

1最大化延迟的最优性2数据结构（又称数据库破解）

3

4 抽象的

5本文研究在事先不知道 $r \leq n$ 的情况下,如何以在线方式（即,仅在前一个查询的结果准备好后才给出下一个查询）最小化回答 n 个元素上的 r 个查询的总成本。传统的索引首先在 n 个元素上构建完整的索引,然后再回答查询,这可能并不合适,因为索引的构建时间通常为 $\Omega(n \log n)$ 会成为性能瓶颈。相反,对于许多问题,对于每个 $r \in [1, n]$, r 个查询的总成本都有一个下限 $\Omega(n \log(1+r))$ 。匹配这个下限是延迟数据结构(DDS)的主要目标,在系统社区中也称为数据库破解。对于广泛的一类问题,我们提出了通用的归约¹³,将传统索引转换为与 r 的长范围的下限匹配的 DDS 算法。对于可分解问题,如果可以在 $O(n \log n)$ 时间内构建数据结构并具有 $15 Q(n)$ 查询搜索时间,则我们的简化方法可产生一个在 $O(n \log(1+r))$ 时间内运行的算法,其中对于所有 $n \log n \log n$,上限在温和约束下渐近最佳。 $\mathcal{A} \leq Q(n)$, $Q(n)$

16

17特别地,如果 $Q(n) = O(\log n)$, 则 $O(n \log(1+r))$ 时间保证扩展到所有 $r \leq n$, ¹⁸我们可以轻松地用它解决大量 DDS 问题。我们的结果可以推广到¹⁹一类“频谱可索引问题”,它包含可分解问题类。

20 2012 ACM 主题分类 计算理论 → 21 数据管理的数据结构和算法

22 关键词和短语 延迟数据结构、数据库破解、数据结构

23 数字对象标识符 10.4230/LIPIcs...

24 1 简介

25 传统索引首先在整个数据集上创建索引,然后才开始回答²⁶个查询。例如,在 $O(n \log n)$ 时间内在由 n 个实值组成的（未排序）集合 S 上构建一棵二叉搜索树 (BST) 后,我们可以使用该树在 $O(\log n)$ 时间内回答每个前导查询^{1 28}。然而,当数据集 S 仅被搜索少数几次²⁹时,此范例就不足为奇了。在极端情况下,如果只需要回答一个查询,那么最好的³⁰方法是完全扫描 S , 这只需要 $O(n)$ 时间。更一般地,如果我们要以在线方式回答 r 个查询- 即在输出前一个查询的结果后给出下一个查询- 则可能支付总成本 $O(n \log(1+r))$ [¹⁹]。

³³ 当 $r \ll n$ （例如 $r = \text{polylog } n$ 或 $2 \sqrt{\log n}$ ）时,成本 $O(n \log(1+r))$ 突破了³⁴ $\Omega(n \log n)$ 的障碍,这表明排序是不必要的。

³⁵ 上述情况发生在许多大数据应用中,这些应用收集了大量数据³⁶,但很少查询。这种现象促使系统社区开展了一项名为数据库破解的研究³⁷; 参见 [12、13、16–18、23、29、31、32] ³⁸及其参考文献。在这些环境中,问题规模 n 非常大,以至于即使是排序也被认为是昂贵的,应该尽可能避免。此外,无法预测需要支持多少个查询（即前面提到的整数 r ）。数据库破解⁴²不会立即构建完整的索引,而是在每次查询期间仅执行构建的计算部分。如果 r 最终

40

¹ 给定一个搜索值 q , 前趋查询将返回 S 中不超过 q 的最大元素。

XX:2 最大化延迟数据结构的最优性（又称数据库破解）

43达到某个阈值时,索引将完全创建,但更有可能的是 r 44会停止在某个明显低于该阈值的值。破解数据库45的挑战在于确保“平滑度”:随着 r 的增长,到目前为止所有 r 查询的总成本应该尽可能缓慢地增加。

在理论领域,数据库破解的核心数据结构问题早在 1988 年就已由 Karp, Motwani 和 Raghavan [19] 以延迟数据结构(DDS)的名义进行了研究。对于一些问题,他们解释了如何在事先不知道 r 的情况下回答 $r \in [1, n]$ 50 个对 n 个元素的查询,总成本为 $O(n \log(1 + r))$ 。对于每个 $r \in [1, n]$,他们证明了这些问题的匹配下限为 $\Omega(n \log(1 + r))$ 。

52通过简化,相同的下限已被证明适用于许多其他 DDS 53问题。

54 本工作探讨了以下主题:如何设计一个通用的归约方法,给定55 个(传统)数据结构,可以自动将其转换为有效的 DDS 算法?

56这种减少可能会显著简化 DDS 算法的设计,并为传统索引和 DDS 之间的复杂联系提供新的启示。

58 数学约定。我们使用 \mathbb{N}^+ 表示正整数集。对于任何整数 $x \geq 1$,让 $[x]$ 表示集合 $\{1, 2, \dots, x\}$ 。默认情况下,每个对数的底数都是 2。定义 $\text{Log}(x) = \log(1 + x)$,对于任何 $x \geq 1$ 。

61 1.1 问题定义

62本小节将形式化要研究的问题。设 S 称为数据集, 63是从域 D 中抽取的一组 n 个元素。设 Q 是一个(可能无限的)集合,其中的元素称为谓词。给定谓词 $q \in Q$,对 S 发出的查询将返回一个答案,表示为 $\text{Ans}_q(S)$ 。我们认为,对于任何 $q \in Q$,答案 $\text{Ans}_q(S)$ 都可以用 $O(1)$ 个词表示,并且可以在 $O(n)$ 时间内计算出来(这本质上意味着可以使用穷举扫描等强力算法来回答查询)。

68 延迟数据结构化 (DDS)。我们现在将 DDS 问题形式化。最初,算法 A 提供数组中的数据集合 S ,其中元素任意排列。对手首先为第一个查询选择谓词 q_1 ,并从 A 中请求答案 $\text{Ans}_{q_1}(S)$ 。迭代地,对于每个 $i \geq 2$,在获得第 $(i - 1)$ 个查询的答案72之后,对手要么决定终止整个过程,要么为下一个查询选择谓词 q_i 并从 A 中请求 $\text{Ans}_{q_i}(S)$ 。对手被允许观察 A 的执行情况,因此能够为 A 选择“坏” q_i 。

71

73

75 如果对于每 $6 r \leq t$,前 r 个查询的处理总成本最多为 $\text{Time}(n, r)$,则算法 A 可以保证 t 个查询的运行时间为 $\text{Time}(n, r)$ 。我们将集中讨论77 $t \leq n$,因为这是数据库破解的重要场景。

78最优性。虽然上述设置对于 DDS 来说是“标准”,但就数据库79破解而言,从另一个新角度进行调查是有意义的。

80由于数据库破解主要在数据集接收相对较少的查询的情况下有用,因此必须设计在查询数量较少时特别有效的 DDS 算法。

83 为了将“特别有效”的概念形式化,我们利用这样一个事实:对于许多84 个DDS 问题(包括本文中考虑的问题),存在硬度障碍85,规定在相关计算模型下,对于每个 $r \in [1, n]$, $\text{Time}(n, r) = \Omega(n \text{Log } r)$ 。

86受此启发,我们说如果算法 A 的运行时间对所有 $r \leq \text{LogStreak}(n)$ 都满足 $\text{Time}(n, r) = O(n \text{Log } r)$,则算法 A 保证了 $\text{LogStreak}(n)$ 的

87 $\log(r)$ -streak。

88 最差的 $\log(r)$ -streak 保证是 $\text{LogStreak}(n) = O(1)$ 这是微不足道的,因为
 89任何查询都可以在 $O(n)$ 时间内得到回答。理想情况下,我们希望 $\text{LogStreak}(n) = n$,
 90,但这并不总是可能的,正如本文后面所讨论的。然而,对于实际使用来说,
 91 足以让算法确保 $\text{LogStreak}(n) = \Omega(n)$) 其中某个常数 > 0
 92 因为 $\Omega(n)$) 查询可能已经太多了,以至于数据库破解不再具有吸引力。

93类问题。上述定义框架可以具体化为各种

94 个问题实例,它们在数据域 D 、谓词域 Q 和结构上有所不同

95 T。接下来,我们介绍两个与我们的讨论相关的问题类别:

96 ■ 如果对于任何不相交集 $S_1, S_2 \subseteq D$ 和任何谓词,则问题实例是可分解的
 97 $q \in Q$,可以从 $\text{Ans}_q(S_1)$ 和 $\text{Ans}_q(S_2)$ 中导出 $\text{Ans}_q(S_1 \cup S_2)$
 98 时间。

99我们定义一个问题实例是 $(B(n), Q(n))$ 谱可索引的,如果它具有

100 对于每个数据集 $S \subseteq D$,都满足以下属性:对于每个整数 $s \in [|S|]$,都有可能
 101 在 $O(|S| \cdot B(s))$ 时间内在 S 上构建一个数据结构,可以回答任何查询
 102 这 $(\frac{|S|}{s} \cdot Q(s))$ 时间。术语“频谱可索引”被选择来反映构建
 103 就函数 $B(n)$ 和 $Q(n)$ 而言,“良好”的索引结构
 104 参数 s 的整个频谱。

105 对于上述定义,有两点很重要:

106 ■ $(B(n), Q(n))$ 谱可索引性意味着我们可以在任何
 107 数据集 $S \subseteq D$ 在 $O(|S| \cdot B(|S|))$ 时间内回答查询,时间为 $O(Q(|S|))$ (对于这个
 108 目的,只需设置 $s = |S|$)。
 109 ■ 考虑具有以下属性的任何可分解问题实例:对于任何数据集
 110 $S \subseteq D$,我们可以在 $O(|S| \cdot B(|S|))$ 时间内构建一个数据结构 T ,以回答查询
 111 $O(Q(|S|))$ 时间。那么,问题实例必须是 $(B(n), Q(n))$ 谱可索引的。
 112 为了了解原因,给定一个整数 $s \in [|S|]$,将 S 任意分为 $m = \lceil |S|/s \rceil$ 不相交
 113 子集 S_1, S_2, \dots, S_m ,其中除 S_m 外,所有子集的大小均为 s 。对于每个 $i \in [m]$,创建
 114 S_i 中的结构 $T(S_i)$ $| \cdot B(s)$ 时间;创建所有 m 个结构的总时间
 115 是 $O(m \cdot s \cdot B(s)) = O(|S| \cdot B(s))$ 。要回答查询 q ,只需搜索每个 $T(S_i)$ 即可
 116 在 $O(Q(s))$ 时间内获得 $\text{Ans}_q(S_i)$,然后组合 $\text{Ans}_q(S_1), \text{Ans}_q(S_2), \dots, \text{Ans}_q(S_m)$
 117 使用 $O(m)$ 时间将查询结果转化为 $\text{Ans}_q(S)$ 。因此,总查询时间为 $O(m \cdot Q(s))$ 。

118 1.2 相关工作

119 Motwani 和 Raghavan 在一篇会议论文 [24] 中提出了 DDS 的概念,

120被合并到与 Karp 合著的期刊文章 [19] 中。他们 [19] 设计

针对以下 DDS 问题的121 种算法:

- 122 ■ 前置搜索,其中 S 由 n 个实数组成,每个查询被赋予一个任意
 123 值 q 并返回 S 中 q 的前身。
- 124 ■ 半平面包含,其中 S 由 R 中的 n 个半平面组成 2 , 并且每个查询都给出
 125 任意点 $q \in R$ 2 并返回 q 是否被 S 中的所有半平面覆盖。
- 126 ■ 凸包包含,其中 S 由 R 中的 n 个点组成 2 , 每个查询都会被赋予
 127 任意点 $q \in R$ 2 并返回 q 是否被 S 的凸包覆盖。
- 128 ■ 二维线性规划,其中 S 由 R 中的 n 个半平面组成 2 , 并且每个查询都给出
 129 二维向量 u ,并返回所有 n 个半平面交点中的点 p
 130 最大化点积 $u \cdot p$ 。
- 131 ■ 正交范围计数,其中 S 由 R 中的 n 个点组成 d 维数为 d
 132 是一个固定常数,查询是给定的任意 d 矩形 q 即

XX:4 最大化延迟数据结构的最优性（又称数据库破解）

133 轴平行盒形式为 $[x_1, y_1] \times [x_2, y_2]$ 被 q 覆盖的 S 中的点。 $\dots \times [x_d, y_d]$ 返回

134

135 对于前四个问题, Karp, Motwani 和 Raghavan 提出了实现

136 对于所有 $r \leq n$, $\text{Time}(n, r) = O(n \log r)$ 。对于正交范围计数, 他们提出了两个

137 算法, 其中第一个确保 $\text{Time}(n, r) = O(n \log d)$ $r)$ 对所有 $r \leq n$, 而

138 另一个确保 $\text{Time}(n, r) = O(n \log n + n \log d - 1)$ $r)$, 其中 $r \leq n$ 。对于所有这些问题,

139 他们证明了, 在比较模型和/或代数模型下, 运行时间

任何算法的 140 都必须满足对于每个 $r \in [1, n]$, $\text{Time}(n, r) = \Omega(n \log r)$ 。

141 Aggarwal 和 Raghavan [1] 提出了一种 DDS 算法, 其时间 $(n, r) = O(n \log r)$

142 对于所有 $r \leq n$ 的最近邻搜索, 其中 S 由 R 中的 n 个点组成 2 , 并且查询是

143 给定任意点 $q \in R$ 2 并返回 S 中距离 q 最近的点。运行时间为

144 对于所有 $r \leq n$, 在代数模型下都是最优的。

145 关于距离中值的 DDS 可以讲述一个“成功的故事”。在问题的离线

146 版本中, 我们给定一个 (未排序的) 数组 A 中的 n 个实数值集合 S 。对于任何 $1 \leq x \leq y \leq n$,

147 设 $A[x : y]$ 表示元素集 $\{A[x], A[x+1], \dots, A[y]\}$ 。此外, 给定 r

148 个整数对 $(x_1, y_1), (x_2, y_2), \dots, (x_r, y_r)$, 其中对于每个 $i \in [r]$, $1 \leq x_i \leq y_i \leq n$ 。目标

149 是找到集合 $A[x_i : y_i]$] 且其中所有 $i \in [r]$ 。在 [15] 中, Har-Peled 和 Muthukrishnan

150 解释了如何在 $O(n \log r + r \log n \cdot \log r)$ 时间内解决这个问题, 并证明了一个更低的

151 在比较模型下, 当 $r \leq n$ 时, $\Omega(n \log r)$ 的界。在 [11] 中 (另见 [5]), Gfeller

152 和 Sanders 考虑了该问题的 DDS 版本, 其中 S 如前所述, 并且

153 查询给定一个任意对 (x, y) , 其中 $1 \leq x \leq y \leq n$, 并返回 $A[x : y]$ 的中位数。

154 他们设计了一种算法, 当 $r \leq n$ 时, $\text{Time}(n, r) = O(n \log r)$ 。很容易

155 看到任何 DDS 算法都可用于解决 $\text{Time}(n, r)$ 中的离线问题

156 次。因此, Gfeller 和 Sanders 的算法改进了 [15] 的离线解决方案, 并且

157 在比较模型下是最优的 (对于 DDS)。

158 与我们的工作更相关的是 [8], 其中 Ching, Mehlhorn 和 Smid 考虑

159 动态 DDS 问题, 除了查询之外, 还需要算法来支持

160 次更新数据集 S 。在另一个方向上, Barbay 等人 [2] 研究了 DDS

161 版本的前身搜索, 但使用更细粒度的参数分析了他们的算法

162 称为“间隙” (而不是仅使用 n 和 r)。这个方向也扩展到

163 动态 DDS; 参见最近的著作 [26, 27]。

164 如上所述, DDS 已在以下领域得到系统界的广泛研究:

165 个名称数据库破解。重点是设计有效的启发式方法来加速

166 在各种实际场景中遇到的查询工作负载, 而不是建立强大的

167 个理论保证。感兴趣的读者可以参考 [12, 13, 16–

168 18, 23, 31, 32] 作为进入文献的切入点。

169 1.3 我们的结果

170 我们的主要结果是具有以下保证的通用减少:

171 ►定理 1. 假设 $B(n)$ 和 $Q(n)$ 是非减函数, 且 $B(n) =$

172 $O(\log n)$ 且 $Q(n) = O(n^{1-\epsilon})$ 其中 $\epsilon > 0$ 是任意小的常数。每个 $(B(n), Q(n))$

173 频谱可索引问题允许使用 DDS 算法

$$174 \quad \text{LogStreak}(n) = \min \left\{ n, \frac{c \cdot n \log n}{Q(n)} \right\} \quad (1)$$

175 对于任意大的常数 c , 即算法实现 $\text{Time}(n, r) = O(n \cdot \log r)$

176 对于所有 $r \leq \min \left\{ n, \frac{c \cdot n \log n}{Q(n)} \right\}$ 。

177 在温和的约束条件下,如我们在第3节中论证的那样,(1)中的条纹界限是最好的。作为一个重要的特例,当 $Q(n) = O(\log n)$ 时,由我们的归约产生的 DDS 算法实现 $\text{LogStreak}(n) = n$,即,对于所有 $r \leq n$, $\text{Time}(n, r) = O(n \cdot \log r)$ 。对于许多对数据库系统具有重要意义的可分解问题,构造时间为 $O(n \log n)$,搜索时间为 $O(\log n)$ 的数据结构已经已知(例如,前驱搜索和最近邻搜索)。由于这类问题必须是 $(\log n, \log n)$ 谱可索引的(如第1.1节所述),定理1立即为所有这些问题的 DDS 版本产生了优秀的解决方案,这些解决方案通常在比较/代数模型下是最优的。第4.1节将给出这些问题的一部分。

186 定理1为数据库破解带来了好消息:即使查询时间很慢的数据结构也可用于破解!例如,对于2D空间中的正交范围计数(在1.2节中定义),kd树的构建时间为 $O(n \log n)$,可在 $Q(n) = O(\sqrt{n})$ 时间内回答查询。定理1表明,该结构可用于在 $O(n \log r)$ 时间内回答任何 $r = O(\sqrt{n \log n})$ 查询,这对于现实中的数据库破解来说可能绰绰有余。这很好地体现了我们研究“对小 r 特别有效的 DDS 算法”的动机(在1.1节中说明)。

193 定理1对DDS算法的设计具有指导意义。我们应该研究底层问题实际上的频谱可索引性。在这个方向上的探索会变得相当有趣,正如我们将在第4.2节中展示的半平面包含、凸包包含、范围中位数和二维线性规划(参见第1.2节了解它们的定义)。我们可以证明,对于适当选择的函数 $Q(n)$,这些问题中的每一个都是 $(\log n, Q(n))$ 频谱可索引的,然后利用定理1获得一个算法,该算法对所有 $r \leq n$ 都使用 $\text{Time}(n, r) = O(n \log r)$ 来解决它。

200 构建需要 $\omega(n \log n)$ 时间的结构,或者说 $B(n) = 201 \omega(\log n)$,会发生什么情况?第5节将展示,在温和的约束下,没有通用的归约可以使用这样的结构来获得 $\text{LogStreak}(n) = \omega(1)$ 的 DDS 算法。换句话说,这些算法只能在 $r = O(1)$ 的简单场景中实现 $\text{Time}(n, r) = O(n \log r)$ 。

204 尽管如此,如果人们接受具有“对于所有 $r \leq n$, $\text{Time}(n, r) \leq n \log O(1)$ ”形式保证的算法,那么只要 $B(n)$ 和 $Q(n)$ 是 $r \log O(1) n$,我们就可以扩展定理1所依据的归约来产生这样的算法,正如第5节将要讨论的那样。

207 **2 热身:前任搜索**

208 前任搜索是 DDS 和209数据库破解社区最关注的问题。本节将回顾[19]中开发的两种方法,用于实现此问题的 $\text{Time}(n, r) = O(n \log r)$ 。这些方法构成了几乎所有现有 DDS 算法的基础。

212 自下而上。回想一下,数据集 S 由 n 个实数组成。我们假设 $w \log$, $2^{13} n$ 是2的幂。在任何时候,集合 S 都被任意划分为不相交的子集(称为运行),每个子集的大小相同, $s = 2^i$,其中 $i \geq 0$ 。每个运行都经过排序并存储在数组中。最初, $s = 1$,即运行包含 S 的单个元素。随着时间的推移,运行大小 s 单调增加。每当 s 需要时从2开始,对于某个值 $j > 1$,就会进行大修以构建新的运行。由于可以通过在 $O(2^j)$ 个大小为 2^{j-1} 的运行中合并 2^{j-1} 来获得大小为 2^j 的运行,因此大修可以在

214 $O(n \cdot (j-1))$ 的时间内完成219。因此,如果当前运行大小为 s ,则生成历史上所有220次运行的总成本为 $O(n \log s)$ 。

图 2 个单位

221 每次运行时只需执行二进制搜索即可回答(前一个)查询。
222 然而,为了控制成本,该算法确保当前大小 s 至少为

XX:6 最大化延迟数据结构的最优性（又称数据库破解）

223 在处理第 i 个 ($i \geq 1$) 个查询之前,先计算 $i \log i$ 。如果不满足条件,则调用大修224以将 s 增加到最接近的 2 的幂至少 $i \log i$ 。此后,查询 $\cdot \log(i \log i) = O(n/i)$ 时间。如果我们将所有 r 个查询都加起来,总和变为 $O(n) = O(n \log r)$ 。由于最终运行大小 s 为 227

225 需要花费 $O(n \log n)$ 对数) = $O(n \log n)$ 我们可以得出结论,该算法在 $O(n \log r)$ 时间内处理 r 个

226 查询。注意228这仅在 $r \log r \leq n$ (最大化运行大小为 n) 时成立。但是,当 r 达到 $229 \quad n / \log n$ 时,算法可以在 $O(n \log n) = O(n \log r)$ 时间内对整个 S 进行排序,并在 $O(\log n) = O(\log r)$ 时间内回答每个后续查询。因此,对于所有 $r \leq n$,该算法实现231

$\text{Time}(n, r) = O(n \log r)$ 。

232 自上而下。此方法模仿以下在 S 上构建二叉搜索树(BST) T 的策略:(i) 找到 S 的中位数,并在中位数处将 S 拆分为 S_1 和 S_2 ;(ii) 将中位数存储为根的键,然后在 S_1 (或 S_2) 上递归构建根的左 (或右) 子树。该算法不是直接进行完整构建,而是在查询处理期间“按需”构建 T 。

237 一开始,只有根存在,并且它处于未扩展模式。通常,未扩展节点 u 还没有子节点,但是与应该存储在其子树中的元素的子集 $S_u \subseteq S$ 相关联。查询的回答方式与普通 BST 中相同,即遍历 T 的根到叶路径 π 。主要区别在于,当搜索到 π 上的未扩展节点 u 时,算法必须先扩展它。扩展 u 意味着为 u 创建两个子节点,在中位数 (u 的关键词) 处拆分 S_u ,在中位数处将 S_u 分成两部分,并将每部分与一个子节点相关联。此后, u 变为扩展状态,其子节点处于未扩展模式。扩展需要 $O(|S_u|)$ 时间245 (查找集合的中位数需要线性时间[4])。

246 经过 r 次查询后,二叉搜索树 T 只被部分构建,因为在查询处理过程中,只有 r 条从根到叶的路径上的节点被扩展。前 $\log r$ 层

247 节点2 248的总扩展成本为 $O(n \log r)$ 。对于每条从根到叶的路径 π ,位于 $\log r$ 层的节点 u 249 的扩展成本为 $O(n/r)$,这决定了 π 上 u 的250 个后代节点的总扩展成本。因此,除了前 $\log r$ 层的节点外, T 中所有其他节点的总扩展成本为 $r \cdot O(n) = O(n)$ 。因此,该算法

251 $\text{Time}(n, r) = O(n \log r)$ 。

253 3 从数据结构到 DDS 算法的简化

254 第 3.1 节将自下而上的方法 (在上一节中进行了回顾) 扩展为通用归约,这可以产生具有较大 $\log(r)$ 条纹边界的 DDS 算法,前提是满足关键要求 线性可合并性 (即将定义)。虽然这个归约将被我们在定理1 下的最终归约 (在第 3.2 节中介绍) 所取代,但其讨论 (i) 加深了读者对该方法的强大功能和局限性的理解,并且 (ii) 阐明了在没有线性可合并性的情况下新想法的必要性。最后,第 3.3 节将建立一个硬度结果,以表明定理 1 中的条纹边界不再能得到显著改善。

262 3.1 第一个简化:推广自下而上的方法

263 本节将重点讨论可分解问题。对于任何数据集 $S \subseteq D$,我们假设存在一个数据结构 $T(S)$,它可以在 $O(Q(n))$ 时间内回答任何查询。此外,

² 节点的级别是从该节点到根的路径上的边的数量。

265 结构是线性可合并的,即,对于任何不相交的 $S_1, S_2 \subseteq D$, 266 $S_1 \cup S_2$ 上的结构可以在 $O(|S_1| + |S_2|)$ 时间内从 $T(S_1)$ 和 $T(S_2)$ 构造出来。请注意,这267意味着 $T(S)$ 可以在 $O(n \log n)$ 时间内构建。

268 ▶引理 2.对于一个可分解问题,存在一个线性可合并结构,其中 $\epsilon > 0$ 为常数,我们可以设计一个 DDS 算法
269 查询时间 $Q(n) = O(n^{1-\epsilon})$ 来
270 $\text{LogStreak}(n) = \min\{n, \frac{c \cdot n \log n}{Q(n)}\}$ 对于任意大的常数 c 。

271 证明。我们假设 $w \log_2 n$ 是 2 的幂。与自下而上的方法一样,在任何时候,我们将 S 任意划分为大小相同的运行 $s = 2^i$,其中 $i \geq 0$ 。对于
272 每个273运行,在其中的元素上构建一个结构。初始运行大小为 1。每当 s 增长到 2^j (对于某个值 $j > i$)时,就会进行一次大修以构建
大小为 2^j 的运行

274 从 2 开始通过 ∞
275 线性可合并性,我们可以通过在 $O(2^j)$ 个大小为 2^{j-1} 的运行中合并 2^{j-1} 的结构来构建大小为 2^j 的运行的结构。通过类似于第 2 节中的分析,如果当前运行大小为 s ,则生成所有运行的结构的总成本
276
277

历史上出现过的278个问题的答案是 $O(n \log s)$ 。

279 带有谓词 q 的查询通过搜索每个运行的结构来回答,然后将所有运行的答案组合成 $\text{Ans}_q(S)$ 。查询成本为 $O(n \cdot Q(s))$ 。我们要求
280 在回答第 i 个 ($i \geq 1$)查询之前,运行大小 s 必须满足 $\frac{n}{s}$

282
$$Q(s) / s \leq 1 / i. \tag{2}$$

283 如果这个要求不满足,我们将进行彻底检查,将 s 增加到满足(2)的最小284次方。这确保第 i 个查询的回答成本为 $O(n/i)$ 。因此,处理 r
个查询的总成本为 $O(n \log r)$ 。

286 剩下的就是限制大修的成本。由于 (2), 最终的运行规模为
287是 2 的最小幂,满足 $s \leq n$ (运行大小不能超过
288 n)和 $s/Q(s) \geq r$ (因为 (2))。

289
290 由于 $Q(s) = O(s^\alpha)$,其中 $\alpha > 0$ 为常数,我们知道 $Q(s) \leq \alpha / \alpha$ 。由于我们的目标是找到 s 的上限,我们要求 s 满足 $s / \alpha \geq r$ (比 $s/Q(s)$
291 $\geq r$ 更严格),或等效地 $s \geq (\alpha \cdot r)^{1/\alpha}$ 。因此,如果 $(\alpha \cdot r)^{1/\alpha} \leq n$,则我们可以声称 $s = O(r^{1/\alpha})$,因此 $\tilde{O}(\log s) = O(\log r)$,在
292 这种情况下,所有大修294总共需要 $O(n \log r)$ 时间。

293
295 如果 $(\alpha \cdot r)^{1/\alpha} > n$,则上述策略不起作用。但是,在这种情况下, $r > n^\alpha / \alpha$,即 r 已经是 n 的多项式。这激发了以下蛮力策略。当 r 达到 n^α / α 时 我们称之为捕捉点
296 我们只需在 $O(n \log n) = O(n \log r)$ 时间内在整个 S 上创建一个结构,并使用它来回答每个后续查询。捕捉点之后的所有查询的总成本为 $O(n \log n) = O(n \log r)$ 。因此,
297 我们获得了一种算法,该算法保证对于所有 $r \leq \min\{n, \frac{c \cdot n \log n}{Q(n)}\}$ $\text{Time}(n, r) = O(n \log r)$ 。

299 在 $O(Q(n))$ 时间内直到 $r = \min\{n, \frac{c \cdot n \log n}{Q(n)}\}$
301 $\frac{c \cdot n \log n}{Q(n)}$

302 上述简化主要依赖于结构 T 是线性可合并的事实。在这种情况下,运行次数将变为 $\Omega(n \cdot \log(2^i))$,否则,创建大小为 2 的所有
303 以削减 $\log r$ 因子的简化方法。 ∞ 大修的总成本将为 $\Omega(n \cdot \log^2 r)$ 。接下来,我们将介绍305另一种可

304

306 3.2 第二个归约:非线性可合并性

307 我们现在将删除第 3.1 节中的线性可合并性要求并建立定理 1。

308 回想一下,底层问题是 $(B(n), Q(n))$ 谱,可用 $B(n) = O(\log n)$ 来索引



XX:8 最大化延迟数据结构的最优性 (又称数据库破解)

309 并且 $Q(n) = O(n^{1-\epsilon})$ 其中某个常数 $\epsilon > 0$ 。目标是设计一个算法

310 对于所有 $r \leq \min\{n, \frac{c \cdot n \log n}{Q(n)}\}$, 其中 $c > 0$ 可以是任何常数。

311 假设 $w \log n$ 是 2 的幂。我们的算法以 epoch 为单位执行。在

第 i 个 ($i \geq 1$) 时期的开始, 我们设置

313 $s = 22^i$

314 并在 $O(n \cdot B(s))$ 时间内在 S 上创建一个结构 T (即 $(B(n), Q(n))$ 频谱可索引性所承诺的结构)。该结构使我们能够在 $O(n \cdot Q(s))$ 时间内

315 回答任何查询。

316 第 i 个时期结束于 $s/Q(s)$ 个查询。在该时期内得到回答之后, 这些查询

317 要求总成本

318 $\frac{s}{\text{问}} \cdot \frac{n}{s} \cdot \text{问题} = O(n)$ 。

319 从上文可以清楚看出, 第 i 个 epoch 的总计算时间为 $O(n \cdot B(s)) =$

320 $O(n \cdot 2^i)$ 。作为 $n \cdot 2^i$ 当增加 1 时, 其总成本翻倍, 回答 r 个查询的总成本为 h , 其中 h 是所需的 epoch 数。确切地说, h 的值是最小的

321 $O(n \cdot 2^h)$

322 整数 $h \geq 1$ 满足两个条件:

323 ■ C1: $2^h \leq n$ (s 的值不得超过 n);

324 ■ C2: $\sum_{h=1}^h \frac{2^{h-1}}{Q(2^{h-1})} \geq r$ (h 个周期可以回答的查询数量必须

325 少为 r)。

326 由于我们的目标是找到 h 的上限, 因此我们用更严格的条件替换条件 C2: 对于某些

327 $\frac{r}{Q(2^h)} \geq r$ 。因为 $Q(n) = O(n^{1-\epsilon})$, 已知 $Q(n) \leq \alpha \cdot n^{1-\epsilon}$ 常数 $\alpha > 0$ 。

由此, 我们进一步将 C2 修改为更为严格的不等式:

329 $\frac{2^{h-1}}{2^{h-1} \cdot (22^{h-1})^{1-\epsilon}} = \frac{(22^h)}{2^{h-1}} \geq r \Leftrightarrow 2^{h-1} \geq (\alpha \cdot r)^{1/\epsilon}$ (3)

330 令 H 表示满足 (3) 式的最小整数 $h \geq 1$ 。这意味着

331 $2^{H-1} < (\alpha \cdot r)^{1/\epsilon} \Leftrightarrow 2^{H-1} < (\alpha \cdot r)^{1/\epsilon}$ (4)

332 当 $2^{H-1} \leq n$, 上述论证保证了迭代次数 h 最多 $\leq n$, 满足条件 C1)。在这种情况下, 所有, 根据 (4) 可知其复杂度为 $O(n \log r)$ 。

333 H (不等式 2 个时期的总成本 $2^H \leq n$ 意味着 2^H

334 本为 $O(n \cdot 2^H) = O(n \cdot 2^H)$

335 如果 $2^H > n$, 上述论证就不成立。然而, 当这种情况发生时, 我们从 (4) 知道 $(\alpha \cdot r)^{1/\epsilon} > 2^H > n$, 从而导致 $r > n^{1/\epsilon} / \alpha$ 。一旦 r 达到

336 $n^{1/\epsilon} / \alpha$ 即捕捉点我们就会在 $O(n \log n) = O(n \log r)$ 时间内在整个 S 上创建一个结构 T , 并使用它回答 $O(Q(n))$ 时间内的每个

后续查询, 直到 $r = \min\{n, \frac{c \cdot n \log n}{Q(n)}\}$ 。捕捉点之后的 339 个查询需要的总成本为 $O(n \log n) = O(n \log r)$ 。因此我们有,

338 $\frac{c \cdot n \log n}{Q(n)}$

340 获得了一个算法, 保证对于所有 $r \leq \min\{n, \frac{c \cdot n \log n}{Q(n)}\}$,

341 完成定理 1 的证明。

342 3.3 条纹边界的紧密性

343 本节将解释为什么对于具有“合理”行为的约简算法来说, 条纹界限 $\frac{c \cdot n \log n}{Q(n)}$ 在定理 1 中渐进地

Ω(344 是最好的)。

³ 在实践中, 我们的算法可以通过使这个数字为 $s \cdot B(s)/Q(s)$ 而略有改进, 但这对于证明定理 1 来说并不是必要的

345 具有受限结构的可分解问题的黑盒约化。

346 为了证明难度结果,我们可以随意专门化问题类,我们通过仅考虑可分解问题来实现这一点。需要一种简化算法A 来处理任何
(B(n), Q(n))谱可索引的可分解问题。如第 1.1 节所述, 349这意味着可以在 $O(|S| \cdot B(|S|))$ 时间内在任何 $S \subseteq D$ 上构建一个数据结
构T,并在 $O(Q(|S|))$ 时间内回答S上的任何查询。

351 只要保留(B(n), Q(n))谱可索引性,我们就可以限制功能

352使得 A 难以解决。具体来说,T 只为 A 提供以下“服务”:给定谓词 $q \in Q$,A 可以创建结构 T (S ;给定谓词q ,A 可以使用

353 ■ T (S ;在A已经获得Ans_q(S₁) 和Ans_q(S₂)对任意子集 $S' \subseteq S$;

354 ■ $S' \subseteq S$ 之后,它可以将答案组合成Ans_q(S₁ ∪ S₂)找到答案Ans_q(S₁ ∪ S₂) 在S₁ ∪ S₂ 中;

355 ■ 子集S₁ 和S₂ 为不相交) 在恒

356 组合算法由 T 提供)。

357 组合算法由 T 提供)。

358尽管功能有限, T 仍然使得问题(B(n), Q(n))谱可索引359,因为该问题是可分解的;参见第 1.1 节中的解释。

360 到目前为止,还没有对 A 的行为施加任何限制,但现在我们准备这样做。要回答查询,算法A需要搜索数字 上的结构)。
该算法可以选择任何(它们不需要不相交)。此外, A还被允许时间。有了所有这些,算法

362 (包括零)S 的子集,比如T (S₁ ≥ 0 和任何S₂), T (S₂),..., T (S_n)

363 另一个子集S₂, ..., t检查

364 扫描 S,支付 $\Omega(|S|)$ 扫描

365必须确保

$$366 \quad \text{年代扫描} \cup \text{小号}_1 \cup \text{小号}_2 \cup \dots \cup S_{n_t} = S. \quad (5)$$

367上述约束是自然的,因为否则, S中至少有一个元素不存在于中。如果算法 A “敢” 无论如何返回查询答案,那么它一定已经获得了底层问题的某些特殊属性。在这项工作中,我们370关注的是

368 年代扫描 $\cup \text{小号}_1 \cup \text{小号}_2 \cup \dots \cup S_{n_t}$ 不知道问题特定属性的通用归纳。

371我们将满足上述要求的约简算法称为黑盒约简。引理 2 和定理 1 中的算法属于黑盒类。

373定理1的严密性。我们将证明,对于任何函数 $Q(n) : \mathbb{N} \rightarrow \mathbb{N}^+$,任何黑箱归纳都只能在 $O(n \log r)$ 时间内回答 $n \log n r = O()$ 个查询,并且Q(n)满足 $Q(n) = O(n)$,且对于任何常数

374 $Q(n) = \Theta(Q(c^{-1/n}))$ 是亚可加性的,即,对于任何 $x, y \geq 1$, $Q(x + y) \leq Q(x) + Q(y)$ 成立。

375 ■

376 ■

377这将证实定理 1 对黑盒类的紧密性。

378我们将设计一个可分解的问题和一个附带的数据结构,这两者在现实中都没有意义,但在数学上却是合理的。数据集S 380由n 个
任意元素组成;给定任何谓词,对S的查询始终返回|S| (元素和谓词的具体形式无关紧要)。每当被要求在 S 上“构建”数据382结
构时,我们都会故意浪费 $|S| \log |S|$ 时间,然后简单地输出数组中S的任意383排列。每当被要求“回答”查询时,我们都会故意浪费384
 $Q(|S|)$ 时间,然后返回|S|。该问题显然是可分解的。

我们认为,任何黑箱归纳算法A都需要 $\Omega(Q(n))$ 的时间来回答386每个查询。考虑一个任意查询,假设A通过搜索结构T (S₁), ..., T (S_n)
388来处理它,查询成本至少为

387 其中 $t \geq 0$ 且扫描S 扫描。通过我们的结构设计,

$$389 \quad \text{哦}(|S| \text{扫描}) + \sum_{i \in [t]} \text{问题}(|S_i|) \geq \Omega(|S| \text{扫描}) + Q(|S|) \quad (\text{通过子可加性})$$

XX:10 最大化延迟数据结构的最优性（又称数据库破解）

390
$$= \Omega(Q(|S| \text{ 扫描})) + Q(|S'|) \quad (\text{由 } Q(n) = O(n), Q(n) = \theta(Q(cn)))$$

391
$$= \Omega(Q|S| \text{ 扫描}) + \sum_{i \in [t]} Q(|S'_i|) \quad (\text{通过子可加性})$$

392
$$= \Omega(Q(n)). \quad (\text{根据 (5)})$$

393 根据log(r)-streak的定义,算法A必须处理 $r = \text{LogStreak}(n)$ 个查询

394 总成本为 $O(n \log r)$,显然不能超过 $O(n \log n)$ (记住 $r \leq n$) 。

395 因此,A 只能处理 $O(\frac{\log n}{Q(n)})$ 查询。

396 4 个针对具体问题的新型 DDS 算法

397 我们现在部署定理 1 来开发针对具体 DDS 问题的算法,重点关注

第 4.1 节中为398 个可分解问题,第 4.2 节中为不可分解问题。

399 4.1 可分解问题的应用

400 如上所述,如果可分解问题具有可以内置的数据结构T

401 $O(n \log n)$ 时间 (即 $B(n) = O(\log n)$)并支持 $Q(n) = O(\log n)$ 时间内的查询,

402 定理 1 直接得出一个 DDS 算法,其中 $\text{Time}(n, r) = O(n \log r)$ (对于所有 $r \leq n$) 。

403 下面列举了对数据库系统具有重要意义的问题的部分列表,

404 之前没有发现有同样保证的算法。

405 ■ 2D 正交范围计数。问题定义见 1.2 节。结构

406 T 可以是持久的“聚合”二叉搜索树 (BST)[30]。

407 ■ 矩形上的正交范围计数。数据集S是一组n 个2 维矩形 (即

408 R中的轴平行框² 给定任意一个 2 矩形q,查询返回有多少个

409 S与 q 相交的矩形。该问题可以简化为四个查询

410 前一个问题 (点的正交范围计数) [33]。T可以再次成为

411 持久性聚合 BST [30]。

412 ■ 点位置。数据集S是R的平面细分² 由n 个线段定义,

413 其中每个部分都与它所属的两个面的 ID 相关联。给定

414 任意点 $q \in R$ ², 查询返回包含 q 的细分面,

415 归根结底就是找到q正上方的线段 (即第一个“命中”的线段

416 射线从q 向上射出)。结构T可以是持久 BST [28] 或

417 柯克帕特里克结构[20]。

418 ■ 二维中 $k = O(1)$ 最近邻搜索。数据集S是R中的一组n 个点². 修复

419 常数整数 $k \geq 1$ 。给定一个点 $q \in R$ ², 查询返回P中最接近的k 个点

420 q. 结构T可以是基于 k 阶Voronoi构建的点位置结构 [20, 28]

421 图 (k 阶Voronoi 图可以在 $O(n \log n)$ 时间内计算出来[6]) 。

422 当 $k = 1$ 时, 已经发现了一个 DDS 算法,其中 $\text{Time}(n, r) = O(n \log r)$ (对于所有 $r \leq n$) [1]。然而,

423 [1] 中的算法很大程度上依赖于合并两个 (1 阶)Voronoi 图的能力

424 在线性时间内,因此不能轻易扩展到更高的k值。

425 ■ 度量空间中的近似最近邻搜索。数据集S由n 个对象组成

426 在具有恒定倍增维度的度量空间中 (这封装了任何欧几里得

427 维数为常数的空间。设 $\text{dist}(o_1, o_2)$ 表示

428 空间中的两个对象 o_1 和 o_2 。给定空间中的任意对象q ,查询

429 返回一个对象 $o \in S$, 使得 $\text{dist}(o, q) \leq (1 + \epsilon) \cdot \text{dist}(o^*, q)$ 对所有 $o^* \in S$, 其中是
430 一个常数。在 [14, 22] 中可以找到满足我们目的的结构 T 。

431 对于 R 中的正交范围计数 $\text{Count}(R)$ 其中 $d \geq 3$ 是一个固定常数 (定义在
432 1.2 节), 我们可以应用定理 1 得到一个不太寻常的结果。有可能 [3]
433 在 $O(n \log n)$ 时间内构建一个结构 T , 以回答 $Q(n) = O(n^{1-\epsilon})$ 中的查询 (时间
434 常数 $\epsilon > 0$ 可以任意小。定理 1 由此得到 DDS 算法
435 对于所有 $r \leq n$, $\text{Time}(n, r) = O(n \log r)^{1-\epsilon}$ 。由于可以任意接近于 0,
436 $\log(r)$ -条纹界限 $\text{LogStreak}(n) = n^{1-\epsilon}$ 比最大值 n 仅低
437 n 中的因子多项式。如果这个差距可以
438 对于所有常数 ϵ 都是封闭的。如果 $\text{LogStreak}(n) = n$ 的 DDS 算法可以
439 被发现, 那么该算法也将在 $O(n \log n)$ 中解决以下离线版本
440 时间: 我们给定 R 中的 n 个点的集合 P 以及 n 个 d 矩形的集合 Q ; 目标是
441 报告, 对于每个 d 矩形 $q \in Q$, P 中有多少个点在 q 中被覆盖。此离线
442 问题已经得到广泛研究, 但最快的算法仍然在 $n \log \Theta(d)$ 中运行 n
据我们所知是 443 次。

444 4.2 不可分解但可谱索引的问题

445 本小节将利用定理 1 来处理不可分解的问题
446 个问题 (至少不是很明显)。关键是要表明问题是 $(\log n, Q(n))$
447 频谱可索引为合适的 $Q(n)$ 。这本身就是一个有趣的话题, 因为我们
448 接下来通过开发新的 DDS 算法来演示, 其中 $\text{Time}(n, r) = O(n \log r)$
449 $r \leq n$ 在半平面包含、凸包包含、范围中值和二维线性
450 编程。这些问题的原始算法 [11, 19] 都是使用
451 自上而下的方法在第 2 节中进行了回顾。我们的算法提出了一个对比, 说明
452 定理 1 如何促进 DDS 算法的设计。

453 半平面约束。该问题可以转换为 [19] 等价于以下问题
454 形式, 利用几何对偶 [7]:
455 \blacksquare 穿过凸包的线。 S 是 R 中的 n 个点的集合 S 。给定 R 中的任意一行 l , 查询
456 确定 l 是否与 S 的凸包相交, 表示为 $\text{CH}(S)$ 。
457 我们将集中讨论上述问题。
458 现在假设我们给定 l 上的一个点 q , 该点位于 $\text{CH}(S)$ 之外。从 q 开始, 我们
459 可以发射两条切射线, 每条都接触 $\text{CH}(S)$, 但不会进入 $\text{CH}(S)$ 的内部。
460 在图 1(a) 中, S 是黑点集, 第一条射线经过点 $p_1 \in S$, 而
461 第二次通过 $p_2 \in S$ 。这两条射线形成一个“楔形”, 将 $\text{CH}(S)$ 包围在其中 (注意
462 楔角小于 180°)。我们称它为 q 在 $\text{CH}(S)$ 上的楔。线 l 通过
463 通过 $\text{CH}(S)$ 当且仅当通过楔形。如果 $\text{CH}(S)$ 的顶点已经
464 按顺时针顺序存储, 可以在 $O(\log n)$ 时间内找到两个切线 [25]。
465 我们将证明 “通过凸包的线” 问题是 $(\log n, \log n)$ 谱
466 可索引。取任意整数 $s \in [1, n]$ 并设置 $m = n/s$ 。将 S 任意分为 S_1, S_2, \dots, S_m 使得 $|S_i| = s$, 其中 $i \in [m-1]$ 且 $|S_m| = n - s(m-1)$ 。构建结构 T ,
468 使用 $O(s \log s)$ 时间计算每个 $i \in [m]$ 的 $\text{CH}(S_i)$, 并按顺时针方向存储其顶点
469 顺序。该结构的构建时间为 $O(n \log s)$ 。让我们看看如何回答查询
470 带有直线 l 。再次假设给定 l 上的 $\text{CH}(S)$ 外的一个点 q 。对于每个 $i \in [m]$,
471 在 $O(\log s)$ 时间内计算 $\text{CH}(S_i)$ 上 q 的楔形。从这 m 个楔形中, 可以得到一个简单的
472 任务是在 $O(m)$ 时间内获得 $\text{CH}(S)$ 上的 q 的楔形 (我们将处理一个更一般的
473 问题在讨论 “凸包包含” 时很快会得到解决)。现在, l 是否与
474 $\text{CH}(S)$ 可以轻松确定。到目前为止, 查询时间为 $O(m \log s)$ 。

XX:12 最大化延迟数据结构的最优性 (又称数据库破解)

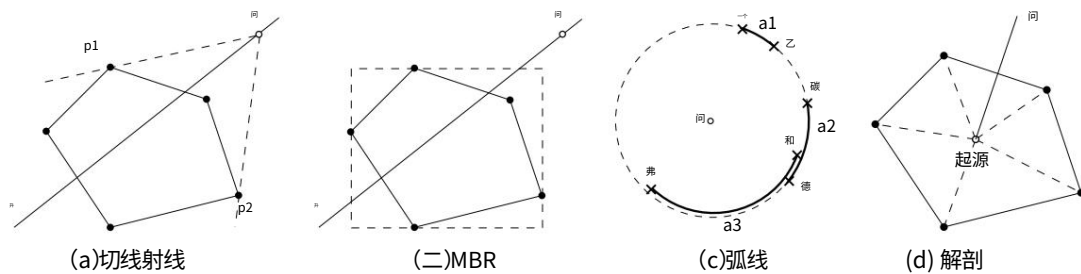


图 1.4.2 节中关键概念的解释

剩下要解释如何找到 q 。如果我们已经

具有 S 的最小轴平行边界矩形,表示为 $MBR(S)$ 。注意

$MBR(S)$ 必定含有 $CH(S)$; 见图 1(b)。显然, $MBR(S)$ 可以得到

在 $O(m)$ 时间内从 $MBR(S_1), MBR(S_2), \dots, MBR(S_m)$ 生成, 而每个 $MBR(S_i)$ ($i \in [m]$) 都可以

在结构 T 的构造过程中可以在 $O(s)$ 时间内计算出来。因此我们得出结论

问题是 $(\log n, \log n)$ 谱可索引的, 现在可以应用定理 1。

凸包包含。在这个问题中, S 是 R 中的一组 n 个点

². 给定任意

点 $q \in R$, 查询确定 q 是否被 $CH(S)$ 覆盖。我们将证明问题

是 $(\log n, \log n)$ 频谱可索引的。

取任意整数 $s \in [1, n]$, 设 $m = \lceil n/s \rceil$ 。将 S 任意分为 S_1, S_2, \dots, S_m

使得 $|S_i| \leq s$ 其中 $i \in [m-1]$ 且 $|S_m| = n - s(m-1)$, 则 $|S_i| \leq s$ 。要构建结构 T , 计算

$CH(S_i)$, 对每个 $i \in [m]$ 进行排序, 并按顺时针顺序存储其顶点; 这需要 $O(n \log s)$ 时间, 因为

之前解释过。让我们看看如何回答给定点 q 的查询。对于每个 $i \in [m]$, 是否

q 是否被 $CH(S_i)$ 覆盖可以在 $O(\log s)$ 时间内得到检验[25]。如果对任意 i 的答案为是,

点 q 一定被 $CH(S)$ 覆盖, 这样就完成了。接下来的讨论假设

对于所有 i , q 都在 $CH(S_i)$ 之外。在 $O(m \log s)$ 时间内, 计算 q 在 $CH(S_i)$ 上的楔形

对于所有 $i \in [m]$, 如上一个问题所述。

剩下的就是根据 m 个楔形确定 q 是否被 $CH(S)$ 覆盖。这可以是

重新建模为以下问题。以 q 为中心放置一个任意圆。对于每个楔形,

它的两条边界射线与圆相交成小于 180° 的圆弧

. 设 a_1, a_2, \dots, a_m 为

以这种方式产生的弧, 并定义一个

是覆盖它们的圆上的最小圆弧

全部。图 1(c) 显示了 $m = 3$ 的示例。圆弧 a_1 由点 A 和 B 所对, 圆弧

点 C 和 D 组成圆弧 a_2 , 点 E 和 F 组成圆弧 a_3 。其中, 最小圆弧 a

顺时针

从点 A 到 F 。至关重要, 当且仅当

跨度至少 180° 。 (这

是图 1(c) 中的情况) 注意, 如果 q 在 $CH(S)$ 之外, 那么 a 必须被

q 在 $CH(S)$ 上的楔形, 必须小于 180° 。

因此, 目标是确定

是否至少为 180° 。

可以在 $O(m)$ 时间内实现目的 (即使 m 个弧以

任意顺序)。为此, 我们逐个处理圆弧, 保持最小圆弧

覆盖已经处理的弧, 并在清楚时停止算法

必须至少为 180° 。具体来说, 对于 $i = 1$, 只需设置一个

到 a_1 。迭代地, 给定

下一个圆弧 a_i ($i \geq 2$), 在常数时间内检查圆弧是否小于 180° 。

可以覆盖人工智能和

如果是, 请更新到该弧; 否则, 停止算法。例如, 在图 1(c) 中,

当处理完 a_2 之后, 我们维护的 a 就从 A 顺时针转到 D 了。当处理 a_3 的时候,

该算法实现了必须至少 180° 。

并因此终止。

因此我们得出结论, 凸包包含问题是 $(\log n, \log n)$ 谱

可转位, 现在可以应用定理 1。

范围中位数。在这个问题中, S 是存储在数组 A 中的一组 n 个实数。给定

513 满足 $1 \leq x \leq y \leq n$ 的任意整数对 (x, y) , 查询返回

514 $A[x : y]$ 。我们将证明该问题是 $(\log n, \log n)$ 谱可索引的。修复任何

515 整数 $s \leq n$ 且 $m = \lfloor n/s \rfloor$ 。定义 $S_i = A[(i-1)s + 1 : i \cdot s]$, 其中 $i \leq m-1$

516 且 $S_m = A[(m-1)s + 1 : n]$ 。接下来, 我们假设 $s \leq \sqrt{n}$; 否则, 只需使用

517 $O(n \log n) = O(n \log s)$ 的时间来创建一个 $[11]$ 的结构在整个 S 上, 它能够

518 在 $O(\log n) = O(\log s)$ 时间内回答 S 上的任何查询。

519 建立一个结构 T , 对于每个 $i \in [m]$, 存储 S_i 按升序排列, 但每个元素

520 S_i 的 520 应该与它在 A 中的原始位置索引相关联。显然, T 可以

521 构建时间为 $O(n \log s)$ 。让我们看看如何回答带有谓词 (x, y) 的查询。首先,

522 确定值 $a, b \in [m]$ 使得 $A[x] \in S_a$ 且 $A[y] \in S_b$, 这可以简单地

523 在 $O(m)$ 时间内完成。扫描 S_a 并识别子集 $S'_a = S_a \cap A[x : y]$ (对于每个元素

524 在 S_a 中, 检查它在 A 中的原始索引是否在 $[x, y]$ 中。因为 S_a 是排序的, 我们可以得出

525 按排序顺序在 $O(s)$ 时间内完成。以同样的方式, 计算 $S'_b = S_b \cap A[x : y]$ 在

526 $O(s)$ 时间。此时, $A[x : y]$ 的所有元素都已在 $b - a + 1$ 中分割

527 排序数组: $S'_a, S_a + 1, S_a + 2, \dots, S_b - 1, S'_b$ 。现在的目标是找到第 $(y - x + 1)/2$

528 这些数组并集中的最小元素。Frederickson 和 Johnson [10] 描述了

529 一种从排序数组的并集中选择给定等级的元素的算法。它们的

530 在我们的场景中, 算法运行时间为 $O(m \log s) = O(m \log s)$ 。总体查询时间为

531 $O(s + m \log s) = O(m \log s)$, 因为 $s \leq \sqrt{n}$ 。

532 我们现在得出结论, 范围中位数问题是 $(\log n, \log n)$ 谱可索引的

533 并准备应用定理 1。

534 二维线性规划。利用几何对偶 [7], 该问题可以表示为

535 将 [19] 转换为“穿过凸包的线”(我们已经解决了), 如下所示:

536 射线离开凸包。其中, S 是 R 中的一组 n 个点 $\{p_i\}_{i=1}^n$ 使得 $CH(S)$ 覆盖

537 原点。给定从原点发出的任何射线 q , 查询将返回边 e_{exit}

538 $CH(S)$, 其中 q 退出 $CH(S)$ 。

539 我们将集中讨论上述问题。在继续之前, 让我们陈述两个事实

540 关于问题:

541 给定任意射线 q , 都有可能找到 e_{exit} , 将此称为基本算法。使用 [21] 中的算法, 在 $O(n)$ 时间内, 我们

542

543 我们可以在 $O(n \log n)$ 时间内创建一个结构, 以在 $O(\log n)$ 时间内回答任何查询。首先

544 计算 $CH(S)$, 然后使用连接原点和所有

545 顶点; 参见图 1(d)。射线 q 的查询可以通过查找

546 q 经过的部分三角形。我们将其称为基本结构。

547 与迄今为止讨论的所有问题不同, 目前我们无法证明“射线退出

548 凸包”是 $(\log n, \log n)$ 谱可索引的。然而, 根据定理 1, 我们

549 不必! 只需表明问题是 $(\log n, nc)$ 频谱可索引 $^{1-c}$ 登录 n 查询

550 任何正常数 $c < 1$ 。定理 1 允许我们回答 $r \leq n$

551 $O(n \log r)$ 时间。当 r 达到 n 时 $^{1-c}$, 我们有能力建立基本结构

552 花费 $O(n \log n) = O(n \log r)$ 的时间回答每个后续查询, 时间为 $O(\log n) = O(\log r)$ 。

553 这使我们能够实现对于所有 $r \leq n$ 的 $\text{Time}(n, r) = O(n \log r)$ 。

554 我们将证明“射线离开凸包”问题是 $(\log n, \sqrt{n})$ 频谱可索引。

555 我们将通过 [19] 的结果来实现这一目的, 其中 Karp, Motwani 和

556 Raghavan 采用自上而下的方法构建了具有以下属性的二叉树 T 。

557 如果节点 u 位于 T 的层级 ℓ , 那么 u 与 $n/2$ 的集合 $S(u)$ 相关联 $^\ell$ S 中的点。

558 树的前 ℓ 层可以在 $O(n \cdot \ell)$ 时间内构建。

XX:14 最大化延迟数据结构的最优性（又称数据库破解）

559 ■ 一个查询的答案最多是遍历T的一条从根到叶的路径。如果搜索过程延伸到节点u,那么通过在 $S(u)$ 上运行基本算法
560 [21],可以在 $O(|S(u)|)$ 时间内找到目标边eexit。

561

562 回到我们的场景,固定任何整数 $s \leq n$ 。在 $O(n \log \sqrt{s}) = O(n \log s)$ 时间内在S上构建T的前 $1 + \log \sqrt{s}$ 层。为了回答查
563 询,如果尚未找到边eexit,我们将T的路径下降到层数为 $\log \sqrt{s}$ 的节点u。集合 $|S(u)|$ 最多有 $n/2 \log \sqrt{s} = n/\sqrt{s}$ 个点。因此,我
564 们可以在 $S(u)$ 上运行基本算法,在 $O(n/\sqrt{s}) = O(n \cdot \sqrt{s})$ 时间内找到eexit。因此,“射线离开凸包”问题是 $(\log n, \sqrt{n})$ 567谱可
565 索引的。

566

568我们用上面的讨论所揭示的关于自上而下的方法和我们的约简之间的内在联系的评论来结束本节。本质上,我们逐步构建[19]的结构: T的第 i ($i \geq 1$)个“时期”(在第3.2节的证明中)级别。这不同于[19]中节点“在一次接触时展开”的方式,如第2节所述。事实上,所有现有的基于自上而下方法设计的DDS算法都可以通过“谱可索引性”的桥梁封装到我们的约简框架中,就像我们为“射线离开凸包”所展示的方式一样。

571 重建前 2 2^{i-1}

576 5 使用具有 $\omega(n \log n)$ 构造时间的结构的DDS

577到目前为止,我们的讨论集中在 $B(n) = O(\log n)$ 的数据结构上。在本节中, 578我们将首先说明黑盒归约的这一条件的必要性,以保证非常量 $\log(r)$ 条纹边界。作为第二步,我们提出了定理1580的扩展,该定理允许部署 $\max\{B(n), Q(n)\} = \text{polylog } n$ 的结构,以生成一个DDS算法,对于所有 $r \leq n$,该算法的 $\text{Time}(n, r)$ 都很好。

582 $B(n) = \omega(\log n)$ 的常数条纹边界。我们随后的硬度论证583要求 $n \cdot B(n)$ 为凸函数。考虑任何黑盒归约算法A。

584回想一下, A需要处理具有受限数据结构的任何可分解问题(读者可能希望在继续阅读之前查看第3.3节)。假设A可以保证所有此类问题的 $\text{LogStreak}(n)$ 的 $\log(r)$ -streak 界限,即,对于所有 $r \leq \text{LogStreak}(n)$, A总是可以在 $O(n \log r)$ 时间内回答 r 个查询。我们将证明,

588 如果 $B(n) = \omega(\log n)$,则 $\text{LogStreak}(n)$ 必须是 $O(1)$ 。

589 以类似于第3.3节的方式,我们将设计一个可分解的问题和一个伴随的数据结构。数据集S包含 n 个任意元素;给定任何谓词,对S的查询始终返回 $|S|$ 。每当被要求在S上构建数据结构T(S)时,我们都会故意浪费 $|S| \cdot B(|S|)$ 时间,然后在数组中输出S的任意排列。

593每当被要求回答一个问题时,我们都会立即在常数时间内返回 $|S|$ 。换句话说,函数 $Q(n)$ 固定为1。

595 此后,我们将 r 固定为A在给定我们设计的数据结构时可以确保的 $\text{LogStreak}(n)$ 的值。我们假设 $r = \omega(1)$;否则, $\text{LogStreak}(n) = O(1)$ 597并且我们的工作已完成。由于它在 $O(n \log r)$ 时间内回答了 r 个查询,因此至少有一个查询)。我们将在其余部分集中讨论这个特定查询

598 必须有一个 $O(\frac{n}{r})$ 的参数成本。

600 回想一下3.3节,为了回答这个查询,算法A需要搜索一定数量(包括0)的结构 $T(S_1), T(S_2), \dots, T(S_{\leq S_n})$ 。由于A需要 $n \log r$ 支付 $\Omega(|S_{\text{scan}}|)$ 的成本来扫描S以获得一些

601 并扫描子集S,则 $|S_{\text{scan}}| \leq \frac{n}{r}$

602 常数 $\alpha > 0$ 。我们将考虑 r ,我们知道它 $\frac{n}{r} \geq \alpha \cdot 603$ 扫描, 是 $\omega(1)$,足够大以使

604 一个 $\frac{n \log r}{r} \leq n/2$ 。因为 A 必须服从 (5), 所以我们可以断言

605 $|S_1 \cup \dots \cup S_t| \geq |S| - |S_{\text{排除}}| \geq n/2。$ (6)

606 这意味着 $t \geq 1$ 。由于算法 A 必须花费常数时间来搜索每个结构

607 (秒 s_i) ($i \in [t]$), 则我们必须有

608 $t = O((n/r) \cdot \log r)。$ (7)

609 然而, 通过我们如何设计 T, 建造 T 的总成本 ($S_1, T(S_2), \dots, T(S_{\text{re}})$ 是

610 $|S_{\text{re}}| \cdot B(|S_{\text{re}}|)。$ (8)

611 设 $\lambda = \min_{i \in [t]} |S_{s_i}|$; 从 (6) 式可知 $\lambda \geq n/2$ 。由于 $n \cdot B(n)$ 为凸函数, 且 $B(n)$ 是非递减的, 我们知道当 $|S_{s_i}| = \lambda/t$ 且对于所有 $i \in [t]$ 。因此:

612 $|S_{s_i}| = \lambda/t$ 且对于所有 $i \in [t]$ 。因此:

613 $(8) \geq \lambda \cdot B(\lambda/t) \geq \frac{n}{2} \cdot B\left(\frac{n}{2t}\right)。$ (9)

614 从 (7) 式可得 $n/(2t) = \Omega(r/\log r)$ 。因为 $B(n) = \omega(\log n)$, 所以基本渐近

615 分析表明 $B(n/(2t))$ 必定是 $\omega(\log r)$ 。

616 我们现在得出结论, (9), 因此 (8), 必须是 $\omega(n \log r)$ 。但这与

617 算法 A 可以在 $O(n \log r)$ 时间内处理 r 个查询。因此, 我们假设

618 $r = \omega(1)$ 必定是错误的。

619 个 DDS 算法, $B(n) = \omega(\log n)$ 。具有 $\omega(n \log n)$ 构造的数据结构

620 时间对于 DDS 仍然有用, 只要我们不沉迷于回答 r 个查询

621 $O(n \log r)$ 时间。为了使其正式化, 我们修改定理 1 背后的技术以获得

622 另一个具有以下保证的通用减少。

623 ► 定理 3. 假设 $B(n)$ 和 $Q(n)$ 是非减函数, 并且都是

$O(\log^\gamma n)$ 其中 $\gamma \geq 1$ 为常数。每个 $(B(n), Q(n))$ 谱可索引问题都承认

625 时间 $(n, r) = O(n \log^\gamma r)$, 其中 $r \leq n$ 。

626 证明与第 3.2 节中介绍的类似, 并移至附录 A。

627 上述内容的一个有趣应用是 R 中的正交范围计数/报告 ^d 其中 d 是

628 固定常数至少为 3 (参见 1.2 节中的问题定义)。范围树增强

629 具有分数级联 [9], 可在 $O(n \log d - 1)$ n 点时间, 回答 a

630 计数/报告查询同样在 $O(n \log d - 1)$ n 时间。该问题可分解,

631 因此 $(\log d - 1, \log d - 1)$ 频谱可索引。定理 3 直接给出了 DDS 算法

632 时间 $(n, r) = O(n \log d - 1)$ r , 对所有 $r \leq n$, 这严格改进了 [19] 的结果, 即

633 在第 1.2 节中提到。

634 参考

635 1 Alok Aggarwal 和 Prabhakar Raghavan。最近邻的延迟数据结构

636 问题。信息处理快报 (IPL), 40 (3):119-122, 1991。

637 2 杰里米·巴尔贝、安库尔·古普塔、斯里尼瓦萨·拉奥·萨蒂和乔纳森·索伦森。接近最优

638 内部和外部存储器中的在线多选。J. 离散算法, 36:3-17, 2016 年。

639 3 Jon Louis Bentley 和 Hermann A. Maurer。范围的有效最坏情况数据结构

640 搜索。Acta Inf., 13:155-168, 1980。

641 4 Manuel Blum, Robert W. Floyd, Vaughan R. Pratt, Ronald L. Rivest 和 Robert Endre

642 Tarjan。选择的时间界限。计算机与系统科学杂志 (JCSS),

643 7(4):448-461, 1973。

XX:16 最大化延迟数据结构的最优性（又称数据库破解）

5格思·斯托廷·布罗达尔 (Gerth Stolting Brodal)、比特·格费勒 (Beat Gfeller)、艾伦·格隆伦德·乔根森 (Allan Gronlund Jorgensen) 和彼得·桑德斯 (Peter Sanders)。向
最优范围中值。理论计算机科学, 412(24):2588–2601,2011。
6 Timothy M. Chan 和 Konstantinos Tsakalidis。二维和三维空间的最优确定性算法
3-d浅岩层。离散与计算几何, 56(4):866–881,2016年。
7 Bernard Chazelle、Leonidas J. Guibas 和 DT Lee。几何对偶的力量。BIT
数值数学, 25 (1) :76-90,1985。
8 Yu-Tai Ching、Kurt Mehlhorn 和 Michiel HM Smid。动态延迟数据结构。
信息处理快报 (IPL) , 35 (1) :37-40,1990。
9马克·德伯格、奥特弗里德·张、马克·范·克雷维尔德和马克·奥维马斯。计算型
几何:算法与应用。Springer -Verlag,第 3 版,2008 年。
654 10 Greg N. Frederickson 和 Donald B. Johnson。
x+y 和带排序列的矩阵。计算机与系统科学杂志 (JCSS),
655 24(2):197–208, 1982。
656 11 Beat Gfeller 和 Peter Sanders。迈向最优范围中值。国际
657 自动机、语言和编程讨论会 (ICALP),第 475–486 页,2009 年。
第 659 章12 Goetz Graefe、Felix Halim、Stratos Idreos、Harumi A. Kuno 和 Stefan Manegold。并发性
660 自适应索引控制。VLDB Endowment 论文集 (PVLDB), 5(7):656–667,
661 2012 年。
662 13 Felix Halim、Stratos Idreos、Panagiotis Karras 和 Roland HC Yap。随机数据库
663 破解:迈向主内存列存储中的稳健自适应索引。
664 VLDB捐赠基金 (PVLDB) , 5 (6) :502–513,2012。
665 14 Sarel Har-Peled 和 Nirman Kumar。低维数据的近似最近邻搜索
666 查询。SIAM计算杂志, 42(1):138–159,2013 年。
667 15 Sarel Har-Peled 和 S. Muthukrishnan。范围中值。欧洲研讨会论文集
668 算法论 (ESA) ,第 503-514 页,2008 年。
669 16 Stratos Idreos、Martin L. Kersten 和 Stefan Manegold。数据库破解。第 68–78 页,
670 2007 年。
671 17 Stratos Idreos、Martin L. Kersten 和 Stefan Manegold。自组织元组重建
672 在列存储中。在ACM 数据管理 (SIGMOD) 论文集,第 297-308 页,
673 2009年。
674 18 Stratos Idreos、Stefan Manegold、Harumi A. Kuno 和 Goetz Graefe。合并已破解的内容,
675 破解合并的内容:主内存列存储中的自适应索引。
676 VLDB捐赠基金 (PVLDB) , 4 (9) :585–597,2011。
677 19 Richard M. Karp、Rajeev Motwani 和 Prabhakar Raghavan。延迟数据结构。
678 SIAM 计算杂志, 17(5):883–902, 1988。
679 20 David G. Kirkpatrick。平面细分中的最佳搜索。SIAM计算杂志,
680 12(1):28–35, 1983。
681 21 David G. Kirkpatrick 和 Raimund Seidel。终极平面凸包算法? SIAM
682 计算机杂志, 15 (1) :287-299,1986。
683 22 Robert Krauthgamer 和 James R. Lee。导航网络:邻近度的简单算法
684 搜索。在ACM-SIAM 离散算法年度研讨会 (SODA) 论文集上,
685 第798–807页,2004年。
第 686 章23康斯坦丁诺斯·兰普罗普洛斯、法特梅·扎尔巴尼、尼科斯·马穆利斯和帕纳吉奥蒂斯·卡拉斯。
687 高维度量空间中的自适应索引。VLDB 捐赠基金论文集
688 (PVLDB) , 16 (10) :2525–2537,2023年。
689 24 Rajeev Motwani 和 Prabhakar Raghavan。延迟数据结构:查询驱动的几何搜索问题预处理。在计算科学研讨会论文集
690
691 几何 (SoCG) ,第 303-312 页,1986 年。
692 25 Mark H. Overmars 和 Jan van Leeuwen。飞机配置的维护。期刊
693 计算机与系统科学学报(JCSS), 23(2):166-204, 1981。
694 26 Bryce Sandlund 和 Sebastian Wild。惰性搜索树。在IEEE 年度论文集上
695 计算机科学基础研讨会 (FOCS),第 704–715 页,2020 年。

696 27 Bryce Sandlund 和 Lingyi Zhang.可选堆和最佳惰性搜索树。在ACM-SIAM 离散算法年度研讨会 (SODA) 论文集,第 1962-1975页,
697 2022 年。

698

699 28 Neil Sarnak 和 Robert Endre Tarjan. 使用持久搜索树进行平面点定位。

700 ACM通讯(CACM), 29(7):669–679,1986。

701 29菲利克斯·马丁·舒克内希特、阿莱克·金达尔和延斯·迪特里希。实验评估和

702 数据库破解分析。VLDB杂志, 25 (1) :27-52,2016年。

703 30 Yufei Tao 和 Dimitris Papadias。空间数据库中的范围聚合处理。IEEE知识与数据工程学报 (TKDE), 16(12):1555–1570,
704 2004 年。

705 31 Fatemeh Zardbani,Peyman Afshani 和 Panagiotis Karras。重温数据库破解的理论和实践。在《扩展数据库技术 (EDBT) 论文集》中,第 415–418 页, 2020 年。

706

707

708 32 Fatemeh Zardbani,Nikos Mamoulis,Stratos Idreos 和 Panagiotis Karras。具有空间范围的对象的自适应索引。VLDB Endowment
709 (PVLDB) 论文集, 16(9):2248– 2260,2023 年。

710

711 33 Donghui Zhang,Vassilis J. Tsotras 和 Dimitrios Gunopulos。对具有范围的对象进行高效聚合。在ACM 数据库系统原理研讨
712 会(PODS) 论文集,第 121-132 页,2002 年。

713

714 定理 3 的证明

715 我们的算法以epoch 为单位执行。在第 i 个($i \geq 1$) epoch开始时,我们设置 $s = 2^{2i}$,并在 $O(n \cdot B(s)) = O(n \cdot 2^i \cdot \gamma)$ 时间内在S上创建一个结
716 构T (即 $B(n), Q(n)$)频谱可索引性所承诺的结构)。该结构允许我们在 $O(n \cdot Q(s))$ 时间内回答任何查询。当s 个查询在该 epoch期间得到回答后,第 i
717 个 epoch 结束。

718 $\frac{n}{s}$
719 这些查询的总成本为

720 仅 $\frac{n}{s} \cdot Q(s) = O(n \cdot Q(s)) = O(n \cdot 2^i \cdot \gamma)$ 。

721 因此,第 i 个时期的总计算时间为 $O(n \cdot 2^i \cdot \gamma)$ 。因为 $\gamma \geq 1$,我们知道当增加1

722 时, $n \cdot 2^i \cdot \gamma$ 至少翻倍。因此,所有时期的总成本为

723 渐近受 $O(n \cdot 2^h)$ 支配, h 的值是满足两个条件的最小整数 $2^h \geq n \cdot \gamma$,其中h 是所需的 epoch 数。确切地说,

724 $h \geq 1: \leq n \cdot \gamma$ (s的值不能超过n) ;

725 ■ C1: $2^{h-1} < n \cdot \gamma$

726 ■ C2: $\sum_{i=1}^{h-1} 2^{2i} \geq r$ (h 个时期可回答的查询数量必须至少为r)。

727 令H表示最小整数 $h \geq 1$,使得 $2^{2H} \geq r$ 。这意味着

728
$$2^{2(H-1)} < r \leq 2^{2H} < 4r \tag{10}$$

729 当 $2^{2H} \leq n$,上述论证保证了迭代次数h 最多为H

730 (不等式 $2^{2H} \leq n$ 必然有 $2 \leq n$,满足条件C1)。此时所有 $\gamma = O(n \cdot 2^H \cdot \gamma) = O(n \log \gamma r)$ 。

731 本为 $O(n \cdot 2^H)$

732 如果 $2^{2H} > n$,上述论证就不成立。然而,当这种情况发生时,我们从(10)知道 $r > n$,从而导致 $r > \sqrt{n}$ 。一旦r达到 \sqrt{n} 即捕捉点 我们就会在 $O(n \log \gamma$
733 $n)$ 时间内在整个S上创建一个结构T,并使用它回答 $O(Q(n)) = O(\log \gamma n)$ 时间内的每个后续查询,直到 $r = n$ 。捕捉点之后的查询总成本为 $O(n \log \gamma n) = O(n$
734 $\log \gamma r)$ 。因此,我们得到了一个算法,该算法保证对于所有 $r \leq n$, $\text{Time}(n, r) = O(n \log \gamma r)$,从而完成了定理 3的证明。

735

736

737