

学校代码 10530

学 号 201230111747

分 类 号 TP311

密 级

湘潭大学

硕士学位论文

最小生成树在不确定图中的应用与研究

学位申请人 唐 杰

指 导 老 师 文 中 华 教授

学 院 名 称 信息工程学院

学 科 专 业 计算机科学与技术

研 究 方 向 不 确 定 图

二〇一四年十一月三十日

最小生成树在不确定图中的应用与研究

学 位 申 请 人_____唐 杰_____

导师姓名及职称_____文 中 华 教授_____

学 院 名 称_____信息工程学院_____

学 科 专 业_____计算机科学与技术_____

研 究 方 向_____不 确 定 图_____

学位申请级别_____工 学 硕 士_____

学位授予单位_____湘 潭 大 学_____

论文提交日期_____2014-11-30_____

Application and Research of Minimum Spanning Tree in Uncertain Graph

Candidate _____ Jie Tang _____

Supervisor _____ Professor Zhonghua Wen _____

College _____ College of Information Engineering _____

Program _____ Computer Science and Technology _____

Specialization _____ Uncertain Graph _____

Degree _____ Master of Science _____

University _____ Xiangtan University _____

Date _____ April 20th, 2012 _____

湘潭大学

学位论文原创性声明

本人郑重声明：所呈交的论文是本人在导师的指导下独立进行研究所取得的研究成果。除了文中特别加以标注引用的内容外，本论文不包含任何其他个人或集体已经发表或撰写的成果作品。对本文的研究做出重要贡献的个人和集体，均已在文中以明确方式标明。本人完全意识到本声明的法律后果由本人承担。

作者签名：日期：年 月 日

学位论文版权使用授权书

本学位论文作者完全了解学校有关保留、使用学位论文的规定，同意学校保留并向国家有关部门或机构送交论文的复印件和电子版，允许论文被查阅和借阅。本人授权湘潭大学可以将本学位论文的全部或部分内容编入有关数据库进行检索，可以采用影印、缩印或扫描等复制手段保存和汇编本学位论文。

涉密论文按学校规定处理。

作者签名：日期：年 月 日

导师签名：日期：年 月 日

摘 要

图论作为数学的一个分支，在计算机科学中也扮演着非常重要的角色。许多问题都可以通过抽象化之后转化为图论问题，然后通过计算机将对应的图论算法实现来得到问题的解。随着大数据时代的到来，图数据的规模也随之增大，因此图算法的效率对于解决图论问题至关重要。

由于现实问题中可能会遇到很多的意外情况，如机器故障，信号传输的不稳定性等，最终导致实际数据存在不确定性。因此不确定图已经成为一个非常热门的研究领域。本文对不确定图中的最小生成树进行了研究，主要内容分为以下三个部分：

- (1) 由于不确定图的边存在不确定性，所以传统的最小生成树定义已经不能完全适用于不确定图。因此提出了最优生成树的概念，它主要是在最小生成树定义的基础上增加了对概率的最优化定义。然后我们借助贪心的思想设计了不确定图中最优生成树求解算法。
- (2) 由于最优生成树在不确定图中不一定完全稳定，所以我们又设计了Top-K最优生成树算法，并提出了基于并查集的优化算法和基于启发式搜索的优化算法。当 K 值较小时，基于启发式搜索的优化算法具有明显的优势。
- (3) 在内容(2)的算法基础上，我们对最小生成树在不确定图中的可靠性进行了分析。

最后，我们对不确定图中最优生成树的研究进行了相关的实验验证，实验结果表明，设计的算法的运行效率与理论分析基本一致，且能够完全正确的得到问题的解。

关键词：不确定图，蕴含图，最小生成树，最优生成树。

Abstract

Key words: uncertain graph, implicated graph, minimum spanning tree, optimized spanning tree.

目 录

摘要	I
Abstract	II
第一章 绪 论	1
1.1 研究背景和意义	1
1.2 研究现状	2
1.3 本文的结构安排	3
第二章 不确定图中最优生成树和次优生成树算法	4
2.1 问题定义	4
2.2 最优生成树算法	6
2.3 次优生成树算法	7
2.4 本章小结	9
第三章 不确定图中Top-K最优生成树算法	10
3.1 问题定义	10
3.2 基本算法	10
3.3 并查集优化	12
3.4 启发式搜索A*优化	17
3.5 实验结果	20
3.6 本章小结	20
第四章 不确定图中最小生成树的可靠性研究	21
4.1 问题定义	21
4.2 最小生成树的可靠性求解算法	21
4.3 实验结果	23
4.4 本章小结	24
第五章 总结与展望	26
参考文献	27

致 谢	30
附录A (攻读硕士学位期间发表的论文)	31
附录B (攻读硕士学位期间参与的科研项目)	32
附录C (攻读硕士学位期间获奖情况)	33

第一章 绪 论

§1.1 研究背景和意义

1736年,瑞士数学家欧拉针对哥尼斯堡七桥问题发表了图论的首篇论文,奠定了图论的基础。此后,图论经过漫长的发展,直到1936年哥尼格发表了第一本图论专著,从此图论成为一门独立的学科。时至今日,图论已经渗透到了许多领域,它在自然科学和社会科学中都有广泛的应用。

在生物信息学中,蛋白质交互网络^[1, 2]用来记录蛋白质之间的交互作用。我们可以用无向图来表示蛋白质交互网络,其中顶点用来表示蛋白质,边用来表示蛋白质之间的交互;在无线传感器网络中^[3, 4],传感器和他们之间的通信可以用图来表示;在社交网络中(例如Facebook、Twitter),图可以用来描述人与人之间的社会关系。这些由不同领域产生的图结构也可以成为“图数据”,由于他们之间具有很多的共性,因此研究图数据本身的性质具有重要的意义。

由于数据获取技术的随机错误与测量误差、数据传输的故障与延时、多源集成数据的不完整性和不一致性等多种原因,大量的图数据具有不确定性。例如,生物信息学中,由于蛋白质交互的高通量生物检测技术存在固有误差,因此实验测得的蛋白质交互是否真实存在是不确定的^[5-8];在道路网络中,由于全球定位系统的数据传输故障和延时,当前网络中的交通流量数据无法确切放映实际情况,因而具有不确定性。在无线传感器网络中,由于物理世界的干扰、无线传感器网络的不稳定性、传感器节点的易失效性等原因,两个传感器之间是否可以正常通信是不确定的^[9, 10]。我们称这些具有不确定因素的图数据为“不确定图数据”或“不确定图”。

在无线传感器网络中,一般都需要在监测区域布置大量的传感器,而每个传感器的能量是有限的,使用最小的能耗保证整个传感器网络的畅通显然具有重要的意义^[11]。在理想情况下,我们可以将传感器网络用一个带权无向图来表示,权值代表两个传感器之间单位时间内通信所需的能耗,这样,可以使用经典的最小生成树算法(Prim、Kruskal)来获得一颗总能耗最小的网络拓扑结构树。然而,由于两个传感器之间的通信可能会出现故障,导致某些传感器之间的通信链路断开,这时我们需要做到能够快速的更新网络拓扑结构来保证网络的畅通。因此对不确定图中最小生成树的研究具有重要意义。

由于信息技术的飞跃发展,图数据的规模也越来越大,不确定因素对图数据的影响也越来越得到重视。国内外许多学者已经开始对不确定图数据展开了一些列研

究,也取得了很多的成果。在国内,邹等人已经出版了不确定图数据挖掘一书,国际上也出版了Managing and Mining Graph Data^[12]和Mining Graph Data^[13]等书籍。在一些国际顶级会议上(如SIGKDD、ICDE等),也专门设立了不确定图数据的相关议程。

综上所述,对不确定图数据的研究具有重要的科学意义和广泛的应用价值。无论是与其它学科相互结合还是将传统的图论问题应用到不确定图模型上,都有很多待解决的问题。

§1.2 研究现状

本节主要对不确定图在国内外已有的研究进行大致分类并进行简要的分析。

(1) 频繁子图模式挖掘

频繁子图模式挖掘也是确定图数据挖掘中的经典问题^[14-17],目的是从一组确定图中发现频繁出现的子图模式。而在不确定图数据中,现有的研究主要集中在期望频繁子图模式挖掘^[18-20]、概率频繁子图模式挖掘^[21, 22]和稠密子图挖掘^[23, 24]。前两者的主要目的是求解支持度和置信度满足一个指定阈值的频繁子图模式,在不确定图中的稠密子图挖掘具有代表性的有Top-K极大团挖掘和紧密顶点子集挖掘。

(2) 分类和聚类

我们知道,分类和聚类算法在数学、统计学、计算机科学和生物学等领域都有很多应用。在不确定图数据中,对分类算法的研究主要是K近邻查询算法^[25-27]。例如,文献[25]提出了基于SimRank^[28]度量的方法来求解给定节点的 k 个最相似的节点,文献[27]也对不确定图中的KNN进行了研究,它提出了三种方法(median-distance, majority-distance, expected-reliable-distance)来度量两个节点之间的距离。然而对聚类算法的研究主要有层次聚类算法^[29, 30]。文献[29]首先将求出的强联通子图作为聚类中心,然后对余下的顶点进行层次聚类,文献[30]主要是基于LinLog能量模型及其Newman-Girvan的模块化聚类,提出了一种多层次模块化聚类算法框架。

(3) 基于最短路、最大流等图论模型

在图论中,Dijkstra、Prim和Edmonds-Karp算法都是用来解决图论中经典问题(最短路、最小生成树和最大流)的方法,然而这些方法并不能直接应用在不确定图数据中。因而,一些学者对不确定图中的最短路径、最大流等问题进行了研究^[31-34],文献[31]首先定义了概率语义下的最短路径问题,然后设计了求大于给定阈值的所有最短路径的算法;文献[33]给出了最大流可靠性的概率计算模型,提出了一种基于简单路径组合思想的算法和基于状态空间划分的最可靠最大流改进算法。

§1.3 本文的结构安排

第一章 简要介绍了不确定图的研究背景、国内外研究现状等相关知识。

第二章 定义了不确定图中的最优生成树和次优生成树等概念，并设计了相应的求解算法。

第三章 针对Top-K最优生成树算法进行了深入分析，且根据K值的不同提出了两种优化算法。

第四章 对不确定图中最小生成树的可靠性进行研究。

第五章 对全文进行了总结，并对未来的相关工作进行探讨。

第二章 不确定图中最优生成树和次优生成树算法

§2.1 问题定义

定义 2.1.1 (不确定图) 不确定图是一个四元组 $\mathcal{G} = (V, E, W, P)$, 其中 V 是顶点集, E 是边集, $W = \{w(e) | e \in E, w(e) \in \mathbb{N}^+\}$ 是边的权重集, $P = \{p(e) | e \in E, p(e) \in (0, 1]\}$ 是边存在可能性的集合。

定义 2.1.2 (蕴含图) 令不确定图 $\mathcal{G} = (V, E, W, P)$, 若确定图 $G = (V_G, E_G, W_G)$ 是 \mathcal{G} 的一个蕴含图, 则必然满足 $V_G = V$, $E_G \subseteq E$ 和 $W_G = \{w(e) | e \in E_G\} \subseteq W$ 。

从定义 2.1.1 可知, 不确定图 $\mathcal{G} = (V, E, W, P)$ 每条边以 $p(e)$ 的概率存在, 在可能世界模型下, 每条边有存在和不存在两种可能性, 所以可以派生出 $2^{|E|}$ 个蕴含图。本文沿用文献[5, 8, 12]对不确定图模型所做的假设, 即不确定图中不同边的概率分布相互独立。将蕴含图 G 和不确定图 \mathcal{G} 之间的关系表示为 $\mathcal{G} \Rightarrow G$ 。基于以上假设, 蕴含图 G 存在的概率为

$$\Pr(\mathcal{G} \Rightarrow G) = \prod_{e \in E_G} p(e) \prod_{e \in E \setminus E_G} (1 - p(e)). \quad (2.1)$$

记不确定图 \mathcal{G} 的所有蕴含图的集合为 $\text{Imp}(\mathcal{G})$ 。由文献[6]可知 \mathcal{G} 中所有蕴含图出现的概率和为1, 即

$$\sum_{G \in \text{Imp}(\mathcal{G})} \Pr(\mathcal{G} \Rightarrow G) = 1. \quad (2.2)$$

例 1 图 2.1(a) 为一个不确定图 \mathcal{G}_1 , 边上的两个数字分别代表权值和概率, 由于该不确定图有五条可能出现的边, 因此该不确定图有 2^5 个蕴含图。图 2.1(b) 为不确定图 \mathcal{G}_1 的一个蕴含图, 显然该蕴含图存在的概率为 $p(e_2) \times p(e_3) \times p(e_4) \times p(e_5) \times (1 - p(e_1)) = 0.03528$ 。

在传统图论中, 图的最小生成树被定义为边的权值和最小的生成树, 一个图可能存在多个最小生成树, 这些最小生成树具有相同的权值和。然而在不确定图中, 每条边都有一个存在的概率, 这样会导致每一颗最小生成树都只有一定的概率存在, 我们可以通过一个稳定性来区分最小生成树的好坏, 即

$$R_T = \prod_{e \in E_T} p(e) \quad (2.3)$$

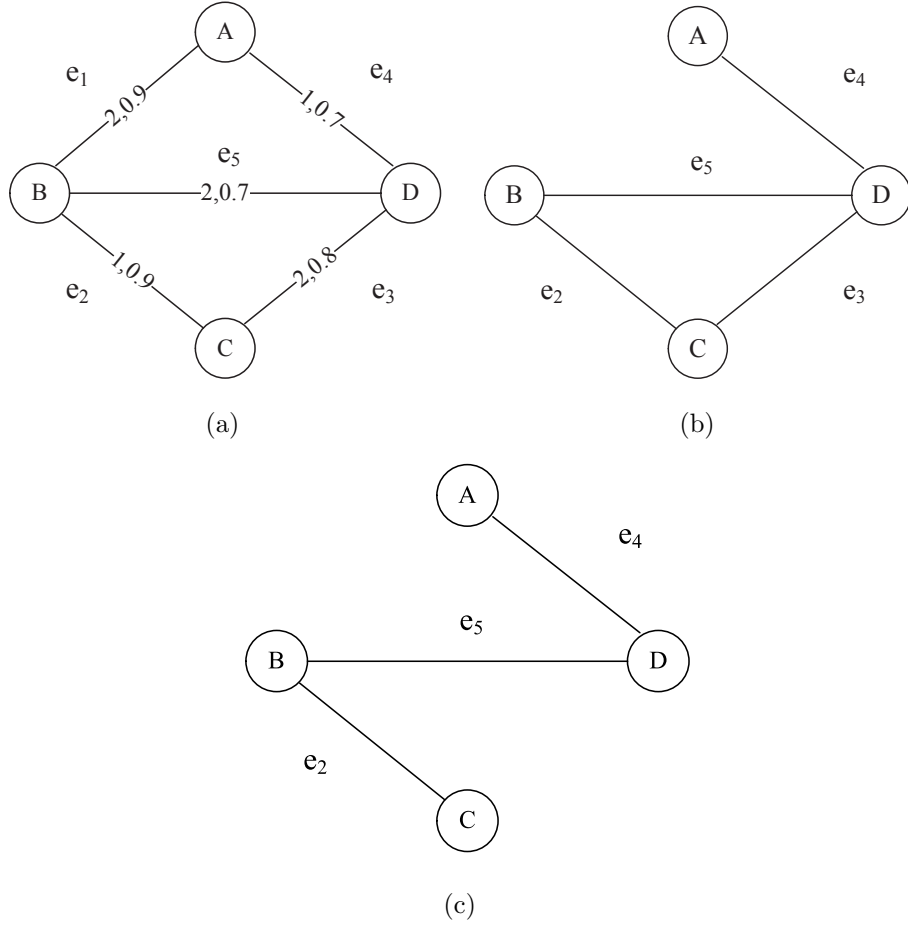


图 2.1 例 1: (a) 不确定图 \mathcal{G} ; (b) 图 \mathcal{G} 的一个蕴含图; (c) 图 \mathcal{G} 的生成树。

公式 2.3 的含义与公式 2.1 有所不同, 公式 2.3 的另一层含义可以解释为所有包含最小生成树 T 的蕴含图存在的概率和。

为了区分每颗生成树, 我们需要给每个生成树进行编号, 假设所有生成树的边都按照编号从小到大进行排序, 我们称字典序较小的生成树具有更小的编号。

定义 2.1.3 (最小生成树) 设不确定图 $\mathcal{G} = (V, E, W, P)$, 若生成树 T 满足 $\forall T' (\sum_{e \in E_T} W(e) \leq \sum_{e' \in E_{T'}} W(e'))$, 则称 T 为不确定图 \mathcal{G} 的最小生成树, 记编号最小的最小生成树为 $T_M^{\mathcal{G}}$, 其边集为 $E_M^{\mathcal{G}}$, 边的权值和为 $W_M^{\mathcal{G}}$ 。

定义 2.1.4 (最小最大乘积生成树) 设不确定图 $\mathcal{G} = (V, E, W, P)$, 若生成树 T 满足 $\forall T' (\prod_{e \in E_T} P(e) \leq \prod_{e' \in E_{T'}} P(e'))$, 则称 T 为不确定图 \mathcal{G} 的最小乘积生成树, 记编号最小的最小乘积生成树为 $T_{PL}^{\mathcal{G}}$, 其边集为 $E_{PL}^{\mathcal{G}}$, 边的权值和为 $W_{PL}^{\mathcal{G}}$ 。若生成树 T 满足 $\forall T' (\prod_{e \in E_T} P(e) \geq \prod_{e' \in E_{T'}} P(e'))$, 则称 T 为不确定图 \mathcal{G} 的最大乘积生成树, 记编号最小的最大乘积生成树为 $T_{PB}^{\mathcal{G}}$, 其边集为 $E_{PB}^{\mathcal{G}}$, 边的权值和为 $W_{PB}^{\mathcal{G}}$ 。

定义 2.1.5 (最优生成树) 设 $Imp(\mathcal{T})$ 为不确定图 \mathcal{G} 中所有最小生成树的集合, 若最小生成树 T 满足 $\forall T' \in Imp(\mathcal{T}) (R_{T'} \leq R_T)$, 则称 T 为不确定图 \mathcal{G} 的最优生成树。记编号最小的最优生成树为 $T_O^{\mathcal{G}}$, 其边集为 $E_O^{\mathcal{G}}$, 边的权值和为 $W_O^{\mathcal{G}}$ 。

显然, 最优生成树不是唯一的, 我们也可以将最优生成树看作不确定图中的最小生成树中的最大乘积生成树。

§2.2 最优生成树算法

定理 2.2.1 设不确定图 $\mathcal{G} = (V, E, W, P)$ 和 $\mathcal{G}' = (V, E, W, P')$, 其满足 $P'(e) = \log_2(P(e))$, 则 $W_{PL}^{\mathcal{G}} = 2^{W_M^{\mathcal{G}'}}$ 。

证明 设 T 为 \mathcal{G} 和 \mathcal{G}' 的任意一颗生成树, 其边集为 E_T , 则 $\sum_{e \in E_T} P'(e) = \log_2 \prod_{e \in E_T} P(e)$, 显然生成树 T 在 \mathcal{G}' 中的边权和与在 \mathcal{G} 中边权的乘积成正比, 当 T 为 \mathcal{G}' 的最小生成树时, 则有 $W_M^{\mathcal{G}'} = \log_2 W_{PL}^{\mathcal{G}}$, 即 $W_{PL}^{\mathcal{G}} = 2^{W_M^{\mathcal{G}'}}$ 。 ■

由定理 2.2.1 可知, 要想求得最小乘积生成树, 可以将图中所有边的权值进行 \log 变换, 然后求新图的最小生成树。我们将在后面使用该思路去求解不确定图的最优生成树。

推论 2.2.1 若生成树的边权和是最小的, 那么其边权的乘积也是最小的; 若生成树的边权和是最大的, 那么其边权的乘积也是最大的。

我们知道使用 *kruskal* 算法求解最小生成树的第一步是需要对图中所有的边进行排序, 因此可以假设不确定图 $\mathcal{G} = (V, E, W, P)$ 中所有的边按照如下优先级进行排序:

- 将权值较小的边排在前面;
- 对于权值相同的边, 则将概率较大的边排在前面;
- 若两者的值都相同, 则将编号较小的边排在前面。

我们称靠前的边具有更高的优先级, 最优生成树的算法步骤如下:

1. 新建一个不确定图 \mathcal{G}' , 其顶点集和不确定图 \mathcal{G} 相同, 边集为空;
2. 假设不确定图 \mathcal{G} 的边集已经按照上面的规则进行排序;

3. 依次处理排序好的边, 若这条边连接的两个顶点不在同一个连通分量中, 则将这条边加入不确定图 \mathcal{G}' ;
4. 重复第3步, 直到不确定图 \mathcal{G}' 中所有的顶点都连通。

这样, 我们得到的不确定图 \mathcal{G}' 既是不确定图 \mathcal{G} 的最优生成树, 显然该最优生成树也是编号最小的最优生成树 $T_O^{\mathcal{G}}$ 。

定理 2.2.2 最优生成树算法是正确的。

证 明 显然最优生成树的算法步骤和 $Kruskal$ 算法的步骤一致, 因此我们得到的生成树一定是最小生成树。根据最小生成树的性质, 所有的最小生成树应该具有相同的边权值序列 (假设边权值按照从小到大排序), 并且相同权值的边构成的连通分量应该相同, 根据最优生成树算法可知, 由相同权值的边构成的连通分量, 其边的概率和是最大的。再由推论 2.2.1可知, 这些边的概率乘积也是最大的, 因此最优生成树算法是正确的。 ■

例 2 图 §2.2 为一个不确定图 \mathcal{G} , 按照优先级排序后的边集为 $\{e_1, e_2, \dots, e_{10}\}$ 。图 §2.2 为图 §2.2 中不确定图的最优生成树求解过程, 红色标记的边为加入生成树中的边, 在图 (2) 到图 (3) 时, 由于顶点 D 和顶点 E 已经在同一个连通分量, 所以边 e_3 没有加入生成树的边集中, 图 (6) 中红线构成的生成树则为最优生成树。从图中可知, $W_O^{\mathcal{G}} = 13$ 且 $R_{T_O^{\mathcal{G}}} = 0.14112$ 。

算法 1 为最优生成树算法的伪代码, 我们使用了并查集数据结构来快速合并两个连通分量, 每个连通分量都会有一个唯一的根节点, 通过函数 $GETROOT()$ 即能快速查询到每个顶点的根节点。合并连通分量只需要一条赋值语句即可, 函数 $OST()$ 中的第 7 行则是合并顶点 u 和顶点 v 所在的连通分量。

§2.3 次优生成树算法

定义 2.3.1 (次优生成树) 设 $Imp(\mathcal{T})$ 为不确定图 \mathcal{G} 中所有最小生成树的集合, 若最小生成树 T 满足 $\forall T' \in (Imp(\mathcal{T}) \setminus \{T_O^{\mathcal{G}}\}) (R_{T'} \leq R_T)$, 则称 T 为不确定图 \mathcal{G} 的次优生成树。记编号最小的次优生成树为 $T_{SO}^{\mathcal{G}}$, 其边集为 $E_{SO}^{\mathcal{G}}$, 边的权值和为 $W_{SO}^{\mathcal{G}}$ 。

由定义 2.3.1 可知, 当最优生成树不唯一时, 次优生成树也是最优生成树。显然, 我们可以借用求次小生成树的思想求解次优生成树, 步骤如下

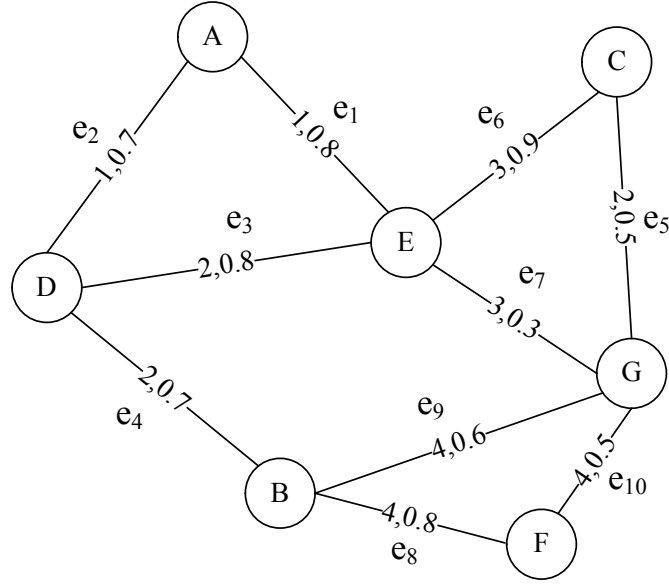


图 2.2 例 2: 不确定图 \mathcal{G} 。

Algorithm 1 最优生成树算法

模块1: 并查集

initialization $RT_i \leftarrow i (i = 1, 2, \dots, |V|)$
1: **procedure** GETROOT(u)
2: **if** $u = RT_u$ **then return** u
3: **else return** ($RT_u \leftarrow \text{GETROOT}(RT_u)$)

模块2: 最优生成树算法

1: **procedure** OST(\mathcal{G}) $\triangleright \mathcal{G} = (V, E, W, P)$
2: $E_O^\mathcal{G} = \emptyset$
3: **for each** $e_i \leftarrow (u, v) \in E$ **do**
4: $ru \leftarrow \text{GETROOT}(u)$
5: $rv \leftarrow \text{GETROOT}(v)$
6: **if** $ru \neq rv$ **then**
7: $RT_{ru} = rv$
8: $E_O^\mathcal{G} = E_O^\mathcal{G} \cup \{e_i\}$
9: **return** $E_O^\mathcal{G}$

1. 选择一条不在生成树 $T_O^\mathcal{G}$ 中的边加入生成树 $T_O^\mathcal{G}$;

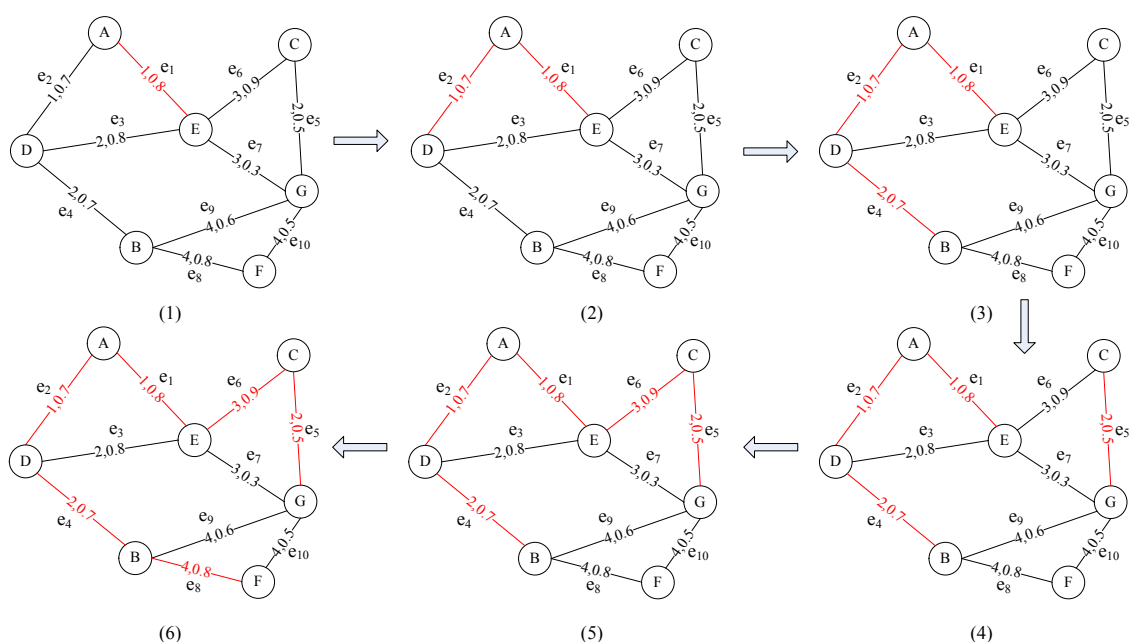


图 2.3 例 2: 最优生成树计算过程。

2. 假设新加入的边为 $e = (u, v)$, 寻找生成树 T_O^G 中从顶点 u 到顶点 v 的路径中优先级最高的边并删除。这样可以得到一颗新的生成树;
3. 重复前面2步, 知道所有的边都已经处理, 这样我们能够得到 $|E| - |V| + 1$ 颗新的生成树, 从这些生成树中求出最优的生成树则是不确定图 G 的次优生成树。

在第二步中, 我们需要知道最优生成树中任意两个顶点之间优先级最高的边, 我们可以在求解最优生成树时进行预处理得到:

1. 设 $F(i, j)$ 为在最优生成树中从顶点 i 到顶点 j 的路径中优先级最高的边;
2. 在求解最优生成树时, 假设边 $e = (u, v)$ 加入了边集 E_O^G , 则对于图中所有顶点 x , 应该进行如下转换:
 - (a) 如果边 e 比边 $F(x, u)$ 优先级更高, 则 $F(x, u) = e$;
 - (b) 如果边 e 比边 $F(x, v)$ 优先级更高, 则 $F(x, v) = e$;

§2.4 本章小结

第三章 不确定图中Top-K最优生成树算法

§3.1 问题定义

§3.2 基本算法

定理 3.2.1 设 $E_T = \{e_1, e_2, \dots, e_{n-1}\}$ 为不确定图 $\mathcal{G} = (V, E, W, P)$ 的任意一颗最小生成树的边集, 其中 $n = |V|$, 图的最小生成树可以通过 *Prim* 或 *Kruskal* 算法求得, 令 $A_i = \{e_1, e_2, \dots, e_i\}$, 其中 $i \in \{1, 2, \dots, n-1\}$ 。我们将蕴含图作以下划分:

$$D_i = \begin{cases} \{G | G \in \text{Imp}(\mathcal{G}) \wedge e_1 \notin E\}, & i = 0 \\ \{G | G \in \text{Imp}(\mathcal{G}) \wedge A_i \subseteq E \wedge e_{i+1} \notin E\}, & 0 < i < n-1 \\ \{G | G \in \text{Imp}(\mathcal{G}) \wedge A_i \subseteq E\}, & i = n-1 \end{cases} \quad (3.1)$$

则等式 $\bigcup_{i=0}^{n-1} D_i = \text{Imp}(\mathcal{G})$ 成立, 且 $D_i \cap D_j = \emptyset$ ($i < j$)。设 $A_0 = \emptyset$ 且不确定图不存在边 e_n , 则公式 3.1 可以统一表示为第二种情况。

证明 (1) $\bigcup_{i=0}^{n-1} D_i = \text{Imp}(\mathcal{G})$ 。设不确定图的所有蕴含图的集合 $\text{Imp}(\mathcal{G}) = \{G_1, G_2, \dots, G_r\}$, 其中 $r = |\text{Imp}(\mathcal{G})|$ 。若 $E_{G_i} \cap E_T = \{e_{x_1}, e_{x_2}, \dots, e_{x_y}\}$, 其中 $x_j < x_{j+1}$ ($j = 1, 2, \dots, y-1$), 则必然存在一个 k 满足:

$$\begin{cases} x_j = x_{j-1} + 1, & j \leq k \\ x_j > x_{j-1} + 1, & j = k \end{cases}$$

可以将蕴含图 G_i 进行如下归类:

$$\begin{cases} G_i \in D_0, & x_1 \neq 1 \\ G_i \in D_k, & x_1 = 1 \end{cases}$$

因此可以得到 $\text{Imp}(\mathcal{G}) \subseteq \bigcup_{i=0}^{n-1} D_i$, 由 D_i 的定义可知 $D_i \subseteq \text{Imp}(\mathcal{G})$ ($i = 0, 1, \dots, n-1$), 由此可知 $\bigcup_{i=0}^{n-1} D_i \subseteq \text{Imp}(\mathcal{G})$, 因此 $\bigcup_{i=0}^{n-1} D_i = \text{Imp}(\mathcal{G})$ 。

(2) $D_i \cap D_j = \emptyset$ ($i < j$)。由公式 3.1 可知, $D_i \cap D_j = \{G | G \in \text{Imp}(\mathcal{G}) \wedge A_j \subseteq E \wedge e_{i+1} \notin E\}$, 由于 $e_{i+1} \in A_j$, 所以 $A_j \subseteq E$ 和 $e_{i+1} \notin E$ 不可能同时成立, 因此 $D_i \cap D_j = \emptyset$ 。

■

引理 3.2.1 设 $E_T = \{e_1, e_2, \dots, e_{n-1}\}$ 为不确定图 $\mathcal{G} = (V, E, W, P)$ 的任意一颗最小生成树的边集, 令边集 $E_I = \{e_1, e_2, \dots, e_k\}$ 为 E_T 的一个子集; 边集 $E_O \subseteq E$ 且 $E_O \cap E_T = \emptyset$; 若集合 $S = \{G \mid G \in \text{Imp}(\mathcal{G}) \wedge E_I \subseteq E_G \wedge E_O \cap E_G = \emptyset\}$, 即集合 S 为包含边集 E_I 中的所有边且不包含边集 E_O 中的所有边的蕴含图集合, 则我们可以将集合 S 进行如下划分:

$$D_i = \{G \mid G \in S \wedge E_{I_i} \subseteq E_G \wedge E_{O_i} \cap E_G = \emptyset\}$$

其中 $E_{I_i} = E_I \cup \{e_{k+1}, e_{k+2}, \dots, e_i\}$, $E_{O_i} = E_O \cup \{e_{i+1}\}$, 显然 $\bigcup_{i=k}^{n-1} D_i = S$ 成立。

由引理 3.2.1 可知, 我们可以通过边集划分的方法来获取所有的最小生成树, 然后将最小生成树的边乘积按照从大到小的顺序排列, 选取前 K 个最小生成树, 即可得到 Top-K 最优生成树。算法 2 为 TOP-K 最优生成树算法, 函数 $\text{GETMST}(E_I, E_O, \mathcal{G})$ 将返回一颗不确定图 \mathcal{G} 中的最小生成树, 且该最小生成树的边集一定包含 E_I 且不包含 E_O 。函数 $\text{FINDALLMST}()$ 中的第 6 行如果成立则表示 E_I 已经是一颗最小生成树的边集, 所以不用继续划分, 直接加入集合 MSTSet 即可。

例 3 图 3.1(a) 为一个不确定图, 图 3.1(b) 为不确定图中主蕴含图的其中一颗最小生成树, 其中边的权值和为 4, 图 3.2 为函数调用 $\text{FINDALLMST}(\emptyset, \emptyset, \mathcal{G})$ 的全过程, 我们使用了树形结构来表示, 节点中左右两边表示传递给被调用函数的参数, 左边表示可以包含的边的集合, 右边表示不可以包含的边的集合, 带有粗边框的节点表示到此为止已经找到了一颗最小生成树, 这颗生成树的边集为节点左侧的边集, 粗边框右上角的数字表示这个生成树中所有边的概率乘积。显然, 最优生成树和次优生成树的边集分别为 $\{e_2, e_3, e_5\}$ 和 $\{e_2, e_3, e_6\}$ 。

定理 3.2.2 算法 2 的时间复杂度为 $O(Nnm)$, 其中 N 为不确定图 \mathcal{G} 中最小生成树的个数, n 和 m 分别代表不确定图 \mathcal{G} 中的顶点数和边数。

证明 显然, 每次调用函数 FindAllMST 时, 需要调用一次函数 GetMST 并且产生 $n - 1$ 次的递归调用。这样, 函数 FindAllMST 最多被调用 $Nn + 1 \approx Nn$ 次。然而, 函数 GetMST 求解最小生成树的时间复杂度为 $O(m)$ 。因此, 算法 2 的时间复杂度为 $O(Nnm)$ 。 ■

Algorithm 2 TOP-K最优生成树算法

模块1: TOP-K最优生成树算法

```

1: procedure TOPKMST( $K, \mathcal{G}$ )  $\triangleright \mathcal{G} = (V, E, W, P)$ 
2:    $MSTSet = \text{FINDALLMST}(\emptyset, \emptyset, \mathcal{G})$ 
3:   将集合 $MSTSet$ 中所有最小生成树按照边乘积的值从大到小排序
4:   return 集合 $MSTSet$ 中前 $K$ 个最小生成树
  
```

模块2: 查找所有最小生成树

```

1: procedure FINDALLMST( $E_I, E_O, \mathcal{G}$ )
2:    $E_T \leftarrow \text{GETMST}(E_I, E_O, \mathcal{G})$ 
3:    $k \leftarrow |E_I|$ 
4:    $MSTSet \leftarrow \emptyset$ 
5:   if  $E_T = E_I$  then
6:      $MSTSet = MSTSet \cup \{T\}$ 
7:   else
8:     for  $i \leftarrow k$  to  $|V| - 1$  do
9:        $E_{I_i} \leftarrow E_I \cup \{e_{k+1}, e_{k+2}, \dots, e_i\}$ 
10:       $E_{O_i} \leftarrow E_O \cup \{e_{i+1}\}$ 
11:       $MSTSet \leftarrow MSTSet \cup \text{FINDALLMST}(E_{I_i}, E_{O_i}, \mathcal{G})$ 
12:   return  $MSTSet$ 
  
```

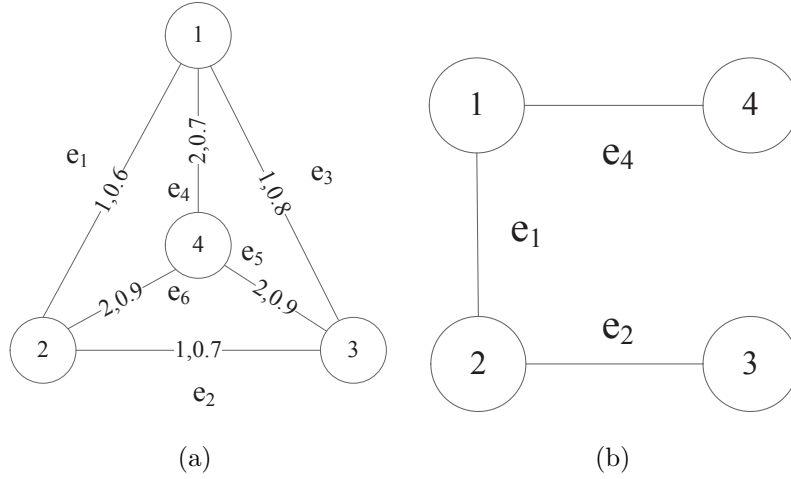


图 3.1 例 3: (a) 不确定图 \mathcal{G} ; (b) 图 \mathcal{G} 中的一颗最小生成树

§3.3 并查集优化

在算法 2 中, 我们每次都需要调用函数 $\text{GETMST}()$ 来获取一颗新的最小生成树,

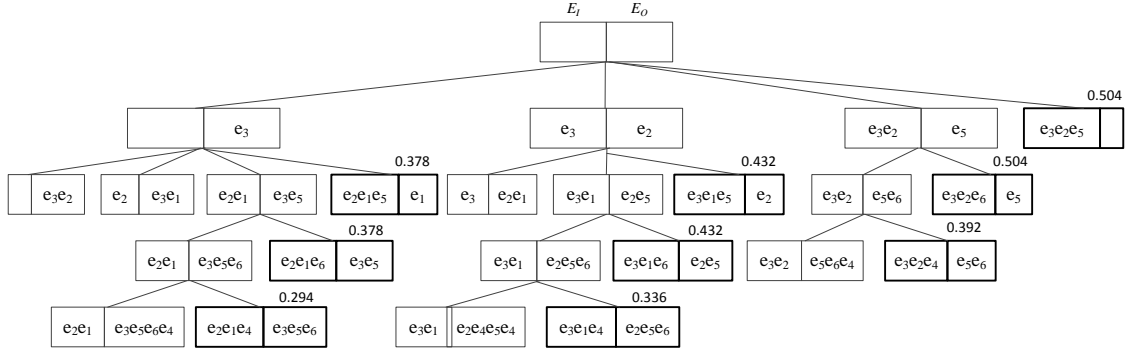


图 3.2 例 3: 函数FindAllMST($\emptyset, \emptyset, \mathcal{G}$)的递归调用过程。

由于新的生成树和上层调用的生成树有很多相似之处, 因此在这一小节中, 我们将对该方法进行改进, 我们希望通过一个边的替换操作, 能够很快的获得一颗新的最小生成树, 假设我们用 $T_1 \xrightarrow[e']{e} T_2$ 表示一个边替换操作, 该操作表示在生成树 T_1 中添加边 e 并删除 e' 后可以得到生成树 T_2 , 即 $T_2 = T_1 \cup \{e'\} \setminus \{e\}$, 我们称边 e' 为边 e 的替换边。

定理 3.3.1 在不确定图 \mathcal{G} 中, 假设 T_a 为一颗生成树, T_b 为一颗最小生成树, 若 $E_{Diff} = \{e_1, e_2, \dots, e_k\}$ 表示这两颗生成树中不同边的集合, 即满足 $e_i \in E_{T_a}$ 并且 $e_i \notin E_{T_b}$ 。则我们可以通过重复执行下面的边替换操作 k 次, 使得 T_a 变为 T_b , 即

$$T_a(T_0) \xrightarrow[e'_1]{e_1} T_1 \xrightarrow[e'_2]{e_2} T_2, \dots, \rightarrow T_{k-1} \xrightarrow[e'_k]{e_k} T_b(T_k),$$

且

$$\sum_{e \in E_{T_{i-1}}} w(e) \geq \sum_{e \in E_{T_i}} w(e).$$

证 明 我们可以重复执行下面的操作 k 次:

1. 假设边 e_i 连接的两个顶点分别为 u 和 v , 若从生成树 T_{i-1} 中移除边 e_i , 则可以得到两个连通分量, 我们将其设为 V_1 和 V_2 ;
2. 若在生成树 T_b 中从顶点 u 到顶点 v 的路径为 $u(u_0) \rightarrow u_1, \dots, u_{x-1}, \rightarrow v(u_x)$ (我们也可以用 P_v^u 来表示这个路径), 那么必然存在 $i (i < x)$ 使得 $u_i \in V_1$ 且 $u_{i+1} \in V_2$ 。令 $e'_i = (u_i, u_{i+1})$, 由于 T_b 是最小生成树, 显然 $w(e'_i) \leq w(e_i)$, 因此我们可以通过边替换 $T_{i-1} \xrightarrow[e'_i]{e_i} T_i$ 得到生成树 T_i 且 $\sum_{e \in T_i} w(e) = \sum_{e \in T_{i-1}} w(e) + w(e'_i) - w(e_i) \leq \sum_{e \in T_{i-1}} w(e)$ 。

定理 3.3.2 设 T 为不确定图 \mathcal{G} 中的最小生成树, 假设边 $e \in E_T$ 连接的两个连通分量为 V_1 和 V_2 。那么对于任意不包含边 e 的最小生成树 T' , 必然存在一条边 $e' = (u', v')$ 满足 $w(e') = w(e)$, 其中 $u' \in V_1, v' \in V_2$ 。

证明 假设边 e 连接的两个顶点分别为 u 和 v , 在生成树 T' 中的路径 P_v^u 为 $u(u_0) \rightarrow u_1, \dots, u_{x-1}, \rightarrow v(u_x)$ 。那么必然存在一个 $i(i < x)$ 满足 $u_i \in V_1$ 且 $u_{i+1} \in V_2$, 设 $e' = (u_i, u_{i+1})$, 则有下列两种情况:

- (1) 如果 $w(e') > w(e)$, 可以设 $E_{T^*} = E_{T'} \cup \{e\} \setminus \{e'\}$, 那么 $\sum_{e \in E_{T^*}} w(e) < \sum_{e \in E_{T'}} w(e)$;
- (2) 如果 $w(e') < w(e)$, 可以设 $E_{T^*} = E_T \cup \{e'\} \setminus \{e\}$, 那么 $\sum_{e \in E_{T^*}} w(e) < \sum_{e \in E_T} w(e)$ 。

显然, 以上两种情况与 T 和 T' 是最小生成树冲突, 因此, $w(e_i) = w(e)$ 。

设 $E(V_1, V_2, \mathcal{G}) = \{e \in E \mid e \in V_1 \times V_2\}$ 。根据定理 3.3.2可知, 如果存在边 e' 满足 $e' \in E(V_1, V_2, \mathcal{G})$ 且 $w(e') = w(e)$, 那么我们可以通过边替换操作 $T \xrightarrow[e']{e} newT$ 得到一颗新的最小生成树, 即 $E_{newT} = T \cup \{e'\} \setminus \{e\}$, 否则该不确定图不存在其它包含边 e 的最小生成树。

为了发现在生成树 T 中每条边的交换边, 我们可以按照权值从小到大遍历不确定图 $\mathcal{G} = (V, E, W, P)$ 中的所有边, 对于每条边 $e^* \in E$, 我们将执行下面两个步骤:

1. 假设边 e^* 连接的两个顶点分别为 u^* 和 v^* , 那么我们可以遍历路径 $P_{v^*}^{u^*}$ 。显然这条路径上所有权值和边 e^* 权值相等的边都可以被边 e^* 替换, 即边 e^* 为他们的替换边;
2. 处理完路径 $P_{v^*}^{u^*}$ 上所有边后, 则可以使用并查集来压缩路径 $P_{v^*}^{u^*}$ 。

为了查找最小生成树 T 中任意两个顶点的路径, 我们需要知道任意两个顶点之间的关系。首先, 我们选择一个顶点作为树的根节点, 然后通过深度优先搜索对每个顶点进行编号, 这样就可以通过编号区间来判定两个顶点之间的关系, 我们定义函数GETCHILD为:

$$\text{GETCHILD}(T, w) = \{v \mid v \text{ is a child-node of } w \text{ in } T\}$$

则算法 3对子树中的节点进行编号, 显然我们可以调用函数SETNUMBER($T, 1, 1$)为生成树 T 中的所有顶点进行编号。

Algorithm 3 子树编号算法

```
1: procedure SETNUMBER( $T, w, num$ )
2:    $L_w \leftarrow num, R_w \leftarrow num$ 
3:    $C \leftarrow \text{GETCHILD}(T, w)$ 
4:   for each  $v \in C$  do
5:      $R_w \leftarrow \text{SETNUMBER}(T, v, R_w + 1)$ 
6:   return  $R_w$ 
```

定义 3.3.1 (最近公共祖先) 假设 T 为不确定图 G 中的最小生成树，编号为1的顶点为根节点，则顶点 u 和顶点 v 的最近公共祖先为生成树 T 中的顶点，该顶点即是顶点 u 的祖先也是顶点 v 的祖先，且该顶点的编号最小。

例 4 图 3.3是通过函数调用SETNUMBER($T, 1, 1$)对顶点编号后的最小生成树 T 。根据算法 3可知， L_i 是通过前序遍历得到的结果，而 $R_i = \max\{L_j \mid \text{顶点} j \text{ 为顶点} i \text{ 的子孙节点}\}$ ，即 R_i 为以顶点 i 为根节点的子树中 L_i 的最大值。

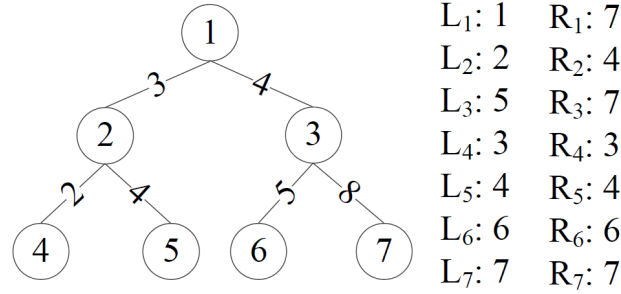


图 3.3 编号后的最小生成树。

由图 3.3可知，如果 $L_j \in [L_i, R_i]$ ，则顶点 j 是顶点 i 的子孙节点，设 u_1, u_2 为生成树 T 的两个顶点，则我们可以通过以下步骤发现他们之间的路径：

1. 从顶点 u_1 开始向父节点进行遍历，直到找到一个顶点 v 满足 $L_{u_2} \in [L_v, R_v]$ ；
2. 从顶点 u_2 开始向父节点进行遍历，直到顶点 v 停止。

这样，路径 $P_{u_2}^{u_1}$ 可以表示为 $u_1 \rightarrow \cdots \rightarrow v \rightarrow \cdots \rightarrow u_2$ ，且顶点 v 为顶点 u_1 和 u_2 的最近公共祖先。

接下来，我们使用并查集来压缩路径 $P_{u_2}^{u_1}$ 。设 RT_i 为顶点 i 每次被处理后记录的最近公共祖先，显然对于每个在路径 $P_{u_2}^{u_1}$ 中的顶点 w ，满足 $RT_w = v$ 。设函数GETFATHER(T, x)返回顶点 x 的父节点。则算法 4为替换边求解算法。

Algorithm 4 查找生成树 T 中每条边的替换边。

Module 1: 并查集

initialization $RT_i \leftarrow i (i = 1, 2, \dots, n)$
1: **procedure** GETROOT(u)
2: **if** $u = RT_u$ **then return** u
3: **else return** $(RT_u \leftarrow \text{GETROOT}(RT_u))$

Module 2: 查找替换边

1: **procedure** FINDSWAPEDGE(G, T)
2: $S = \emptyset$
3: **for each** $e \leftarrow (u, v) \in E_G$ **do**
4: $lca \leftarrow ru \leftarrow \text{GETROOT}(u)$
5: $rv \leftarrow \text{GETROOT}(v)$
6: **while** $rv \notin [L_{lca}, R_{lca}]$ **do**
7: $f \leftarrow \text{GETFATHER}(T, lca)$
8: $lca \leftarrow \text{GETROOT}(f)$
9: Let P_{rv}^{ru} be $ru \rightarrow \dots \rightarrow lca \rightarrow \dots \rightarrow rv$.
10: **for each** edge e' in path P_{rv}^{ru} **do**
11: **if** $w(e') = w(e)$ **then**
12: $S \leftarrow S \cup \{(e', e)\}$
13: **for each** vertex w in path P_{rv}^{ru} **do**
14: $RT_w \leftarrow lca$
15: **return** S

例 5 如图 3.4所示, 设 $e_1 = (2, 5), e_2 = (3, 6), e_3 = (4, 5), e_4 = (5, 6), w(e_3) = 4, w(e_4) = 5$ 且 $e_3, e_4 \in V$ 。由于 $w(e_3) < w(e_4)$, 所以我们首先查找顶点4和顶点5之间的路径, 由于 $w(e_1) = w(e_3)$, 所以边 e_3 为边 e_1 的替换边。显然顶点2为顶点4和顶点5的最近公共祖先, 我们用顶点2来替换路径 P_5^4 , 则 $RT_4 = 2$ 且 $RT_5 = 2$ 。之后我们再处理边 e_4 , 我们知道顶点5已经被其最近公共祖先替换, 这样 $e_4 = (5, 6) = (2, 6)$, 又因为 $w(e_2) = w(e_4)$, 所以边 e_4 为边 e_2 的替换边, 因此集合 $S = \{(e_1, e_3), (e_2, e_4)\}$ 。

算法 5为优化后的算法。

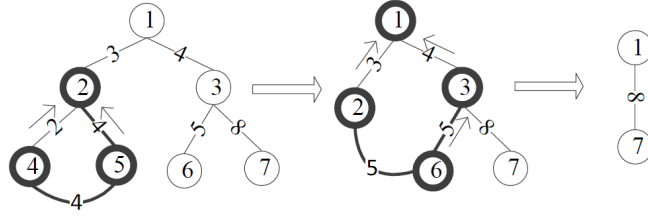


图 3.4 路径压缩过程。

定理 3.3.3 算法 5 的时间复杂度为 $O(Nm)$, 其中 N 为不确定图 \mathcal{G} 中最小生成树的个数, m 为不确定图 \mathcal{G} 中的边数。

证明 首先, 函数 `SETNUMBER` 只需要对生成树搜索一遍, 所以函数 `SETNUMBER` 的时间复杂度为 $O(n)$; 其次, 函数 `FINDSWAPEDGE` 需要遍历不确定图 \mathcal{G} 中的每一条边, 由于使用了路径压缩, 生成树中的每条边则只会被访问一次, 所以函数 `FINDSWAPEDGE` 的时间复杂度为 $O(m)$ (假设 $m > n$); 这样, 使用边替换操作查找一颗新的最小生成树的平均时间复杂度为 $O(m/n)$, 因此算法 5 的时间复杂度为 $O(Nn * (m/n)) = O(Nm)$ 。

■

§3.4 启发式搜索 A* 优化

由算法 2 和算法 5 可知, 当参数 K 很小时, 该算法也需要获取所有的最小生成树, 时间复杂度显然很高, 因此我们将使用启发式搜索来进行优化。

定理 3.4.1 设 T 为不确定图 $\mathcal{G} = (V, E, W, P)$ 中的一颗最小生成树且 $e \in E_T$, 边 e 划分的两个连通分量分别为 V_1 和 V_2 , 设 $e' \neq e$ 为集合 $E(V_1, V_2, \mathcal{G})$ 中边权值第二大的边, 则边集 $E_{T'} = E_T \cup \{e'\} \setminus \{e\}$ 对应的生成树 T' 为不确定图 $\mathcal{G}' = (V, E \setminus \{e\}, W, P)$ 中的最小生成树。

证明 我们使用反证法来证明, 假设 T' 不是不确定图 \mathcal{G}' 中的最小生成树。那么, 必然存在两条边 $c \notin T'$ 和 $d \in T'$ 满足: (1) 存在生成树 T'' , 其边集 $E_{T''} = E_{T'} \cup \{c\} \setminus \{d\}$; (2) $w(c) < w(d)$ 。显然 $c \in E(V_1, V_2, \mathcal{G})$, 否则边集 $E_{\tilde{T}} = E_T \cup \{c\} \setminus \{d\}$ 对应的生成树 \tilde{T} 的边权值和比最小生成树 T 更小, 这样与 T 是最小生成树矛盾, 因此我们可以得到

$$w(e) \leq w(e') \leq w(c) < w(d)$$

Algorithm 5 TOP-K最优生成树算法-并查集优化。

模块1： TOP-K最优生成树算法-并查集优化

```

1: procedure TOPKMST( $K, \mathcal{G}$ )  $\triangleright \mathcal{G} = (V, E, W, P)$ 
2:    $E_T \leftarrow \text{OST}(\mathcal{G})$ 
3:    $MSTSet = \text{FINDALLMST}(\emptyset, \emptyset, E_T, \mathcal{G})$ 
4:   将集合  $MSTSet$  中所有最小生成树按照边乘积的值从大到小排序
5:   return 集合  $MSTSet$  中前  $K$  个最小生成树

```

模块2： 查找所有最小生成树

```

1: procedure FINDALLMST( $E_I, E_O, E_T, \mathcal{G}$ )
2:    $MSTSet \leftarrow \{T\}$ 
3:    $G \leftarrow (V, E \setminus E_O, W)$ 
4:   SETNUMBER( $T, 1, 1$ )
5:    $S \leftarrow \text{FINDSWAPEDGE}(G, T)$ 
6:   for  $i \leftarrow k$  to  $|V| - 2$  do
7:      $E_{I_i} \leftarrow E_I \cup \{e_{k+1}, e_{k+2}, \dots, e_i\}$ 
8:      $E_{O_i} \leftarrow E_O \cup \{e_{i+1}\}$ 
9:     if  $\exists e'((e_{i+1}, e') \in S)$  then
10:       $E_{T'} \leftarrow E_T \cup \{e'\} \setminus \{e_{i+1}\}$ 
11:       $MSTSet \leftarrow MSTSet \cup \text{FINDALLMST}(E_{I_i}, E_{O_i}, E_{T'}, \mathcal{G})$ 
12:   return  $MSTSet$ 

```

设边集 $E_T \cup \{e'\}$, $E_T \cup \{c\}$, $E_{T'} \cup \{c\}$ 组成的环路分别为 C_0, C_1, C_2 。显然 $C_2 = C_0 \otimes C_1$ 且 $d \in C_2$ ，这样， $d \in C_0$ 和 $d \in C_1$ 至少有一个成立，我们定义：

$$E_{\hat{T}} = \begin{cases} E_T \cup \{e'\} \setminus \{d\}, & d \in C_0 \\ E_T \cup \{c\} \setminus \{d\}, & d \in C_1 \end{cases} \quad (3.2)$$

那么，生成树 \hat{T} 的权值和比最小生成树 T 更小，产生矛盾，因此原定理成立。 ■

显然，定理 3.4.1 也同样适用于最小(最大)乘积生成树。由于启发式搜索一般都需要使用一个估价函数来评估当前状态被处理的优先级，由算法 5 可知，函数 $\text{FINDALLMST}()$ 每次最多产生 $n - 2$ 次递归调用，而每次的边替换操作都会选择一个权值相等且概率最大的边进行替换，由定理 3.4.1 可知，递归调用后求得的最小生成树的边概率乘积必然比当前层的最小生成树 T 的边概率乘积小，因此，我们只要使用边替换操作后求

的最小生成树的边概率乘积作为估价值即可。算法 6 为使用启发式搜索的TOP-K最优生成树算法，*PriorQue*按照传入的第一个参数进行从大到小排序，*elem.T*为以*elem.E_T*为边集的最小生成树。

Algorithm 6 TOP-K最优生成树算法-启发式搜索优化。

```

1: procedure TOPKMST( $K, \mathcal{G}$ )
2:    $TopKMST \leftarrow \emptyset$ 
3:    $E_T \leftarrow \text{OST}(\mathcal{G})$ 
4:    $PriorQue \leftarrow \emptyset$ 
5:    $PriorQue.push(\prod_{e \in E_T} p(e), \emptyset, \emptyset, E_T)$   $\triangleright elem : [priority, E_I, E_O, E_T]$ 
6:   while  $PriorQue \neq \emptyset$  do
7:      $elem = PriorQue.pop()$ 
8:      $TopKMST = TopKMST \cup \{elem.E_I\}$ 
9:     if  $|TopKMST| \geq K$  then
10:      return  $TopKMST$ 
11:      $G \leftarrow (V, E \setminus elem.E_O, W)$ 
12:      $\text{SETNUMBER}(elem.T, 1, 1)$ 
13:      $S \leftarrow \text{FINDSWAPEDGE}(G, elem.T)$ 
14:     for  $i \leftarrow k$  to  $|V| - 2$  do
15:        $E_{I_i} \leftarrow elem.E_I \cup \{e_{k+1}, e_{k+2}, \dots, e_i\}$ 
16:        $E_{O_i} \leftarrow elem.E_O \cup \{e_{i+1}\}$ 
17:       if  $\exists e'((e_{i+1}, e') \in S)$  then
18:          $E_{T'} \leftarrow elem.E_T \cup \{e'\} \setminus \{e_{i+1}\}$ 
19:          $PriorQue.push(\prod_{e \in E_{T'}} P(e), E_{I_i}, E_{O_i}, E_{T'})$ 
20:   return  $TopKMST$ 

```

定理 3.4.2 算法 6的时间复杂度为 $O(Km \log N)$ ，其中 K 为需要查找的Top-K最优生成树， N 为不确定图 \mathcal{G} 中最小生成树的个数， m 为不确定图 \mathcal{G} 中的边数。

证明 显然*PriorQue*中元素最多为 N 个，如果*PriorQue*使用最大堆实现，则插入和删除的时间复杂度为 $\log N$ ，再由定理 3.3.3可知，算法 6的时间复杂度为 $O(Km \log N)$ 。

■

由定理 3.4.2可知，当 K 较小时，*PriorQue*中的元素也会较少，这样，算法 6的时间复杂度将接近 $O(m)$ ，所以算法 6非常适合 K 值较小的情况。

§3.5 实验结果

图 3.5为算法 5和算法 6的运行时间对比,该不确定图的顶点数为34,边数为196,最小生成树的个数为575624。由图 3.5可知,当 K 值小于270000时,算法 6优于算法 5。由于最小生成树的个数最多为575624,所以当 K 大于600000时,算法 6的运行时间不再增长。

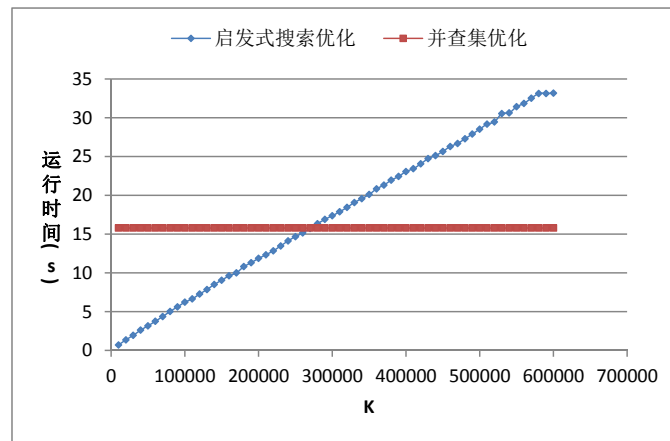


图 3.5 并查集优化和启发式搜索优化的实验对比。

§3.6 本章小结

第四章 不确定图中最小生成树的可靠性研究

§4.1 问题定义

定义 4.1.1 (可靠蕴含图) 设 G 为不确定图 $\mathcal{G} = (V, E, W, P)$ 的任意蕴含图, T_M^G 为图 G 的最小生成树, 若满足 $\sum_{e \in E_{T_M^G}} w(e) = W_M^G$, 则称 G 为不确定图 \mathcal{G} 的可靠蕴含图。

定义 4.1.2 (不确定图中最小生成树的可靠性) 设不确定图 $\mathcal{G} = (V, E, W, P)$, 令 $S = \{G | G \in \text{Imp}(\mathcal{G}) \wedge W_M^G = W_M^{\mathcal{G}}\}$ 为不确定图 \mathcal{G} 中所有可靠蕴含图的集合, 则我们定义不确定图 \mathcal{G} 中最小生成树的可靠性为:

$$r = \sum_{G \in R} \Pr(\mathcal{G} \Rightarrow G).$$

例 6 表 4.1 的第一列显示了图 2.1(a) 中的不确定图的蕴含图边集, 该不确定图一共有三颗最小生成树(其边集分别为 $\{e_1, e_2, e_4\}$, $\{e_2, e_3, e_4\}$ 和 $\{e_2, e_4, e_5\}$), 因此一共有7个可靠蕴含图(表 4.1 的前7行对应的蕴含图), 由于非可靠蕴含图出现的概率和为0.37378, 因此该不确定图的最小生成树的可靠性为 $1 - 0.37378 = 0.62622$ 。

§4.2 最小生成树的可靠性求解算法

定理 4.2.1 设不确定图 $\mathcal{G} = (V, E, W, P)$ 且不确定图 $\mathcal{G}' = (V, E', W, P)$, 其中 $E' = E \setminus \{e\}$, 则 $\sum_{G \in \text{Imp}(\mathcal{G})} \Pr(\mathcal{G} \Rightarrow G) = \sum_{G \in \text{Imp}(\mathcal{G}')} \Pr(\mathcal{G}' \Rightarrow G)$ 。

证 明 由公式 2.1 可知

$$\begin{aligned} \sum_{G \in \text{Imp}(\mathcal{G}) \wedge e \in E} \Pr(\mathcal{G} \Rightarrow G) &= \sum_{G \in \text{Imp}(\mathcal{G}) \wedge e \in E} p(e) \Pr(\mathcal{G}' \Rightarrow G \setminus e) \\ &= p(e) \sum_{G \in \text{Imp}(\mathcal{G}')} \Pr(\mathcal{G}' \Rightarrow G) \end{aligned}$$

同理

$$\sum_{G \in \text{Imp}(\mathcal{G}) \wedge e \notin E} \Pr(\mathcal{G} \Rightarrow G) = (1 - p(e)) \sum_{G \in \text{Imp}(\mathcal{G}')} \Pr(\mathcal{G}' \Rightarrow G)$$

表 4.1 图 2.1(a)的可靠性评估表.

蕴含图边集	包含的最小生成树	出现概率
$\{e_1, e_2, e_3, e_4, e_5\}$	$\{e_1, e_2, e_4\}, \{e_2, e_3, e_4\}, \{e_2, e_4, e_5\}$	0.31752
$\{e_1, e_2, e_3, e_4\}$	$\{e_1, e_2, e_4\}, \{e_2, e_3, e_4\}$	0.13608
$\{e_1, e_2, e_4, e_5\}$	$\{e_1, e_2, e_4\}, \{e_2, e_4, e_5\}$	0.07938
$\{e_2, e_3, e_4, e_5\}$	$\{e_2, e_3, e_4\}, \{e_2, e_4, e_5\}$	0.03528
$\{e_1, e_2, e_4\}$	$\{e_1, e_2, e_4\}$	0.03402
$\{e_2, e_3, e_4\}$	$\{e_2, e_3, e_4\}$	0.01512
$\{e_2, e_4, e_5\}$	$\{e_2, e_4, e_5\}$	0.00882
others	none	0.37378

因此

$$\begin{aligned}
 \sum_{G \in \text{Imp}(\mathcal{G})} Pr(\mathcal{G} \Rightarrow G) &= \sum_{G \in \text{Imp}(\mathcal{G}) \wedge e \in E} Pr(\mathcal{G} \Rightarrow G) + \sum_{G \in \text{Imp}(\mathcal{G}) \wedge e \notin E} Pr(\mathcal{G} \Rightarrow G) \\
 &= \sum_{G \in \text{Imp}(\mathcal{G}')} Pr(\mathcal{G}' \Rightarrow G)
 \end{aligned}$$

引理 4.2.1 设不确定图 $\mathcal{G} = (V, E, W, P)$, 则 $\sum_{G \in \text{Imp}(\mathcal{G})} Pr(\mathcal{G} \Rightarrow G) = 1$.

证明 设 $E = \{e_1, e_2, \dots, e_m\}$, 且

$$\begin{aligned}
 \mathcal{G}_0 &= \mathcal{G} \\
 \mathcal{G}_1 &= \mathcal{G}_0 \setminus \{e_1\} \\
 \mathcal{G}_2 &= \mathcal{G}_1 \setminus \{e_2\} \\
 &\dots\dots \\
 \mathcal{G}_m &= \mathcal{G}_{m-1} \setminus \{e_m\}
 \end{aligned}$$

则由定理 4.2.1可知

$$\sum_{G \in \text{Imp}(\mathcal{G})} Pr(\mathcal{G} \Rightarrow G) = \sum_{G \in \text{Imp}(\mathcal{G}_1)} Pr(\mathcal{G}_1 \Rightarrow G) = \dots = \sum_{G \in \text{Imp}(\mathcal{G}_m)} Pr(\mathcal{G}_m \Rightarrow G)$$

由于不确定图 \mathcal{G}_m 的边集为空, 所以 $\sum_{G \in \text{Imp}(\mathcal{G}_m)} Pr(\mathcal{G}_m \Rightarrow G) = 1$.

引理 4.2.2 设不确定图 $\mathcal{G} = (V, E, W, P)$, 且 $E_I, E_O \in E$, 则

$$\sum_{G \in Imp(\mathcal{G}) \wedge E_I \subseteq E_G \wedge E_O \cap E_G = \emptyset} Pr(\mathcal{G} \Rightarrow G) = \prod_{e \in E_I} p(e) \prod_{e \in E_O} (1 - p(e))$$

在引理 3.2.1中, 我们知道可以通过边集划分的方法来获取最小生成树。同样的, 这种划分也适用于不确定图中所有的蕴含图, 也就是说, 我们可以将图 3.2看成是不确定图中所有蕴含图集合的划分过程, 节点的左边表示蕴含图必然包含的边的集合, 即 E_I , 右边为蕴含图一定不包含的边的集合, 即 E_O 。当这个节点不能继续划分时, 如果边集 E_I 为不确定图中某颗最小生成树的边集, 那么划分到这个节点中所有的蕴含图即为可靠蕴含图。因此, 由引理 4.2.2可以得到改节点中所有可靠蕴含图的存在概率之和。算法 7 为不确定图中最小生成树的可靠性求解算法。

Algorithm 7 不确定图中最小生成树的可靠性求解算法

Comment Calls from $IMST(\emptyset, \emptyset, OST(\mathcal{G}), \mathcal{G})$

```

1: procedure  $IMST(E_I, E_O, E_T, \mathcal{G})$   $\triangleright \mathcal{G} = (V, E, W, P)$ 
2:    $r \leftarrow \prod_{e \in E_I} p(e) \times \prod_{e \in E_O} (1 - p(e))$ 
3:    $G \leftarrow (V, E \setminus E_O, W)$ 
4:    $SETNUMBER(T, 1, 1)$ 
5:    $S \leftarrow FINDSWAPEDGE(G, T)$ 
6:   for  $i \leftarrow k$  to  $|V| - 2$  do
7:      $E_{I_i} \leftarrow E_I \cup \{e_{k+1}, e_{k+2}, \dots, e_i\}$ 
8:      $E_{O_i} \leftarrow E_O \cup \{e_{i+1}\}$ 
9:     if  $\exists e'((e_{i+1}, e') \in S)$  then
10:       $E_{T'} \leftarrow E_T \cup \{e'\} \setminus \{e_{i+1}\}$ 
11:       $r \leftarrow r + IMST(E_{I_i}, E_{O_i}, E_{T'}, \mathcal{G})$ 
12:   return  $r$ 

```

§4.3 实验结果

实验 1 首先我们针对不确定图为完全图且所有边的权值都相等时的数据进行实验, 其所有的生成树都为最小生成树。根据 *Cayley* 定理^[35], 最小生成树的个数为 n^{n-1} , 其中 n 为不确定图中顶点的个数。

表 4.2 不确定图为完全图时算法的执行效率对比

G	N	运行时间(s)		时间比	可靠性
		$FERM$	$IERM$		
K_5	125	0.00100	0.00100	1.0	0.729287
K_6	1296	0.00500	0.00500	1.0	0.949277
K_7	16807	0.08100	0.07000	1.1	0.812293
K_8	262144	1.35800	1.06700	1.2	0.981199
K_9	4782969	29.2060	22.1580	1.3	0.990147
K_{10}	100000000	671.767	521.091	1.3	0.990147

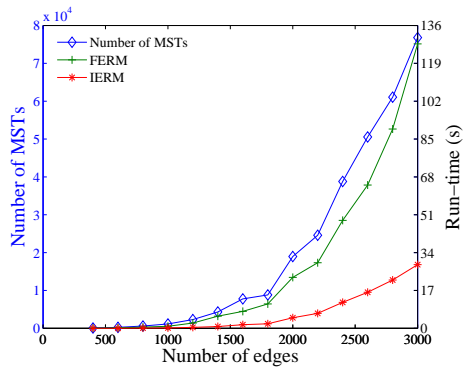
表 4.2为算法 $FERM$ 和 $IERM$ 的运行结果，其中算法 $FERM$ 没有采用并查集优化。表中的 N 表示不确定图中最小生成树的数量， K_n 表示该不确定图是一个顶点数为 n 的完全图，时间比即为算法 $FERM$ 与算法 $IERM$ 的执行时间之比，从表中可以看出，随着图规模的增大，时间比也越来越大，这是因为算法 $IERM$ 的程序时间相对复杂，导致常数级处理较多，所以当数据规模较少时效率不容易体现出来。

实验 2 为了验证顶点数和边数对算法的影响，我们将边的权值控制在1到100之间。图 4.1为算法 $FERM$ 和 $IERM$ 在不同顶点数下的时间对比图。

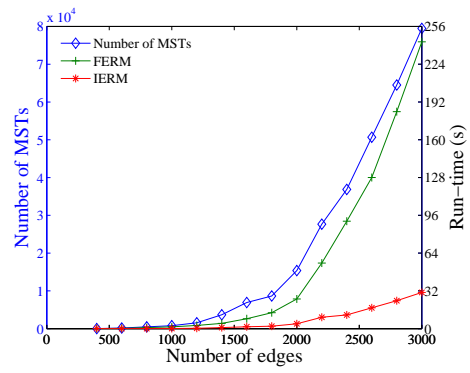
从图 4.1可以得出算法 $IERM$ 的时间复杂度与顶点数无关，这也验证了我们在第三章对并查集优化算法的理论分析。从图中还可以看出，这两个算法的运行时间都会随着边数的增多而增大。

§4.4 本章小结

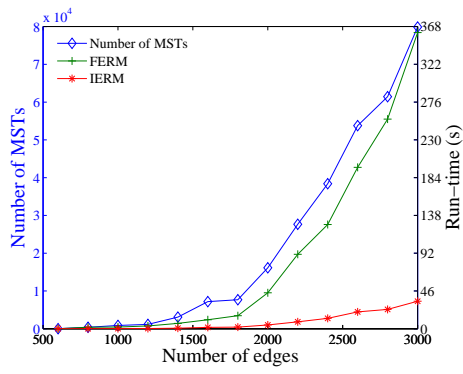
本章提出了在不确定图中最小生成树的可靠性问题，并设计了基于边集划分的可靠性求解算法，为了避免该算法在每次搜索产生分支后都要重新求最小生成树，使用并查集的思想在原有的最小生成树上进行边的删除和替换来更快地获得新的最小生成树。实验结果验证了算法的高效性和理论分析的正确性。



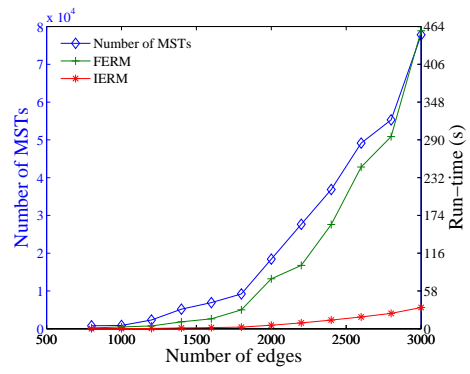
(a) $n = 100$



(b) $n = 200$



(c) $n = 300$



(d) $n = 400$

图 4.1 顶点数变化时算法的执行效率对比

第五章 总结与展望

参考文献

- [1] JENSEN L J, KUHN M, STARK M, et al. (2009) String 8-a global view on proteins and their functional interactions in 630 organisms. *Nucleic acids research*, **37**, 412–416.
- [2] CEOL A, ARYAMONTRI A C, LICATA L, et al. (2009) Mint, the molecular interaction database: 2009 update. *Nucleic acids research*, **38**, 532–541.
- [3] 李建中, 李金宝, 石胜飞. (2003) 传感器网络及其数据管理的概念、问题与进展. *软件学报*, **14**, 1717–1727.
- [4] SOHRABY K, MINOLI D, ZNATI T. (2007) *Wireless sensor networks: technology, protocols, and applications*. John Wiley & Sons.
- [5] BADER J S, CHAUDHURI A, ROTHBERG J M, et al. (2003) Gaining confidence in high-throughput protein interaction networks. *Nature biotechnology*, **22**, 78–85.
- [6] ASTHANA S, KING O D, GIBBONS F D, et al. (2004) Predicting protein complex membership using probabilistic network reliability. *Genome research*, **14**, 1170–1175.
- [7] SUTHRAM S, SHLOMI T, RUPPIN E, et al. (2006) A direct comparison of protein interaction confidence assignment schemes. *BMC bioinformatics*, **7**, 360.
- [8] JIANG R, TU Z, CHEN T, et al. (2006) Network motif identification in stochastic networks. *Proceedings of the National Academy of Sciences*, **103**, 9404–9409.
- [9] GHOSH J, NGO H Q, YOON S, et al. (2007) On a routing problem within probabilistic graphs and its application to intermittently connected networks. pp. 1721–1729. IEEE.
- [10] LIU Y, HE Y, LI M, et al. (2013) Does wireless sensor network scale? a measurement study on greenorbs. *Parallel and Distributed Systems, IEEE Transactions on*, **24**, 1983–1993.

- [11] MURUGANATHAN S D, MA D C F, BHASIN R I, et al. (2005) A centralized energy-efficient routing protocol for wireless sensor networks. *Communications Magazine, IEEE*, **43**, 8–13.
- [12] COOK D J, HOLDER L B. (2006) *Mining graph data*. John Wiley & Sons.
- [13] AGGARWAL C C, WANG H. (2010) *Managing and mining graph data*. Springer.
- [14] BORGELT C, BERTHOLD M R. (2002) Mining molecular fragments: Finding relevant substructures of molecules. pp. 51–58. IEEE.
- [15] YAN X, HAN J. (2002) gspan: Graph-based substructure pattern mining. pp. 721–724. IEEE.
- [16] KURAMOUCHI M, KARYPIS G. (2001) Frequent subgraph discovery. pp. 313–320. IEEE.
- [17] HUAN J, WANG W, PRINS J, et al. (2004) Spin: mining maximal frequent subgraphs from graph databases. pp. 581–586. ACM.
- [18] ZOU Z, LI J, GAO H, et al. (2009) Frequent subgraph pattern mining on uncertain graph data. pp. 583–592. ACM.
- [19] 邹兆年, 李建中, 高宏等. (2009) 从不确定图中挖掘频繁子图模式. *软件学报*, **20**, 2965–2976.
- [20] ZOU Z, LI J, GAO H, et al. (2010) Mining frequent subgraph patterns from uncertain graph data. *Knowledge and Data Engineering, IEEE Transactions on*, **22**, 1203–1218.
- [21] ZOU Z, GAO H, LI J. (2010) Discovering frequent subgraphs over uncertain graph databases under probabilistic semantics. pp. 633–642. ACM.
- [22] LI J, ZOU Z, GAO H. (2012) Mining frequent subgraphs over uncertain graph databases under probabilistic semantics. *The VLDB Journal*, **21**, 753–777.
- [23] ZOU Z, LI J, GAO H, et al. (2010) Finding top-k maximal cliques in an uncertain graph. pp. 649–652. IEEE.

- [24] 邹兆年, 朱镭. (2013) 大规模不确定图上的top-k极大团挖掘算法. 计算机学报, **36**, 2146–2155.
- [25] 张应龙, 李翠平, 陈红等. (2011) 不确定图上的knn 查询处理. 计算机研究与发展, **48**, 1850–1858.
- [26] 张旭, 何向南, 金澈清等(2011) 面向不确定图的k 最近邻查询. 计算机研究与发展, **48**, 1871–1878.
- [27] POTAMIAS M, BONCHI F, GIONIS A, et al. (2010) K-nearest neighbors in uncertain graphs. *Proceedings of the VLDB Endowment*, **3**, 997–1008.
- [28] JE H G, WIDOM J. (2002) Simrank: a measure of structural-context similarity. pp. 538–543. ACM.
- [29] 李俊辉(2013) 基于不确定图的层次聚类算法研究. 中国管理信息化 , **?**, 79–80.
- [30] 丁悦(2012) 不确定图聚类分析研究. Master's thesis. 西北农林科技大学.
- [31] YUAN Y, CHEN L, WANG G. (2010) Efficiently answering probability threshold-based shortest path queries over uncertain graphs. pp. 155–170. Springer.
- [32] 李鸣鹏, 邹兆年, 高宏等. (2012) 不确定图上期望最短距离的计算. 计算机研究与发展, **49**, 2208–2220.
- [33] 蔡伟, 张柏礼, 吕建华等. (2012) 不确定图最可靠最大流算法研究. 计算机学报, **35**, 2371–2380.
- [34] 张柏礼, 吕建华, 生衍等. (2014) 一种不确定图最可靠最大流问题的解决方案. 计算机学报, **10**, 2084–2095.
- [35] KELLY P J. (1957) A congruence theorem for trees. *Pacific J. Math*, **7**, 961–968.

致 谢

本研究及学位论文是在我的导师文中华教授的亲切关怀和悉心指导下完成的。他严肃的科学态度，严谨的治学精神，精益求精的工作作风，深深地感染和激励着我。从课题的选择到论文的最终完成，文老师都始终给予我细心的指导和不懈的支持。三年来，文老师不仅在学业上给我以精心指导，同时还在思想、生活上给我以无微不至的关怀，在此谨向文老师致以诚挚的谢意和崇高的敬意。

感谢信息工程学院的各位老师，从本科入学到研究生毕业的七年里，无私地把知识传授给我们，并教我们求知、做人的真理。

感谢我的母校，给了我良好的学习环境。在这个美丽的校园里，我结识了许多的良师益友。

感谢在一起愉快的度过研究生生活的工科楼506各位同门，正是由于你们的热情和帮助，才使得我很快的融入这个研究环境。你们也是我在三年的科研道路上最坚实的伙伴！

感谢生我养我，含辛茹苦的父母，是你们为我的学习创造了条件，并在我的身后默默的支持着我。谢谢你们！

在论文即将完成之际，我的心情无法平静，从开始进入课题到论文的顺利完成，有多少可敬的师长、同学给了我无言的帮助，在这里请接受我诚挚的谢意！

附录A (攻读硕士学位期间发表的论文)

1. 唐杰, 文中华, 汪泉等. 不确定可逆规划的强循环规划解. 计算机研究与发展, 2013, 50(9): 1970-1980.
2. 唐杰, 文中华, 黄海平等. 不确定规划中一种观察信息高效约简算法. 计算机工程, 2013, 39(12): 162-166.
3. Jie Tang, Yuansheng Liu, Zhonghua Wen. Reliability evaluation of the minimum spanning tree of uncertain graphs. Journal of Computers(已录用).
4. Yuansheng Liu, Jie Tang, Tao Xie. Cryptanalyzing a RGB image encryption algorithm based on DNA encoding and chaos map. Optics & Laser Technology, 2014, 60: 111-115.
5. Xie Tao, Yuansheng Liu, Jie Tang. Breaking a novel image fusion encryption algorithm based on DNA sequence operation and hyper-chaotic system. Optik-International Journal for Light and Electron Optics, 2014, 125(24): 7166-7169.

附录B (攻读硕士学位期间参与的科研项目)

1. 快速更新不确定规划解的算法研究, 湖南省研究生创新项目(项目编号: CX2014B276), 2014 - 2015年。主持人。
2. 基于模型检测的不确定规划的状态可达性及其应用研究, 国家自然科学基金项目(项目编号: 61070232), 2011 - 2013年。参与人。
3. 不确定规划中的观察信息约简方法及其应用研究, 国家自然科学基金项目(项目编号: 61272295), 2013 - 2016年。参与人。

附录C（攻读硕士学位期间获奖情况）

1. 2012 - 2013年湘潭大学三好学生。
2. 2012 - 2013年国家奖学金。
3. 2014年纪念毛泽东同志诞辰120周年环校长跑二等奖。