

Answer Queries Using Views

概要

- 背景
- 查询包含问题
- 使用视图回答查询
- 分析和总结

背景：使用物化视图加速查询

- Naïve方法：仅当以下条件满足时才使用物化视图
 - 查询与视图定义完全匹配 ✕
 - 视图被完整物化 ✓ ← 我们保留此条件
- 利用物化视图的机会远不止于此
 - 视图与查询逻辑等价 ($A=B \text{ AND } B=5$ 与 $A=5 \text{ AND } A=B$)
 - 视图与查询的选择条件可能重叠 ($A>10$ 与 $A>20$)
 - 视图可作为查询中JOIN操作的输入 ($R \text{ JOIN } S \text{ JOIN } T \rightarrow V \text{ JOIN } T$)
 - 查询的聚合操作可在视图上执行 ($\text{GROUP BY } A, B, C$ 与 $\text{GROUP BY } A, B$)
 -

概要

- 背景
- 查询包含问题
- 使用视图回答查询
- 分析和总结

基础问题: 查询包含与查询等价

- **数据库理论**领域的经典问题
- 查询包含: 给定查询 Q_1 和 Q_2 , 若对于任意数据库实例 D 都有 $Q_1(D) \subseteq Q_2(D)$, 则称 Q_2 包含 Q_1 , 记为 $Q_1 \sqsubseteq Q_2$ 。
- 查询等价: 给定查询 Q_1 和 Q_2 , 若 $Q_1 \sqsubseteq Q_2$ 且 $Q_2 \sqsubseteq Q_1$, 则称 Q_1 与 Q_2 等价, 记为 $Q_1 \equiv Q_2$ 。
- 考虑只有一个视图的情形
 - 给定视图定义 V_1 和查询 Q , 在 V_1 包含 Q 时, 可考虑使用 V_1 回答 Q

查询包含问题：复杂度

- Q_1 和 Q_2 均为合取查询时，查询包含问题被证明为**NP完全**问题
 - 合取查询：SPJ型查询，仅包含“属性=属性”或“属性=常量”型条件
- 实际应用中通常不是瓶颈：
 - 问题规模通常较小
 - 一定约束下存在多项式时间解法

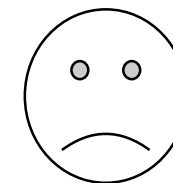
合取查询 (Conjunctive Query, CQ)

- 考虑表 $r(a, b)$ 和 $s(c, d, e)$
- SPJ型查询，仅包含“属性=属性”或“属性=常量”型条件
 - `SELECT r.b, s.d FROM r JOIN s WHERE r.a = s.c AND s.e = 5`
- 理论分析中通常使用Datalog语言表示
 - $(X, Y) :- r(W, X), s(W, Y, 5)$
 - `<head> :- <body>`
 - `body`包含一个或多个子目标
 - r, s 表示关系， W, X, Y 表示变量，共享的变量可理解为JOIN
 - 翻译：如果在 r 和 s 中存在记录 (W, X) 和 $(W, Y, 5)$ ，则将 (X, Y) 包含在结果中

包含映射

- 包含映射：给定查询 Q_1 和 Q_2 ，一个将 Q_2 中的变量映射到 Q_1 中的变量和常量的函数 h 被称为包含映射，如果：

- 对于 Q_2 中的每个关系 $R(x_1, x_2, \dots)$ ，在 Q_1 中有 $R(h(x_1), h(x_2), \dots)$
- $h(\mathbf{head}(Q_2)) = \mathbf{head}(Q_1)$ ，其中 $\mathbf{head}(Q)$ 表示 Q 的头部变量



- 定理： $Q_1 \sqsubseteq Q_2 \Leftrightarrow$ 存在从 Q_2 到 Q_1 的包含映射

- 例：

$Q_2:$	(X, Y)	$:-$	$r(W, X),$	$s(W, Y, Z)$
	\downarrow		\downarrow	\downarrow
$Q_1:$	(X, Y)	$:-$	$r(W, X),$	$s(W, Y, 5)$



CQ的查询包含问题

- 为什么这样一个看上去很简单的问题是NP完全的？
- 考虑布尔查询（头部无变量的查询）：
 - $Q_2() :- R(a, b), R(b, c), R(c, d), R(c, a)$
 - $Q_1() :- R(x, y), R(y, z), R(z, x)$
 - 问题： Q_2 是否包含 $Q_1 \Leftrightarrow$ 是否存在从 Q_2 到 Q_1 的包含映射？
- Three Color Problem
 - 给定图 $G(V, E)$ ，其中 $V = \{a, b, c, d\}$ ， $E = \{(a, b), (b, c), (c, d), (c, a)\}$
 - 是否可用三种不同的颜色 $\{x, y, z\}$ 为图 G 染色？ \leftarrow NP完全问题
 - 使 G 中任意相邻两点颜色不同

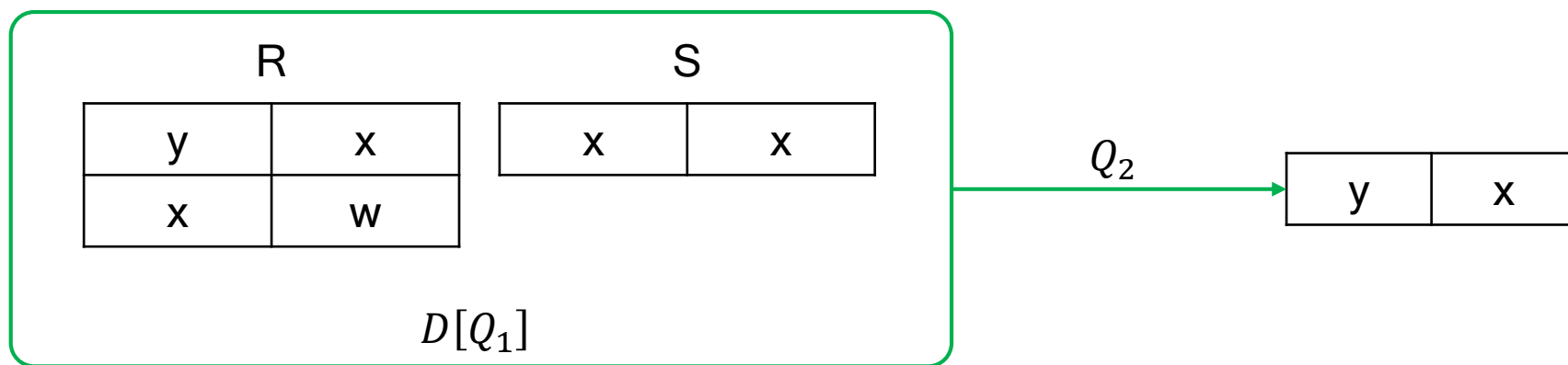
解法： 典型数据库

- **典型数据库**： 给定CQ查询 Q ， 其**典型数据库** $D[Q]$ 是 Q 的主体部分各关系实例化的结果， 实例化过程中将**变量名**作为常量。
- 例： 典型数据库
 - $Q_1(y, x) :- R(y, x), S(x, x), R(x, w)$
 - $D[Q_1] = \{ \{(y, x), (x, w)\}, \{(x, x)\} \}$

R		S	
y	x	x	x
x	w		

解法：典型数据库（续）

- 定理： $Q_1 \sqsubseteq Q_2 \Leftrightarrow \mathbf{head}(Q_1) \in Q_2(D[Q_1])$
- 解法：在 Q_1 的典型数据库上执行 Q_2 ，若 Q_1 的头部元组在结果中，则 Q_2 包含 Q_1
 - $Q_1(y, x) :- R(y, x), S(x, x), R(x, w)$
 - $Q_2(x, y) :- R(x, y), S(y, z), R(z, w)$



无环CQ查询

- 现实中环状JOIN较为少见*
 - 无环CQ查询的包含问题有多项式时间的解法
 - 基于超图分解

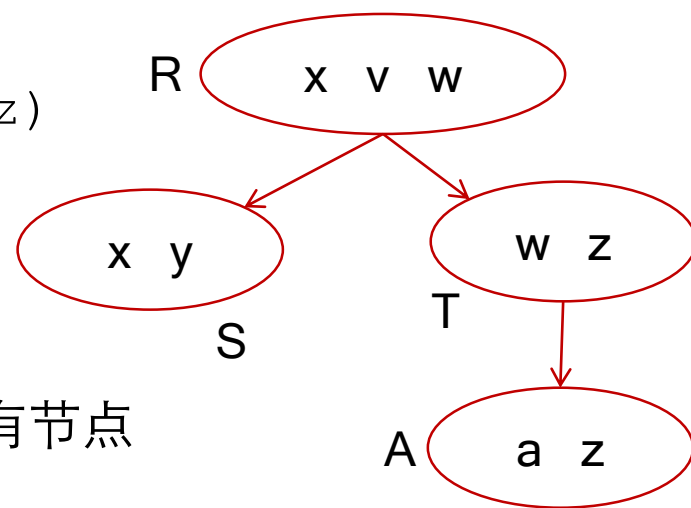
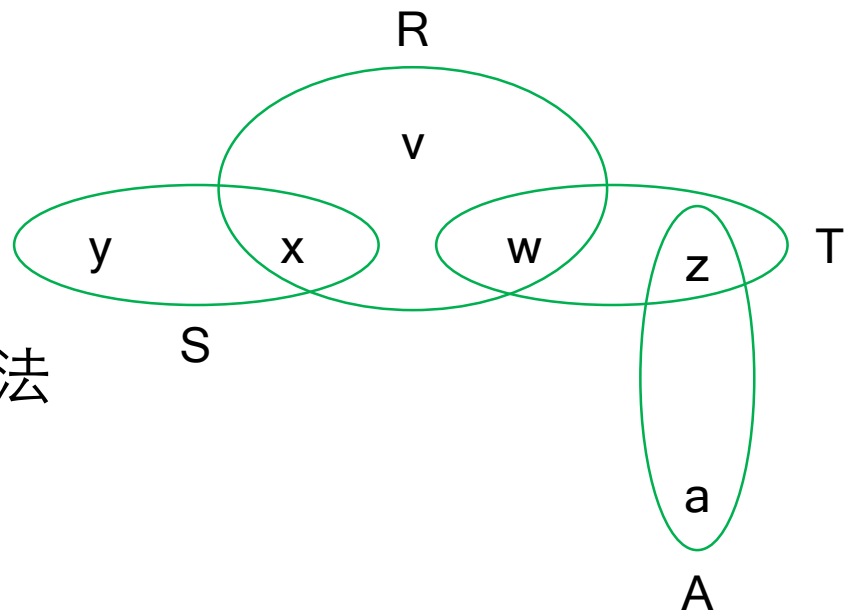
- 超图：一条边可连接多个节点

- 边：关系，节点：属性

$$Q_1(y) :- R(x, v, w), S(x, y), T(w, z), A(a, z)$$

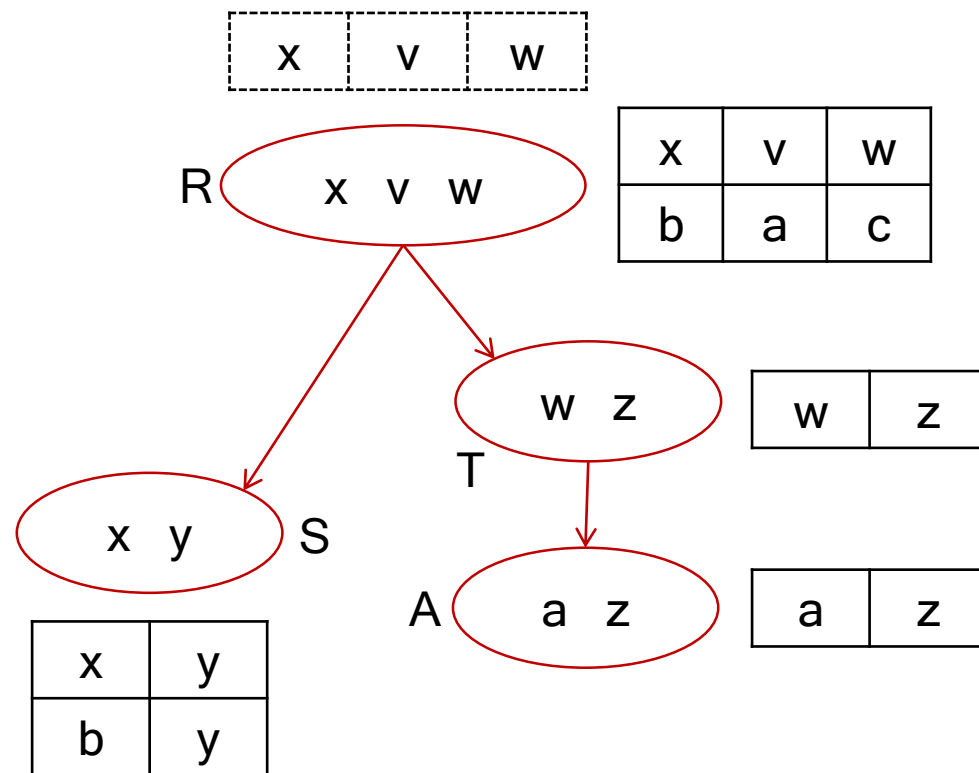
- 构建消除树

- 从超图中不断移除“耳朵”边
 - 耳朵：除和另外一条边的公共节点之外的节点都是私有节点
- 无环：可将超图消除至空



无环CQ查询： 包含问题

- 问题： Q_2 是否包含 Q_1
- 解法：
 - 构建 Q_2 的消除树
 - 使用 Q_1 的主体部分实例化树节点
 - 变量名作为常量
 - 自底向上地计算自然半连接
 - 根节点不为空时， $Q_1 \sqsubseteq Q_2$



$Q_2(y) \text{ :- } R(x, v, w), S(x, y), T(w, z), A(a, z)$

$Q_1(y) \text{ :- } R(x, v, w), S(x, y), T(w, z), A(a, z), S(b, y), R(b, a, c)$

合取查询的并

- 如何处理SQL中的UNION操作?
 - 部分通过OR连接的条件可转化为UNION操作
- 可表示为合取查询的并
 - `SELECT id FROM t WHERE name = 'alice' OR name = 'bob'`
 - $Q(x) :- t(x, \text{'alice'})$
 - $Q(x) :- t(x, \text{'bob'})$
- 定理：令 Q_1 为一个合取查询， $Q_2 = Q_2^1 \cup \dots \cup Q_2^n$ 为一个合取查询的并。 Q_2 包含 Q_1 ，当且仅当存在 $1 \leq i \leq n$ 使得 $Q_1 \sqsubseteq Q_2^i$ 。

等于条件之外

- 以上讨论限制CQ的选择条件为“属性=属性”或“属性=常量”
- 现实中比较 ($=, <, >, \leq, \geq, \neq$) 和非 (\neg) 也很常见
 - $(w) :- R(u, v, w, x), S(x, y), \neg T(y, z), u < 10, v > 5$
- 加入此类条件会显著提高充要条件的问题复杂度
 - $Q_1 \sqsubseteq Q_2 \Leftrightarrow \dots$
- 现实中充分条件足够满足需求
 - $\dots \Rightarrow Q_1 \sqsubseteq Q_2$
- 充分条件：包含映射
 - 比较条件：考虑蕴含关系和等价类

包语义，分组聚合

- 上述讨论基于集合语义（消除重复）
- SQL基于包语义（默认不消除重复）
- 基于包语义的查询包含问题是不可判定的
- 幸运的是：大部分表都有主键
- 聚合函数
 - count, sum基于包语义；max, min基于集合语义
 - 一般处理方法：先考虑底层CQ，再考虑分组和聚合
 - GROUP BY A, B, C可用于计算GROUP BY A, B

外连接 (OUTER JOIN)

- 以上讨论基于内连接 (INNER JOIN)
- 现实中存在大量外连接操作
- 外连接可分解为SPJ型查询的组合 (子查询和UNION)

概要

- 背景
- 查询包含问题
- 使用视图回答查询
- 分析和总结

问题定义

- 给定：
 - 数据库模型 $\mathcal{R} = \{R_1, \dots, R_n\}$
 - 定义在 \mathcal{R} 上的视图定义集合 $\mathcal{V} = \{V_1, \dots, V_m\}$
 - 对 \mathcal{R} 的查询 Q
- 输出查询 Q' （如存在），使得：
 - Q' 的 FROM 子句只引用 \mathcal{V} 中的视图 ← 基础表也可作为视图
 - Q' 等价于 Q
- NP 完全问题

例：使用视图回答查询

- 数据库

- `Movie(ID, title, year), Director(ID, director), Actor(ID, actor)`

- 视图定义

- `V1(I, T, Y) :- Movie(I, T, Y), Y > 1950`
 - `V2(I, D) :- Director(I, D), Actor(I, D)`

- 查询

- `Q(I, T, Y) :- Movie(I, T, Y), Director(I, D), Actor(I, D), Y > 1980`

- 查询重写

- `Q'(I, T, Y) :- V1(I, T, Y), V2(I, D), Y > 1980`

问题解法

- 定理结论：对于有 n 个子目标的合取查询 Q ，无需考虑子目标数超过 n 的重写
 - Naïve解法：枚举所有子目标数不超过 n 的重写，判断重写是否等价于 Q

定理：当存在变量映射 ϕ_i 将 $V_i \in \mathcal{V}$ 的子目标映射到 Q 的子目标时，称视图 V_i 对 Q 是有意义的。令 V_1, \dots, V_k 表示对 Q 有意义的视图， $\mathbf{Y}_i = \mathbf{head}(V_i)$,

$$C(\mathbf{head}(Q)) : - V_1(\phi_1(\mathbf{Y}_1)), \dots, V_k(\phi_k(\mathbf{Y}_k))$$

当且仅当 C 的展开 $C' \sqsubseteq Q$ 时， Q 存在基于 \mathcal{V} 的重写。

- 实际可行的解法：使用上述定理及其它（启发式）规则剪枝。

问题解法（续）

- 桶算法
 - 为 Q 的每个子目标初始化一个桶
 - 将对每个子查询有意义的视图加入其对应的桶中
 - 计算所有桶的笛卡尔积，对产生的每条查询检测其是否与 Q 等价
- 存在其它更高效的算法
 - （下一步工作）

概要

- 背景
- 查询包含问题
- 使用视图回答查询
- 分析和总结

总结和分析

- 上面仅讨论了生成重写的方法
- 但是：
 - 可能存在多种重写
 - 可能直接从源表计算更高效
- 还是需要代价模型
- 进一步细化对具体算法的了解
 - 外连接，查询重写