，到逐步将这些理论研究应用在现实应用上，并衍生出更为复杂的图算法。在这个过程中，图的规模不断扩大，复杂性也不断提升，原有的网络建模假设不足以支持后续更复杂的研究。因而在规则图网络和随机图网络之后，提出了更为接近现实且复杂的复杂网络，聚焦于其小世界、无尺度等特性进行研究。而随着互联网的兴起，社交网络成为了研究热点。在这一时期，以 PageRank 为代表的分析算法大放异彩，为社交网络、社区研究提供了有力支持。同时，图的规模呈指数性增大，直接在原有图上进行研究成本及复杂度变高。因而有学者开始研究如何将复杂高维的图数据降维到低维向量空间，同时又尽可能保持原有图特征及性质，即图嵌入技术。之后，以 RNN、CNN 为代表的神经网络崭露头角，学者开始思考如何将神经网络引入到图机器学习中以提高模型性能，从而形成了图神经网络的研究。而针对持续出现的新技术与应用，我们提出几点未来可能的研究热点与展望。接下来将按照“图论-图算法-复杂网络-社交网络分析-图嵌入-图神经网络-未来展望”的顺序介绍图机器学习的发展历程。通过对每一个时期的研究背景、研究内容、研究方法、优缺点等方面进行介绍，希望能帮助读者了解其发展历程，对图机器学习有一个大致的认识，并对未来的研究有所启发。

1.3.1 图论时期

关于图的研究最早可以追溯到 1736 年莱昂哈德·欧拉解决七桥问题^[19]，自此之后，图论开始成为数学家们热衷于研究的问题。在这个时期，图的研究还仅限于对图的定义及基本性质的研究，包含了对点、线、面的定义及相互之间的联系，诞生了欧拉定理、托兰定理等意义深远的定理。此外，数学家们还发现了现实世界中的许多问题实际上可以抽象为图论问题进行解决，例如欧拉问题、哈密顿问题、拉姆赛问题等。尽管学者在这个阶段的研究对于“巡回售货

员”、“周游世界棋盘”等现实问题都有着启发性的指导，但该时期图论的研究大部分还是局限于理论层面的讨论。

1.3.2 图算法时期

自 19 世纪中叶起，学者们对图研究有着进一步深化，图的定义越发丰富与严谨，更是为点、线、图赋予了特殊的属性。针对这些特殊属性也提出了图着色、网络流、图匹配等问题。

通过为点赋予标签属性，学者们得到了经典的图着色问题：对于无向连通图，给定一定数量的颜色对点进行着色（即赋予标签），是否能使得任意相连的两个节点具有不同的颜色的判定问题；以及为使得任意相连的两个节点具有不同的颜色，该提供多少种颜色的优化问题。此外 2007 年被以色列数学家艾夫拉汉·特雷特曼解决的道路着色问题^[20]也是经典的图着色问题。

通过为边赋予权重的属性，学者们有了以网络流问题为代表的路径规划问题：在这些问题中，重点关注于图中的部分路径，按照需求，求得路径权重总和最大/最小的规划。若希望总和最小，则可以使用 Dijkstra、Bellman-Ford^[21]、Floyd-Warshall^[22]等最短路径搜索算法，Prim 等最小生成树算法；若希望总和最大，则可以使用最大流最小割定理、Ford-Fulkerson^[23]等算法。

通过赋予子图结构特征，学者们得到了图匹配问题：例如在二分图匹配问题中，要求在二分图中找到一个匹配，使得两个子图间的连接最多。对于这样的问题，匈牙利算法、Hopcroft-Karp^[24]算法都是较为高效的解决方法。

除此之外，深度优先搜索算法、广度优先搜索算法、Tarjan^[25]算法、Kosaraju^[26]算法、A* 搜索算法等一系列图遍历算法也是图算法时期不可忽视的重要研究内容。对此也将衍生出许多现实问题。

图算法时期的研究比起图论时期的研究定义更为的丰富、严谨，研究的问题不仅多，更有许多现实应用。图着色问题对路径规划中的判定性问题有着重要影响；网络流问题的解决可以对交通拥挤、网络流量、任务分派等问题有所启发；二分图匹配问题的解决则对交友匹配、多任务流分配有所帮助。至于图遍历算法则对导航设计等规模较大且需要一定计算量的问题起着优化作用。

1.3.3 复杂网络时期

在对图的研究早期，学者们对图的建模普遍分为两种，规则图模型和随机图模型。在规则图模型中，图是按照一定的规则形成的，如完全图、星型图、平衡树等。此类模型起到基准模型的作用，方便于学者们对图基本理论进行研究。随机图模型则使用概率函数的方式随机连接节点以形成随机图，比较经典的有 Erdős-Rényi 模型^[27]。而随着图算法在现实应用的增多，学者们发现现实中的图有着更为复杂的性质，既不是完全依照规则，也不是完全依照概率函数的方式进行建模，而是位于两者之间，因而于 20 世纪 90 年代提出了复杂网络模型。其中 1998 年 Watts 提出的小世界特性^[28]与 1999 年 Barabasi 提出的无尺度网络^[29]成为了复杂网络模型的核心特点。

复杂网络是指在结构和动态特性上表现出复杂性特征的网络。这类网络通常具有自组织、自相似性、吸引子、小世界效应和无标度特性等性质。

小世界特性：对于大多数大规模的网络，在任意两个节点间应当有一条很短的路径相连，以形成小世界。其中六度分离理论就是一种很好的实例：你和任意一个陌生人最多只间隔六个人即可相识。该类特性让学者们意识到节点间关联的紧密性，并针对此有了更多的假设、研究。

无尺度网络：指的是一种大部分节点只与少部分节点相连，少部分节点与大部分节点相连的网络。这些少部分的节点成为了网络中的“枢纽”、“集散节点”，对于整个网络的结构稳定性起到了至关重要的作用。这样的结构也导致了无尺度网络是既具有鲁棒性又具有脆弱性的，当非枢纽节点损坏时，网络的结构基本不会受影响，此时网络是具有鲁棒性的；而当枢纽节点

损坏时，网络很容易因此崩溃，此时网络具有脆弱性。但考虑到枢纽节点与非枢纽节点数量分布与单节点负载的不同，该问题的考虑则会进一步复杂化。若是随机损坏则具有鲁棒性，若是蓄意损坏则具有脆弱性。此类网络所具有的幂律分布比起随机网络中的分布更为符合现实的大部分网络结构。大量的“二八效应”与图学习中的“长尾效应”都应证着这点。

复杂网络的引入，使得学者们在图上的研究更为贴近现实建模，为后续的研究提供了合理性假设。同时该建模方法也让学者们对图的研究不止局限于数学的研究方法，生物学、社会学、物理学等领域的统计方法也都能应用于复杂网络上，拓宽了研究的道路。

1.3.4 社交网络分析时期

在复杂网络建模方法逐步得到认可的同时，社交网络分析也因互联网的广泛兴起而成为研究热点。从传统相识社交网络的“六度分离理论”到 1996 年 Google 提出的对页面重要性分析的 PageRank 算法^[30]，社交网络分析的算法及应用正快速增加，并聚焦于两个问题：“节点重要性”和“社区发现”。

节点重要性：在复杂网络模型中，枢纽节点有着更为重要的作用。那么是否还有着其他重要节点，如何评估并找出这些节点成为了学者们研究的热点。在这里学者们提出了各种衡量指标来找出网络中的重要性节点，例如：基于节点度的度中心性，基于经过的最短路径个数的中介中心性，以及基于距离所有节点最短路径和的连接中心性等。此外 Google 提出的 PageRank 算法也能找出重要性节点。通过找出重要性节点，可以实现对重要性节点的特殊保护，诸如在搜索引擎中保护重要页面免受网络攻击；或者通过研究重要性节点，了解社交网络、交通网络、金融网络的结构及流动方式，甚至于找出生物网络中的核心基因、蛋白质分布。

社区发现：基于复杂网络的小世界效应，可以知道图中任意两个节点之间都存在着某种关联，但关联的强度不同。如何找出图中潜在的、有特定关联的、连接紧密的节点集合就是所谓的社区发现问题。这一问题的研究算法可以分为自上而下的分裂式算法与自下而上的聚类式算法。Girvan-Newman 算法^[31]是一种典型的分裂式算法，通过使用边介数，衡量边对网络连通性的影响，分割边介数最大的边以进行社区的划分。而聚类式算法中，普遍是将初始节点作为一个单独社区，通过衡量节点的模块度逐渐向上合并以形成社区，典型的有 Louvain 算法^[32]，模拟退火算法。但考虑到某一节点可能同时属于多个社区，原有的方法较难处理，因而提出了 Clique 渗透算法，通过完全子图的相邻关系来赋予节点多个社区标签。又由于社交网络可能存在动态的变化，因而出现了 PageRank、标签传播、InfoMap^[33]等沿边更新的、可以随时适应的算法。通过这一系列社区发现算法，可以发掘图中的潜在关系节点集合，对不同的集合采取不同的策略。诸如在购物网络中，按照社区的倾向推销不同的产品；或在金融风控网络中，对高危社区进行监管控制等。

1.3.5 图嵌入时期

随着图规模的不断扩大，直接在原有图上处理的传统方法有着较大的计算成本，学者们开始考虑如何通过降维、映射的方式来减少计算量，同时保证处理后的图能较大程度地保持原有图的信息。并且此前对图中节点特征大多是通过人工设置模版、特征工程等方法进行提取的，缺乏泛化性与效率。因而学者们于 21 世纪提出了一系列图嵌入算法，旨在设计优化目标，以优化算法的方式自动地学习节点表示。

通常的图嵌入模型会考虑如何对图的结构信息进行保持，这里的结构信息可以包括邻域结构信息、结构角色、社区信息、全局信息等。拉普拉斯映射的方式就是一种合理的结构信息保持的降维方法，该方法通过处理图的拉普拉斯矩阵，优化节点的向量表示使得相似的节点有着尽可能相似的向量表示。类似地，基于多种图表示矩阵，学者们可以使用 LINE^[34]、HOPE^[35]等矩阵分解的方法进行降维。这些方法在降维的过程中，除了保持邻近节点的相似性，还会注意对邻域相似的节点进行相似的向量表示。基于随机游走的图嵌入算法作为一种特殊的矩阵分解方法，也能对邻域结构信息进行保持。Node2Vec^[36]、DeepWalk^[37]等方法通过让节点进行随

机游走，在图上得到了随机游走路径的采样，作为学习的语料库。该语料库中的一条随机游走路径相当于一个句子，其中节点则相当于单词。随后通过 skip-gram 等方法对语料库内的随机游走路径进行共现概率的估计，以实现邻域相似的节点具有相似向量表示的目的。

图嵌入模型除了会考虑对图的结构信息进行保持，还会考虑对图中属性、标签之类的侧信息进行保持。Text-Associated DeepWalk 模型^[38]就是一种运用在属性图上的图嵌入模型。该方法首先证明了随机游走其实就是一种特殊的矩阵分解方式，然后通过结构信息矩阵对图的拓扑结构信息进行表示，再通过文本信息矩阵对图中的文本特征等属性进行捕捉，最后将两种矩阵进行分解、融合以得到保持了图中属性的节点特征表示。图嵌入方法因为往往是无监督学习，所以直接应用在分类任务上的性能较为薄弱。而 Max-Margin DeepWalk 模型^[39]就考虑保持原本图中的节点标签信息。该方法通过引入 SVM 等分类器，在对节点特征进行嵌入表示后，以节点标签为监督指标训练 SVM 模型。在优化过程中，该模型因为结合了特征嵌入模型与 SVM 模型的损失函数，所以采取控制变量的方法对两个模型分别进行优化，以得到更加适应分类任务的图嵌入表示。

1.3.6 图神经网络时期

随着深度学习在很多领域取得广泛的成功，神经网络成为了一种主流的技术手段。如何将 RNN、CNN 等经典算法引入图学习中，并适应图的特殊结构性性质成为了研究热点，因而于 21 世纪初期开始了图神经网络的研究。

图神经网络的起源最早可以追溯到 2005 年 Gori 等人提出了 GNN 的概念^[40]，基于不动点理论，使用 RNN 的思想，通过不断的迭代更新以收敛结果。后续随着 CNN 的出现，学者们将 CNN 的卷积算子推广到图神经网络上。图神经网络模型比起传统的图嵌入模型，具有多层网络结构与复杂连接方式，对节点特征有着更强的表达能力。下面我们将简单介绍几个代表性的图神经网络模型。

受传统 CNN 的启发，学者们将欧式几何空间推广到非欧几何空间上，在图的原始空间上进行卷积操作，得到了空域卷积图神经网络。空域卷积图神经网络的核心在于其消息传递范式，通过邻居节点特征变换、聚合邻居节点信息、更新目标节点特征来得到新的节点特征表示。

在 GCN^[41]中，目标节点的邻域节点会将相关特征通过映射函数进行转换；之后目标节点沿边收集邻域节点的特征，通过聚合函数将邻域信息进行整合；最后将整合信息与自身原有的节点特征一起进行目标节点的特征更新。

在消息传递范式的设计中，由于各类模型的核心算子不同，产生了各种变种。图注意力网络 GAT^[42]通过给每个边赋予注意力权重，在聚合过程中进行邻域信息的加权求和，使得目标节点在邻域信息的参考上有所偏重。而 GraphSAGE^[43]为了处理大规模图，在聚合过程中不对所有邻域节点进行消息接收，而是通过采样的方式，得到部分邻域信息以减缓模型的计算压力。

除此之外，由于图中加入了新的信息，学者们还为此设计出了一些特定模型。通过给图增加时间的维度，可以得到了动态图。由于动态图的图结构以及图节点输入会随着时间而变化，传统固定的图神经网络模型无法捕捉随时间变化的特征，因而提出了动态图神经网络模型。该类模型会对动态图中的时间信息与空间信息分别使用合适模型进行处理，之后将捕捉到的时间、空间信息进行结合以完成任务。具体的模型有 DCRNN^[44]、ST-GCN^[45]等。通过增加图中节点、边的种类，得到了异质图。由于异质图复杂的语义关系及多样的节点、边种类，传统图神经网络的参数、特征空间共享、语义理解等可能存在的问题，因而提出了异质图神经网络模型。针对复杂的语义关系，学者们提出了元路径的方法，通过元路径的学习来获取丰富的语义信息以弥补传统同质图神经网络的缺陷。具体的模型有 HAN^[46]、GTN^[47]、HetGNN^[48]等。

1.3.7 未来展望

在过去的十多年时间里，图机器学习不断产生新的研究方向与技术，并应用到很多领域。与此同时，图神经网络也存在很多问题亟待解决，在此我们简单提出几点未来可能的研究热点与展望。

1. 可信图神经网络

虽然图神经网络在处理图数据上已取得较大成果并将其应用于日常生活中，但现实中的应用对图神经网络的可信度方面有着更高的要求，因而提出了可信图神经网络^[49]。可信图神经网络主要聚焦于隐私性、鲁棒性、公平性及可解释性等性能的提升。通过提高模型的隐私性，避免泄露训练、使用过程中的隐私数据（如患者信息）。而鲁棒性的提高，则带来更为稳定的图神经网络系统，减少因遭受恶意攻击而崩溃的情况。同样的，由于图数据中不可避免地存在各种歧视，图神经网络又会加深其影响并做出有偏见的决策，因而公平性也是不容忽视的点。最后，由于神经网络的高度非线性性，其模型输出的可解释性往往遭人诟病。通过提高模型的可解释性，才能更好地让人放心地运用于现实应用中。

2. 以数据为中心的图模型

在图神经网络的发展历程中，学者们往往是通过设计更复杂的神经网络结构以提高模型能力。但该类方法随着模型的增大，正逐渐陷入瓶颈，深受训练成本、数据迁移等问题的困扰，因而兴起了以数据为中心的图模型^[50]。该类模型强调使用合适的数据增强模型能力，在数据准备、预处理、训练以及推理过程中对图数据的拓扑、特征、标签信息进行适当修改以适应模型，借此增强模型能力。同时由于现实的图数据可能存在脆弱性、不公平性、异质性问题，该类模型也会考虑如何从数据的角度去避免这些问题影响模型性能。

3. 多模态图模型

随着 ChatGPT 的横空出世，大语言模型的热潮席卷了各个研究领域，学者们都在研究如何将大语言模型与自己的研究领域进行结合，因而以大语言模型为纽带的通用人工智能再次获得了关注。得益于大语言模型强大的文本理解、逻辑推理能力、模型迁移能力，将各个领域的特定信息进行规范，再一同与大语言模型进行结合也成为了可能，这也就是所谓的多模态模型。在多模态模型中，不仅可以通过统一的模型框架分别处理不同领域的特定知识，还能将不同领域的特定知识进行融合以解决新的问题。而图数据作为其中一种模态数据，一方面可以让其与其他模型进行融合，如文本与图的逻辑推理能力融合，另一方面可以利用图的高度抽象性作为多模态数据规范化的桥梁，通过将文本、图像、生物分子等多模态数据建模成图的形式进行融合。

4. 图基础模型

图的增长与更新正不断地加快，但现如今大部分的图模型均是针对特定任务设计的，无法轻易地适应新的图。如何提高图的迁移能力，降低模型的训练时间与成本是值得研究的问题。受大语言模型强大的迁移能力启发，产生了图基础模型^[51]。图基础模型旨在设计一套通用性质的框架以处理图中的各种类型的问题，通过大量不同数据、任务的预训练，使模型适应各类下游任务。这些下游任务不仅可以是处理同一类型任务的不同数据集，还可以是处理预训练中未曾见过的任务类型。凭借在大量不同数据、任务的预训练，图基础模型学到图中的通用性知识，即使遇到未知任务，也能在适当的引导下高效解决。

1.4 本章小结

在引言中介绍了图机器学习的重要性和广泛应用，阐述了其对于复杂交互场景下的优秀建模和分析能力。第一节介绍图相关的基础知识，对于图的定义和表示，介绍了三种主要的表示方法：邻接矩阵表示法，邻接表表示法，关联矩阵表示法，并对不同的表示法进行了比较分析。

针对图的类型，介绍了无向图，有向图，无权图，带权图，同质图，异质图，无属性图，属性图，静态图和动态图这些常见图类型。第二节引入图机器学习，介绍了机器学习和图机器学习的基本概念，以及图机器学习的特点和难点。对于图机器学习的主要任务和应用，从节点级任务，边级任务和图级任务三个层面进行介绍，其中分别对节点级任务层面的节点分类、社区发现和异常检测，边级任务层面的链接预测和边分类，图任务层面的图分类和图生成进行了详细介绍，并说明了这些任务对应的真实应用场景。第三节介绍图机器学习的发展历程，按照时间顺序，分别针对图论时期，图算法时期，复杂网络时期，社交网络分析时期，图嵌入时期，图神经网络时期的主要研究内容和成果进行了介绍；另外，也对图机器学习未来的发展进行了展望，介绍了可能的未来研究热点，包括：可信图神经网络，以数据为中心的图模型，多模态模型和图基础模型。

拓展阅读材料

- (1) 克劳迪奥·斯塔迈尔等著，马京京译.《图机器学习》。
- (2) 刘忠雨，李彦霖，周洋著.《深入浅出图神经网络:GNN 原理解析》。
- (3) William Hamilton 著.《Graph Representation Learning》。
- (4) Albert-László Barabási 著.《Network Science》。
- (5) David Easley, Jon Kleinberg 著.《Networks, Crowds, and Markets: Reasoning About a Highly Connected World》。
- (6) 吴凌飞，崔鹏，裴健等著.《Graph Neural Networks: Foundations, Frontiers, and Applications》。

习题

- (1) 查找资料，举出一个图的现实应用案例，并指出其中点、线、图的含义。
- (2) 查找文献并思考，举例说明节点级、边级、图级任务。
- (3) 查找文献并思考，规范化拉普拉斯矩阵的意义是什么？
- (4) 计算图 1-19 的邻接矩阵、邻接表、关联矩阵、拉普拉斯矩阵以及标准化拉普拉斯矩阵。

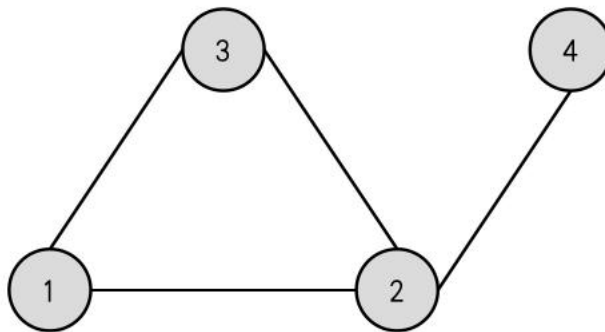


图 1-19 习题（4）图

参考文献

- [1] Gilmer J, Schoenholz S S, Riley P F, et al. Neural message passing for quantum chemistry[C]//International conference on machine learning. PMLR, 2017: 1263-1272.
- [2] Kipf T N, Welling M. Semi-supervised classification with graph convolutional networks[J]. arXiv preprint arXiv:1609.02907, 2016.
- [3] Ying R, He R, Chen K, et al. Graph convolutional neural networks for web-scale recommender systems[C]//Proceedings of the 24th ACM SIGKDD international conference on knowledge

- discovery & data mining. 2018: 974-983.
- [4] Belkin M, Niyogi P. Laplacian eigenmaps for dimensionality reduction and data representation[J]. *Neural computation*, 2003, 15(6): 1373-1396.
 - [5] He X, Niyogi P. Locality preserving projections[J]. *Advances in neural information processing systems*, 2003, 16.
 - [6] Ng A, Jordan M, Weiss Y. On spectral clustering: Analysis and an algorithm[J]. *Advances in neural information processing systems*, 2001, 14.
 - [7] Zhu X, Ghahramani Z, Lafferty J D. Semi-supervised learning using gaussian fields and harmonic functions[C]//*Proceedings of the 20th International conference on Machine learning (ICML-03)*. 2003: 912-919.
 - [8] Defferrard M, Bresson X, Vandergheynst P. Convolutional neural networks on graphs with fast localized spectral filtering[J]. *Advances in neural information processing systems*, 2016, 29.
 - [9] Jordan M I, Mitchell T M. Machine learning: Trends, perspectives, and prospects[J]. *Science*, 2015, 349(6245): 255-260.
 - [10] Voulodimos A, Doulamis N, Doulamis A, et al. Deep learning for computer vision: A brief review[J]. *Computational intelligence and neuroscience*, 2018, 2018(1): 7068349.
 - [11] Chowdhary K R, Chowdhary K R. Natural language processing[J]. *Fundamentals of artificial intelligence*, 2020: 603-649.
 - [12] Goyal A, Bengio Y. Inductive biases for deep learning of higher-level cognition[J]. *Proceedings of the Royal Society A*, 2022, 478(2266): 20210068.
 - [13] Bengio Y, Courville A, Vincent P. Representation learning: A review and new perspectives[J]. *IEEE transactions on pattern analysis and machine intelligence*, 2013, 35(8): 1798-1828.
 - [14] Xiao S, Wang S, Dai Y, et al. Graph neural networks in node classification: survey and evaluation[J]. *Machine Vision and Applications*, 2022, 33(1): 4.
 - [15] Fortunato S. Community detection in graphs[J]. *Physics reports*, 2010, 486(3-5): 75-174.
 - [16] Ma X, Wu J, Xue S, et al. A comprehensive survey on graph anomaly detection with deep learning[J]. *IEEE Transactions on Knowledge and Data Engineering*, 2021, 35(12): 12012-12038.
 - [17] Kumar A, Singh S S, Singh K, et al. Link prediction techniques, applications, and performance: A survey[J]. *Physica A: Statistical Mechanics and its Applications*, 2020, 553: 124289.
 - [18] Zhu Y, Du Y, Wang Y, et al. A survey on deep graph generation: Methods and applications[C]//*Learning on Graphs Conference*. PMLR, 2022: 47: 1-47: 21.
 - [19] Euler L. Leonhard Euler and the Königsberg bridges[J]. *Scientific American*, 1953, 189(1): 66-72.
 - [20] Trahtman A N. The road coloring problem[J]. *Israel Journal of Mathematics*, 2009, 172: 51-60.
 - [21] Bellman R. On a routing problem[J]. *Quarterly of applied mathematics*, 1958, 16(1): 87-90.
 - [22] Floyd R W. Algorithm 97: shortest path[J]. *Communications of the ACM*, 1962, 5(6): 345-345.
 - [23] Ford L R, Fulkerson D R. Maximal flow through a network[J]. *Canadian journal of Mathematics*, 1956, 8: 399-404.
 - [24] Hopcroft J E, Karp R M. An $n^{5/2}$ algorithm for maximum matchings in bipartite graphs[J]. *SIAM Journal on computing*, 1973, 2(4): 225-231.
 - [25] Tarjan R. Depth-first search and linear graph algorithms[J]. *SIAM journal on computing*, 1972, 1(2): 146-160.
 - [26] Atallah M J, Kosaraju S R. An efficient algorithm for maxdominance, with applications[J]. *Algorithmica*, 1989, 4(1): 221-236.
 - [27] Erdős P, Rényi A. On the evolution of random graphs[J]. *Publ. Math. Inst. Hungar. Acad. Sci*, 1960, 5: 17-61.
 - [28] Watts D J, Strogatz S H. Collective dynamics of ‘small-world’ networks[J]. *nature*, 1998, 393(6684): 440-442.
 - [29] Barabási A L, Albert R. Emergence of scaling in random networks[J]. *science*, 1999, 286(5439): 509-512.
 - [30] Page L. The PageRank citation ranking: Bringing order to the web[R]. *Technical Report*, 1999.
 - [31] Despalatović L, Vojković T, Vukicević D. Community structure in networks: Girvan-Newman algorithm improvement[C]//*2014 37th international convention on information and*

- communication technology, electronics and microelectronics (MIPRO). IEEE, 2014: 997-1002.
- [32] Blondel V D, Guillaume J L, Lambiotte R, et al. Fast unfolding of communities in large networks[J]. *Journal of statistical mechanics: theory and experiment*, 2008, 2008(10): P10008.
 - [33] Rosvall M, Bergstrom C T. Maps of information flow reveal community structure in complex networks[J]. *arXiv preprint physics.soc-ph/0707.0609*, 2007, 3.
 - [34] Tang J, Qu M, Wang M, et al. Line: Large-scale information network embedding[C]//*Proceedings of the 24th international conference on world wide web*. 2015: 1067-1077.
 - [35] Ou M, Cui P, Pei J, et al. Asymmetric transitivity preserving graph embedding[C]//*Proceedings of the 22nd ACM SIGKDD international conference on Knowledge discovery and data mining*. 2016: 1105-1114.
 - [36] Grover A, Leskovec J. node2vec: Scalable feature learning for networks[C]//*Proceedings of the 22nd ACM SIGKDD international conference on Knowledge discovery and data mining*. 2016: 855-864.
 - [37] Perozzi B, Al-Rfou R, Skiena S. Deepwalk: Online learning of social representations[C]//*Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining*. 2014: 701-710.
 - [38] Yang C, Liu Z, Zhao D, et al. Network representation learning with rich text information[C]//*IJCAI*. 2015, 2015: 2111-2117.
 - [39] Tu C, Zhang W, Liu Z, et al. Max-margin deepwalk: Discriminative learning of network representation[C]//*IJCAI*. 2016, 2016: 3889-3895.
 - [40] Gori M, Maggini M, Sarti L. Exact and approximate graph matching using random walks[J]. *IEEE transactions on pattern analysis and machine intelligence*, 2005, 27(7): 1100-1111.
 - [41] Kipf T N, Welling M. Semi-supervised classification with graph convolutional networks[J]. *arXiv preprint arXiv:1609.02907*, 2016.
 - [42] Velickovic P, Cucurull G, Casanova A, et al. Graph attention networks[J]. *stat*, 2017, 1050(20): 10-48550.
 - [43] Hamilton W, Ying Z, Leskovec J. Inductive representation learning on large graphs[J]. *Advances in neural information processing systems*, 2017, 30.
 - [44] Li Y, Yu R, Shahabi C, et al. Diffusion convolutional recurrent neural network: Data-driven traffic forecasting[J]. *arXiv preprint arXiv:1707.01926*, 2017.
 - [45] Yan S, Xiong Y, Lin D. Spatial temporal graph convolutional networks for skeleton-based action recognition[C]//*Proceedings of the AAAI conference on artificial intelligence*. 2018, 32(1).
 - [46] Wang X, Ji H, Shi C, et al. Heterogeneous graph attention network[C]//*The world wide web conference*. 2019: 2022-2032.
 - [47] Yun S, Jeong M, Kim R, et al. Graph transformer networks[J]. *Advances in neural information processing systems*, 2019, 32.
 - [48] Zhang C, Song D, Huang C, et al. Heterogeneous graph neural network[C]//*Proceedings of the 25th ACM SIGKDD international conference on knowledge discovery & data mining*. 2019: 793-803.
 - [49] Dai E, Zhao T, Zhu H, et al. A comprehensive survey on trustworthy graph neural networks: Privacy, robustness, fairness, and explainability[J]. *Machine Intelligence Research*, 2024, 21(6): 1011-1061.
 - [50] [50]Yang C, Bo D, Liu J, et al. Data-centric graph learning: A survey[J]. *arXiv preprint arXiv:2310.04987*, 2023.
 - [51] [51]Liu J, Yang C, Lu Z, et al. Towards graph foundation models: A survey and beyond[J]. *arXiv preprint arXiv:2310.11829*, 2023.