

KDD '22, 2022 年 8 月 14-18 日, 美国华盛顿特区

陈光方等。

## 证据

### A.5 命题 8 的证明

鉴于  $X = +\infty$ , 即问题 1 中没有模型约束, 推导出附录 A.2 中的相似比对。

当同时给定  $\lambda = 0$  时, 时间约束  $\Theta(r) \leq \lambda = 0$  在对齐的元组中强制执行相同的时间戳, 即  $r[U_1] = r[U_2] = \dots = r[U_m]$ , 这正是附录 A.2 中等式对齐所要求的。

### A.6 定理 1 的证明

我们首先表明问题出在 NP 中。给定一个对齐的实例, 这三个条件都可以在多项式时间内得到验证。条件 (1) 通过将元组相互比较来验证  $O(|R|^2)$  时间。回想一下  $|R| \leq \text{分钟}[\text{吨}1, \text{吨}2], \dots, \text{吨}k] = G$ 。对于条件 (2), 时间约束  $\lambda$  可以通过遍历中的元组来检查  $O(|R|)$  时间和模型约束  $X$  通过搜索每个元组的最近邻居来检查  $r \in R$  在  $O(|R| \log |M|)$  时间。条件 (3) 可以简单地通过比较  $\lambda$  来检查。  $R$  和  $\lambda$ 。

接下来, 我们将从  $\lambda \leq 3$  的场景开始, 并从最大 3 维匹配 [17, 19], 卡普的 21 个 NP 完全问题之一 [20], 到我们的对齐问题。最后, 我们可以证明缩减也适用于  $\lambda \geq 4$ , 从而证明 NP-硬度。结合问题在 NP 中, 我们可以得出问题的 NP 完全性。

让甲、乙、丙是有限的不相交集, 并且  $P \subseteq A \times B \times C$ ,  $I \in P$ ,  $P = \{(a, b, c) \mid a \in A, b \in B, c \in C\}$ 。为了  $p_1 \in P$  ( $A_1, b_1, C_1$ ),  $p_2 \in P$  ( $A_2, b_2, C_2$ )  $\in$ , 我们说  $p_1$  和  $p_2$  在某个坐标上相交, 表示为  $p_1 \leftrightarrow p_2$ , 如果  $A_1 = A_2$ ,  $b_1 = b_2$  或者  $C_1 = C_2$ 。我们用  $P \subseteq P$  表示一个 3 维匹配, 如果没有元素  $P$  与其他人相交, 即  $\forall p_1 \in P, p_2 \in P, p_1 \neq p_2$ 。这最大 3 维匹配就是找到匹配  $P^*$  在所有可能的值中  $P$  拥有最多的三元组。判定问题是, 给定一个整数  $\lambda$ , 判定是否存在 3 维匹配  $P$  这样  $|P| \geq \lambda$ 。

举个例子最大 3 维匹配问题  $P$ , 为了在模型约束下构建相似性对齐问题, 我们首先设置  $\lambda = +\infty$ , 即, 任何两个对齐的元组都满足时间约束。然后, 我们创建  $\text{吨}1, \text{吨}2, \text{吨}3$

在我们的问题中, 通过为它们分配唯一的时间戳。同时, 让每个  $\text{吨}1 \in \text{吨}1$  相当于  $A \in A$ , 每个  $\text{吨}2 \in \text{吨}2$  相当于  $b \in B$  和每个  $\text{吨}3 \in \text{吨}3$  相当于  $C \in C$ 。接下来, 如果  $(A, b, c) \in P$ , 我们指定  $\text{吨}1 \in [V_1], \text{吨}2 \in [V_2]$  和  $\text{吨}3 \in [V_3]$  值的独特组合。然后, 我们设置回归模型以满足  $\text{吨}1 \in [V_1], \text{吨}2 \in [V_2] = \text{吨}3 \in [V_3]$ 。让模型约束  $X = 0$ , 因为每个  $(\text{吨}1 \in [V_1], \text{吨}2 \in [V_2], \text{吨}3 \in [V_3])$  是独一无二的, 我们有  $R = \{(\text{吨}1 \in [V_1], \text{吨}2 \in [V_2], \text{吨}3 \in [V_3]) \mid \text{吨}1 \in [V_1], \text{吨}2 \in [V_2] = \text{吨}3 \in [V_3]\}$ 。因此, 我们得到  $(\text{吨}1 \in [V_1], \text{吨}2 \in [V_2], \text{吨}3 \in [V_3]) \in R$  当且仅当  $(A, b, c) \in P$ 。

接下来, 我们将证明, 对于每个满足的 3 维匹配, 我们有  $|P| \geq \lambda$ , 当且仅当对齐的实例  $R = \{(\text{吨}1 \in [V_1], \text{吨}2 \in [V_2], \text{吨}3 \in [V_3]) \mid \text{吨}1 \in [V_1], \text{吨}2 \in [V_2] = \text{吨}3 \in [V_3]\}$  对应于  $P$  在我们的问题中也是满足 (1) 三个候选键的集合  $\text{吨}1, \text{吨}2, \text{吨}3$ , (2) 时间约束  $\Theta(r) \leq \lambda$  和模型约束  $\Delta(R, M) \leq X$  每个都满意  $r \in R$  和 (3)  $|R| \geq \lambda$ 。

首先, 根据定义, 我们假设  $|P| \geq \lambda$ 。自从  $\forall p_1 \in P, p_2 \in P, p_1 \neq p_2$ , 对于相应的, 我们会有  $\forall r_1 \in R, r_2 \in R, r_1 \neq r_2$ , 因此满足条件 (1)。对于条件 (2), 首先

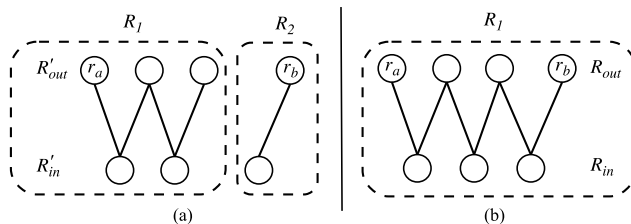


图 10: oa 路径的 (a) 断开连接图和 (b) 连接图示例

回想一下, 我们通过假设时间约束足够大来构建对齐问题。此外, 我们生成  $R$  通过满足模型约束, 从而自动满足条件 (2)。最后, 由于我们有  $(\text{吨}1 \in [V_1], \text{吨}2 \in [V_2], \text{吨}3 \in [V_3]) \in R$  当且仅当  $(A, b, c) \in P$ ,  $|R| = |P| \geq \lambda$  成立, 满足条件 (3)。

相反, 假设  $R$  满足条件 (1)、(2)、(3)。按照类似的步骤, 根据条件 (1),  $\forall r_1 \in R, r_2 \in R, r_1 \neq r_2$ , 相应的  $P$  有  $p_1 \in P, p_2 \in P, p_1 \neq p_2$ 。接下来, 根据条件 (3),  $|R| \geq \lambda$ 。自从  $(\text{吨}1 \in [V_1], \text{吨}2 \in [V_2], \text{吨}3 \in [V_3]) \in R$  当且仅当  $(A, b, c) \in P$ , 我们会得到  $|P| = |R| \geq \lambda$ 。满足 3 维匹配问题的条件。

### A.7 命题 2 的证明

我们首先表明模型约束下问题 1 相似性对齐的过程等同于  $\lambda \leq$  的 Set Packing ( $k$ -SP) 问题。给定一个地面集  $\mathcal{U}$  和一个集合小号的集合, 每个集合包含: 来自  $\mathcal{U}$  的元素,  $k$ -SP 问题是找到的最大子集  $\mathcal{S}$  所以它们成对不相交。在我们的场景中, 让来自原始数据的所有单元组的集合为

$\mathcal{U}$ , 然后让  $\text{吨}2$  相当于  $S$ 。这两个问题的等价性是直观的, 因为它们都最大化了目标集。

接下来, 我们的对齐搜索 (算法 2) 遵循以下框架  $d$ -SP 问题 [18, 24] 的最优局部搜索算法, 它被证明具有有界近似比。根据 [18] 中的证明, 对于  $\lambda$ -维时间序列和参数  $d$  对于对齐搜索, 近似比率的界限  $b$  获得。

### A.8 命题 3 的证明

条件 (4) 要求任何元组  $A \in \mathcal{U}$  必须至少与  $\mathcal{U}$  中的一个元组重叠  $\lambda$ 。假设一个元组  $A \in \mathcal{U}$  不与  $\mathcal{U}$  中的任何元组重叠  $\lambda$ 。结合条件 (3),  $A \in \mathcal{U}$  不与  $\mathcal{U}$  中的任何元组重叠  $\lambda$ 。这是不可能的, 因为  $\lambda \in \mathcal{U}$  已确保  $\mathcal{U}$  中的每个元组  $d \geq 2$  与  $\mathcal{U}$  中的至少一个元组重叠  $\lambda$ 。此外, 每次交换也将确保这一点, 因为它只交换有冲突的元组。因此, 任何元组  $A \in \mathcal{U}$  必须至少与  $\mathcal{U}$  中的一个元组重叠  $\lambda$ 。

条件 (5) 要求任何元组  $A \in \mathcal{U}$  必须至少与  $\mathcal{U}$  中的一个元组重叠  $\lambda$ 。这是因为对齐搜索算法的第 8 行确保我们遍历  $\mathcal{U}$  通过增加子集大小?。如果  $A \in \mathcal{U}$  不与  $\mathcal{U}$  中的任何元组重叠  $\lambda$ , 在里面以前的迭代较小?, 子集  $\mathcal{S} = \mathcal{U} \setminus \{A\}$  应该已经与当前的  $\mathcal{U}$  交换  $\lambda$ 。

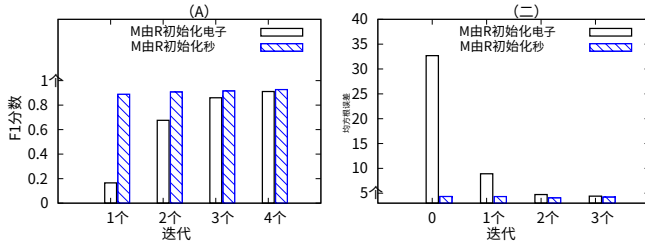


图11: 时间约束下相似性对齐多次迭代的对齐精度和模型预测性能80和模型约束X=8个在房子

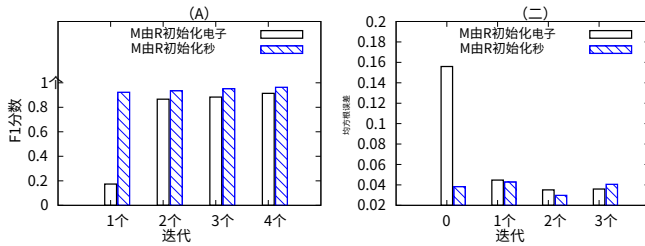


图12: 时间约束下相似性对齐多次迭代的对齐精度和模型预测性能 \ =85和模型约束X=0.5在水上

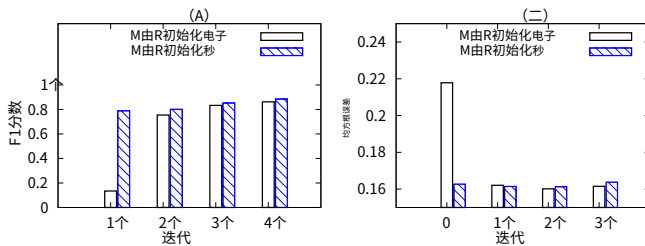


图13: 时间约束下相似性对齐多次迭代的对齐精度和模型预测性能75和模型约束X=0.5关于遥测

## A.9 引理 4 的证明

为了简单起见, 这里我们稍微滥用吨:8表示吨:8[\*:], 即时间戳吨:8. 让 $C < 0G$ ,  $C < 8=$ 表示的最大和最小时间戳 $A_8$ 个, 连续 $A_8$ 个和 $A_{8+1}$ 个在 $oa$ 路径中. 自从 $A_8 \cap A_{8+1}$ ,  $A_8$ 个和 $A_{8+1}$ 个必须重叠一些时间戳 $C$ :

, IE,  $C:8=C:(8+1)$ . 根据定义 2 中的时间约束, 我们有 $C:8 \leq C < 0G \leq C:8+$ 和 $C:(8+1) \leq C_{AG}$ .

因此, 我们得到 $C < 0G$   $8+1 \leq C:(8+1)+ \setminus = C:8+ \setminus \leq C < 0G$   $8+1$ , IE,  $C_{AG} \leq C_{AG} + \setminus$ . 通过迭代应用这个公式, 我们有 $C < 0G \leq C < 0G + \setminus$ .

同样, 我们也可以证明 $C < 0G$   $1 \leq C:8=+$ ;  $\setminus$ . 因此, 我们得到 $C_{AG} \leq C:8=+$ ;  $\setminus$ 和 $C < 0G \leq C:8=+$ ;  $\setminus$ 指示 $|C < 0G$   $1 \setminus - C:8=| \leq \setminus$ 和 $|C < 0G - C:8=| \leq \setminus$ .

最后, 我们有 $C:1-C:| \leq \text{最大} (|C < 0G$   $1 \setminus - C:8=|, |C < 0G - C:8=|) \leq \setminus$ ;  $\forall: \in \{1, 2, \dots, <\}$  即 $|C:1[*:] - C:|[*:]| \leq \setminus$ .

## A.10 引理 5 的证明

在第 3.3.1 节的条件 (2) 和 “中选候选人的要求” 中 $\setminus <$ , 中的元组 $>$ 直流电和 '8= 不应重叠

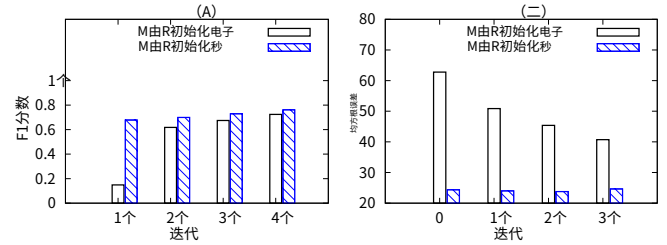


图14: 时间约束下相似性对齐多次迭代的对齐精度和模型预测性能 \ =35和模型约束X=60关于空气质量

同一组中的其他元组。因此, 如果 $A_0$ 和 $A_1$ 个由路径连接, 路径上的元组必须从 $>$ 直流电和 '8=, 即 $oa$ 路径。

## A.11 引理 6 的证明

首先, 我们将证明一个 $oa$ -path ( $A_0, \dots, A_1$ 个) 必须存在。自从 '8= 可以被交换 ' $>$ 直流电, '8=和 ' $>$ 直流电应该满足命题3的条件 (4) 和 (5), 即每个元组 $A$ 在 ' $>$ 直流电和 '8=与至少一个其他元组重叠。如果不存在 $oa$ -path ( $A_0, \dots, A_1$ 个), 这个图是不连通的, 即图的至少两个顶点没有通过路径连接。不失一般性, 我们可以假设断开的路径可以分为我连接部分 $R = \{ '1$ 个, '2个,  $\dots, '我$ 个, 因为每个元组 $A$ 至少与另一个元组重叠,  $\forall '8 \in R$ , 我们有 $|'8 \cap > \setminus| \geq 2$ . 回想一下 ' $>$ 直流电 $> |'8=|$ ?, 也就是说, 必须存在 ' $8 \in R$ ,

$|'8 \cap > \setminus| > |'8 \cap '8=|$  (示例参见图 10(a))。让

' $>$ 直流电 $= '8 \cap > \setminus$ 直流电和 '8= $= '8 \cap '8=$ . 我们可以安全地交换 '8=和 ' $>$ 直流电, 因为元组在 ' $>$ 直流电不与外面的其他元组重叠

'8=, 和 ' $>$ 直流电 $> |'8=|$ . 因此, 我们发现 ' $|'8=| < |'8 \cap > \setminus|$  可以与'交换 ' $>$ 直流电. 它与引理中的假设相冲突

6 算法已找到所有可能与 ' $|$  的交换 ' $|'8=| < ?$ 。

总而言之, 一个 $oa$ 路径 ( $A_0, \dots, A_1$ 个) 必须存在, 即这个图是一个连通图 (例子见图 10 (b))。

接下来, 我们将显示 $oa$ -path % 的长度 $> 0$ 小于 $2? - 1$ . 回想一下 ' $|'8=|$  = ?。最多有 ? 来自 ' $|$  的元组8=

在 $oa$ 路径中, 即 $|'8= \cap \% > 0| \leq ?$ . 由于元组来自 ' $>$ 直流电和 '8= 或者出现在 $oa$ 路径中, 以及开始和结束顶点 $A_0$

和 $A_1$ 个都来自 '8=, 最多有 ? 来自 ' $|$  的 1 个元组 $>$ 直流电在 $oa$ 路径中, 即 $|' > \setminus \text{直流电} \cap \% > 0| \leq ? - 1$ . 因此, 结合 $\% > 0 \subseteq > \setminus \text{直流电} \cup '8=$ , 我们有 $|\% > 0| = |' > \setminus \text{直流电} \cap \% > 0| + |'8= \cap \% > 0| \leq 2? - 1$ .

总之, 我们证明 $\forall A_0, A_1 \in '8=$ , 必须存在一个 $oa$ -path  $\% > 0 = (A_0, \dots, A_1)$  有长度 $\leq 2? - 1$ 。

## A.12 命题 7 的证明

通过结合引理 4 和 6, 一条来自 $A_0$ 到 $A_1$ 个存在, 有 $|C:0[*:] - C:1[*:]| \leq (2? - 1) \setminus, \forall: \in \{1, 2, \dots, <\}$ . 也就是说, 对于任何元组对 $A_0, A_1 \in '8=$ , 时间戳之间的差异 $A_0$ 和 $A_1$ 个不大于 $(2? - 1) \setminus$ 。

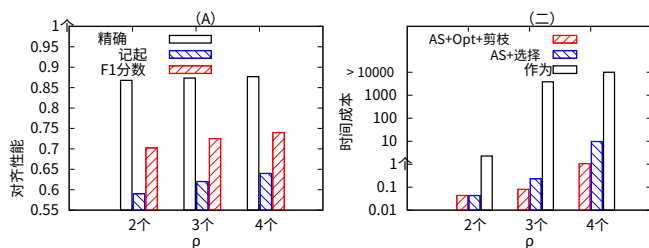


图 18: (a) 对齐性能和 (b) SAMC 的时间成本通过改变d在空气质量数据集上使用  $\lambda=35$  和  $X=60$

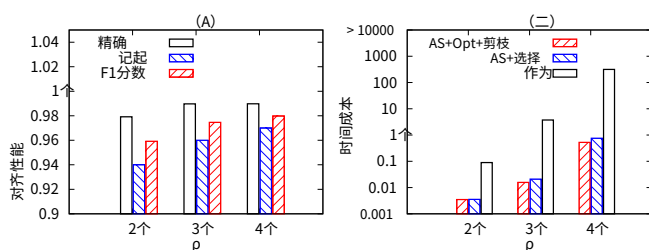


图 19: (a) 对齐性能和 (b) SAMC 的时间成本通过改变d使用  $\lambda=$  在 Fuel 数据集上120和  $X=35$

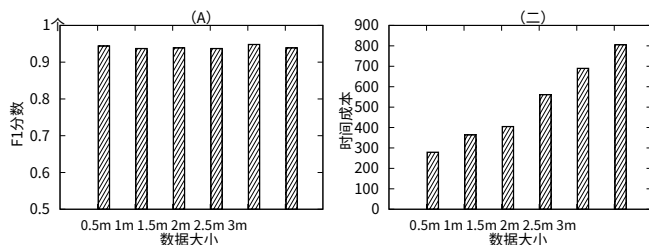


图 15: 在 Apache Spark 上比 Fuel 对齐的可扩展性

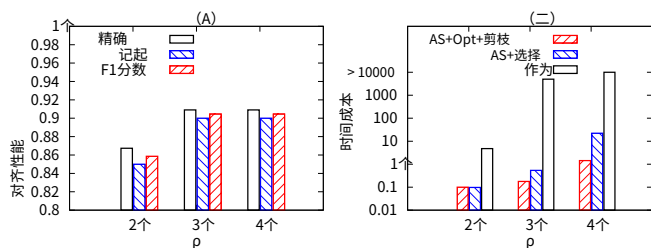


图 16: (a) 对齐性能和 (b) SAMC 的时间成本doverWater 数据集与  $\lambda=85$  和  $X=0.5$

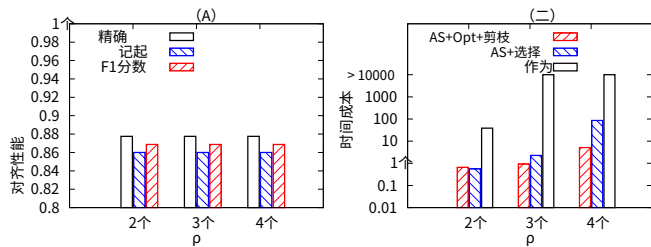


图 17: (a) 对齐性能和 (b) SAMC 的时间成本通过改变d使用  $\lambda=$  在遥测数据集上75和  $X=0.5$

## 其他结果

### A.13 第 4.1 节的附加结果

为了进一步评估大规模数据的可扩展性，我们在 Apache Spark 上实施了建议的 SAMC。图 15 说明了 Fuel 数据集的结果，范围从 0.1 到 300 万行。如图所示，F1-score 总体上是稳定的，而相应的时间成本几乎呈线性增长。

### A.14 第 4.2.1 节的附加结果

图 11-14 显示了在对其他数据集进行不同次数的迭代后该提议的结果。

### A.15 第 4.2.2 节的附加结果

图 16-19 通过不同的方式反对提案的结果丁，与其他数据集相比，我们将考虑交换的最大集合。