

# 从函数依赖综合第三范式关系

菲利普·伯恩斯坦

多伦多大学

**摘要：**有人提出，关系数据库的描述可以表述为数据库属性之间的一组函数关系。然后可以使用这些函数关系通过算法合成关系模式。本文的目的是提出一种进行此类合成的有效程序。该过程产生的模式被证明是 Codd 的第三范式，并且包含尽可能少的关系。还讨论了早期尝试构建此类程序的问题。

**关键词：**数据库模式、函数依赖、关系模型、数据语义、第三范式

## 1 简介

对关系数据库模型的研究表明，当人们考虑如何将属性分组为关系时，函数关系是一个重要的概念[3, 8-12]，有人提出可以制定数据库的基本描述纯粹作为一组这样的函数关系，可以通过算法合成关系模式 [3, 12]。本文的目的是开发一种可证明合理且有效的程序，用于从一组给定的函数关系中合成满足科德第三范式的关系。此外，我们的过程合成的模式将显示包含最少数量的关系。

该方法假设最多存在一个将任何一组属性连接到另一组属性的函数关系。所有早期方法都需要这种唯一性假设，它提出了将详细讨论的困难语义问题。

前三节介绍了从函数依赖性合成关系的问题。本文的第 2 节回顾了关系模型、函数依赖的概念和 Codd 范式。引入了两个新概念：“超级密钥”和“体现”。第 3 节概述了一般综合问题，并提出了一种从函数依赖性综合关系的简单算法。第四节介绍了阿姆斯特朗的函数依赖公理化以及对该理论的唯一性假设和其他语义考虑的评论。主要的综合算法在第 5 节中描述。第 6 节检查了应用于由两个版本的综合算法综合的关系的 Codd 的第三范式属性。本节最后介绍了一种用于合成可证明的第三范式关系的新算法。在第 7 节中，合成的模式显示为最小尺寸。

## 2 关系模型

### 2.1 关系

在 Codd 的关系数据库模型中，一组域上的数学关系用于描述数据项之间的连接 [7]。然而，并非所有关系都同样能很好地描述这些联系 [8]。为了判断各类关系的有效性，我们首先回顾与关系模型相关的术语。

从概念上讲，关系是一个表，其中每一列对应一个不同的属性，每一行对应一个不同的实体（或元组）。对于每个属性，都有一组可能的关联值，称为该属性的域。不同的属性共享一个域是很常见的。例如，属性

QUANTITY\_IN\_STOCK 和 SIZE\_OF\_CLASS 均采用来自称为 NONNEGATIVE INTEGERS（非负整数域）的值。

关系中的（实体，属性）条目是与从属性域中选择的实体关联的值。形式上，关系是与关系属性关联的域的笛卡尔积的（有限）子集。用于描述数据库关系结构的符号包括关系名称（例如 R）和 R 中的一组属性（例如

$\{A_1, A_2, \dots, A_n\}$ ），并写作  $R\{A_1, A_2, \dots, A_n\}$ ；例如见图(1a)。属性的顺序并不重要，因为属性名称在关系中是不同的。（这是区分属性和域的原因之一。）从符号上讲，我们将在字母表开头附近使用大写字母表示简单（即单例）属性，在字母表末尾附近使用大写字母表示复合（即组）属性。

组成关系的实体集通常随着实体的插入、删除和修改而变化。这是数据库关系不同于数学关系的重要方式之一。文献中经常使用“关系”一词来描述关系的结构（例如  $R\{A_1, A_2, \dots, A_n\}$ ），称为其静态的意图，以及关系中的元组集合，称为其扩展。在下文中，除非另有明确说明，“关系”一词将指意图。也就是说，我们通常指的是关系的结构，而不是元组本身。

<b>EMPLOYEE</b> ( <u>EMP#</u> , NAME, DEPT#)	$EMP\# \rightarrow NAME$
<b>DEPARTMENT</b> ( <u>DEPT#</u> , MGR#)	$EMP\# \rightarrow DEPT\#$
<b>INVENTORY</b> ( <u>STOCK#</u> , <u>DEPT#</u> , QTY)	$DEPT\# \rightarrow MGR\#$
	$STOCK\#, DEPT\# \rightarrow QTY$
(a)	(b)

**Fig. 1. Relations and functional dependencies: (a) an example of a relational schema (underlined attributes are keys); (b) functional dependencies for the schema of (a).**

## 2.2 函数依赖

正如我们将在后面的部分中看到的，在选择如何将属性分组为关系时，考虑功能关系非常重要。数据库属性之间的函数关系以函数依赖的概念形式化。

令  $A$  和  $B$  为属性，令  $DOM(A)$  为  $A$  的域， $DOM(B)$  为  $B$  的域，令  $f$  为时变函数，使得  $f: DOM(A) \rightarrow DOM(B)$ 。  $f$  不是精确数学意义上的函数，因为我们允许  $f$  的扩展随时间变化，就像我们允许数据库关系的扩展随时间变化一样。也就是说，如果  $f$  被认为是一组有序对  $\{(a, b) \mid a \in DOM(A) \text{ 和 } b \in DOM(B)\}$ ，那么在每个时间点，对于给定的  $a \in DOM(A)$  值，最多有一个  $b \in DOM(B)$  值。为了区分  $f$  和数学函数，我们将  $f$  称为函数依赖（缩写为 FD）。为了符号方便，我们通常省略“DOM”并写为  $f: A \rightarrow B$ 。如果存在 FD  $f: A \rightarrow B$ ，则称  $B$  在功能上依赖于  $A$ 。

上述定义以明显的方式概括为复合属性的函数依赖性。如果  $X = \{A_1, A_2, \dots, A_n \mid \text{且 } Y = \{B_1, \dots, B_m\}$  是属性集，则  $f: X \rightarrow Y$  表示  $f: DOM(A_1) \times \dots \times DOM(A_n) \rightarrow DOM(B_1) \times \dots \times DOM(B_m)$ 。我们通常会在 FD 中省略集合符号并写成  $f: \{A_1, \dots, A_n\} \rightarrow \{B_1, \dots, B_m\}$  简单地表示为  $f: \{A_1, \dots, A_n\} \times \{B_1, \dots, B_m\}$ 。例如，图 1(a) 中属性的函数依赖关系如图 1(b) 所示。

在本文中，我们将假设对于任意两组属性  $X$  和  $Y$ ，至多有一个 FD  $X \rightarrow Y$ 。属性可能需要重命名以保证这一假设。这一限制很重要，将在 4.2 节中详细讨论。稍后我们还将证明，非功能关系不需要满足这种唯一性假设。

考虑到这个假设，如果  $f: A \rightarrow B$  我们经常将  $A \rightarrow B$  写为缩写。符号  $\nrightarrow$  表示不存在感兴趣的  $FDA \rightarrow B$ （尽管在某种关系中的给定时间点， $A$  的任何值都可能不具有多个对应的值  $B$ ）。

设  $R$  为关系， $X$  为  $R$  的子集。如果  $A_i$  中的每个属性  $A_n$  不在  $X$  中，则在功能上依赖于  $X$ ，并且如果  $X$  的子集不具有此属性。显然，一个关系可以有很多键。 $R$  的超级键是  $R$  中包含  $R$  的键的任何属性集。（每个键也是一个超级键。）

引入超级键的概念主要是为了简化后面章节中的证明，

## 2.3 关系操作

在他对关系模型的最初描述中，Codd 引入了关系代数作为关系数据库模型的数据操作语言 [7]。我们会对两种基本的关系代数运算感兴趣：投影和连接。

关系在其属性  $X$  的子集上的扩展的投影是通过删除  $X$  中不存在的属性而获得的元组集合。如果两个元组现在无法区分，因为它们仅在被删除的属性上有所不同，那么它们是“合并”到一个元组中。也就是说，投影的结果必须是与  $X$  的属性相关的域的笛卡尔积的子集。

连接操作用于在不同关系中出现属性之间建立连接。这里我们要考虑的唯一连接操作是自然连接（即相等连接）。域  $B$  上关系  $R(A, B)$  的外延与关系  $S(B, C)$  的外延的自然连接（表示为  $R \bowtie S$ ）定义为。也就是说，它将与公共  $B$  值相关的所有  $A$  和  $C$  值链接在一起。

## 2.4 模式

任何数据模型（无论是关系数据模型还是其他数据模型）的目的都是允许模型的用户描述和操作他打算存储在数据库中的现实世界中的对象之间的关系。在关系模型中，这样的关系集合以关系模式表示。关系模式由一组数据库关系以及每个关系的一个或多个键的规范组成（例如参见图 1(a)）。

我们会说，如果  $X$  是  $R$  的键且  $A$  是  $R$  的任何其他属性，则函数依赖  $X \rightarrow A$  体现在关系  $R$  中。模式中体现的 FD 集合是所有关系中体现的 FD 的并集。模式的关系。

请注意，这种模式表述是 Codd 的 [8] 的修改，其中 FD 作为关系及其键的附加信息给出。我们采用这种修改后的模式概念有几个原因。首先，我们所知的所有数据定义语言都只允许关系和键的规范。其次，我们对模式的定义消除了明确谈论 FD 的需要。FD 凭借我们对键和实施例的定义而隐式存在。第三，我们将看到第三范式组织了一个关系

模式，使得作为外部信息给出的每个FD要么体现在某种关系中，要么可以通过连接和投影操作从体现的FD中恢复。将体现的概念作为一个原始概念可以大大简化随后的理论。

## 2.5 标准化

Codd 观察到，某些关系具有不适合描述数据库的结构属性。这导致他定义了一系列关系的三种范式。

首先，关系值域被排除在关系之外。如果每个域都包含简单值，则关系处于第一范式（缩写为 1NF）。1NF 有两个主要优点[7]。首先，它允许将数据库视为表的集合——一种非常简单且易于理解的结构。其次，它允许定义一小类原始运算符，这些运算符能够操纵关系以获得属性之间的所有必要的逻辑连接。

引入第二范式和第三范式是为了纠正某些函数依赖引起的问题。要检查这些问题，请考虑重新通过加入 DEPT\_\_INV (STOCK#, DEPT#, QTY, MGR#) 获得图 1(a) 属性 DEPT# 上的 DEPARTMENT 和 INVENTORY 关系。将特定 DEPT# 的第一个库存项目插入到 DEPT\_\_INV 的扩展在 DEPT# 和它的 MGR#。删除特定 DEPT# 的最后一个清单项目会丢失该 DEPT# 与其 MGR# 之间的连接。这些副作用称为插入删除异常，仅在插入或删除 DEPT# 的第一个或最后一个元组时发生。此外，如果允许对各个元组进行任意更新，则对于 DEPT# 中的每个 STOCK#，重复 DEPT# 与其 MGR# 之间的连接可能会导致不一致的关系。出现这些问题的原因是 MGR# 在功能上仅依赖于键 STOCK#、DEPT# 的一部分。为了消除 DEPT\_\_INV 的这些问题，必须将 DEPT\_\_INV 放入第二范式。

当某个属性在功能上依赖于一组属性的子集时，就会出现部分依赖性。让  $f$  是函数依赖，其中  $m < n$ 。属性在  $f$  中是无关的，因为  $A_1, \dots, A_m$  足以在功能上确定  $B$ 。在这种情况下，据说  $B$  部分依赖于  $f$ 。如果对于给定的  $f$ ，不存在具有上述属性的  $g$ ，则  $B$  完全依赖于  $f$ 。也就是说， $f$  域中没有无关的属性。

如果属性  $A_j$  出现在  $R$  的任意键中，则称其在  $R$  中为素数。否则，在  $R$  中为非素数。如果在 1NF 中且其每个非素数属性为，则关系处于第二范式（缩写为 2NF）完全依赖于每个键。关系 DEPT\_\_INV (STOCK#, DEPT#, QTY, MGR#) 不在 2NF 中

因为 MGR# 是非主属性，并且部分依赖于键 STOCK#、DEPT#。图 1(a) 中的关系 DEPT 和 INVENTORY 属于 2NF。

现在考虑关系 EMP\_\_DEPT (EMP#, NAME, DEPT#, MGR#) 通过加入 EMPLOYEE 和 DEPARTMENT 关系得到如图 1(a) 关于 DEPT#。尽管 EMP\_\_DEPT 处于 2NF，但它显示相同的问题作为 DEPT\_\_INV。在特定 DEPT# 中插入或删除第一个 EMP# 会产生异常，因为在此过程中创建或销毁了 DEPT#-MGR# 连接。DEPT# 中每个 EMP# 的 DEPT#-MGR# 连接的重复会产生与 DEPT\_\_INV 中相同的一致性问题。在这种情况下，会出现问题，因为 MGR# 在功能上通过属性 DEPT# 依赖于密钥 EMP#。为了消除这个问题，必须将 EMP\_\_DEPT 关系转化为第三范式。

令  $R(A_1, \dots, A_n)$  为关系。如果存在一组属性使得，则属性  $A_i$  传递依赖于属性  $X$ ，与  $A_j$  不是  $X$  或  $Y$  的元素。

如果一个关系的非素数属性都不传递依赖于任何键，则该关系处于第三范式（缩写为 3NF）[8]。3NF 关系也在 2NF 中，因为如果属性  $A_i$  部分依赖于键  $X$ ，则  $A_i$  传递依赖于  $X$ ，因为关系 EMP\_\_DEPT 不在 3NF 中，因为 MGR# 是非素数且是过渡的

完全依赖于密钥 EMP#。考虑到图 1(b) 的 FD，图 1(a) 中的所有关系都属于 3NF（因此也属于 2NF）。范式关系和周围问题的更多例子可以在 [7-9] 中找到，

## 3 综合关系模式

### 3.1 综合问题和非函数关系

Codd 表明，通过对 FD 已知的 1NF 关系应用简单的分解步骤，该关系可以分解为 3NF 中包含所有 FD 的一组关系 [8]，在 [3] 中提出：由于 FD 完全决定了一个关系是否属于 3NF，因此可以选择 FD 作为基本概念并从中构建 3NF 关系。在推进该提案的过程中，提出了一种有效的算法技术来实际从 FD 构建关系。在本文中，我们提出了一种改进的算法，然后讨论了该算法合成的模式的属性。

从 FD 构建关系模式的方法完全依赖于将所有数据关系表示为 FD 的能力。但显然，并不是世界上的每一个逻辑连接都是有效的。尽管如此，我们声称数据库描述中属性之间的所有联系都可以用 FD 来表示。只要连接有效，当然就没有问题。非功能性连接需要特殊处理。

一组属性  $A_1, A_2, \dots, A_n$  之间的非功能性连接  $f$  将被表示为以下 FD:  $f: A_1, A_2, \dots, A_n \rightarrow \theta$ 。  $\theta$  是  $f$  独有的属性；它不会出现在任何其他 FD 中。代表非功能关系的每个 FD 都有其自己的私有属性。所有这  $n$  个属性的基础域是集合  $\{0, 1\}$ 。对于每个元素  $(a_1, a_2, \dots, a_n) \in \text{DOM}(A_1) \times \text{DOM}(A_2) \times \dots \times \text{DOM}(A_n)$ ,  $f(a_1, a_2, \dots, a_n) = 1$  当且仅当  $(a_1, a_2, \dots, a_n)$  在  $f$  下相关。因此， $f$  的扩展完全定义了  $A_1, \dots, A_n$  之间的非函数关系。例如，DRIVER 和

AUTOMOBILE 之间的非功能关系，其中每个 AUTOMOBILE 可以由多个 DRIVER 驱动，每个 DRIVER 可以驱动多个 AUTOMOBILE，由 FD DRIVER,AUTOMOBILE 表示。

请注意，一组属性之间可以存在多个非函数关系，而不会违反 FD 的唯一性假设。例如，我们可以在 DRIVER 和 AUTOMOBILE 之间建立第二个关系来指示所有权： DRIVER,AUTOMOBILE 。通过为每个非函数关系分配唯一的 theta，保留了 FD 的唯一性假设。

这个 9 表示法允许我们将所有非功能关系表示为 FD。综合算法将为这些非函数关系中的每一个生成大约一个关系。在第 5 节中，我们将准确展示每个“非功能性 FD”如何在综合关系模式中体现。

### 3.2 综合问题的形式化

尽管综合问题的动机来自数据库管理，但我们可以用纯粹的符号术语将问题形式化，如下所示。我们得到一组符号（即属性）和一组符号到符号（即 FD）的映射集合 F。问题是找到一个集合  $C = \{C_1, \dots, C_m\}$  的

我们得到以下一组 FD：

我们根据共同的左侧对 FD 进行分组，得到四组：

对于每个组，我们构建一个由组中所有属性组成的关系：

其中带下划线的属性是键。

We are given the following set of FDs:

$$\begin{array}{ll} f_1: A \rightarrow B & f_4: B \rightarrow D \\ f_2: A \rightarrow C & f_5: D \rightarrow B \\ f_3: B \rightarrow C & f_6: ABE \rightarrow F \end{array}$$

We group the FDs according to common left-hand sides, obtaining four groups:

$$\begin{array}{ll} g_1 = \{f_1, f_2\} & g_3 = \{f_5\} \\ g_2 = \{f_3, f_4\} & g_4 = \{f_6\} \end{array}$$

For each group we construct a relation consisting of all of the attributes in the group:

$$\begin{array}{ll} R_1(\underline{A}, B, C) & R_3(\underline{D}, B) \\ R_2(\underline{B}, C, D) & R_4(\underline{A}, \underline{E}, B, F) \end{array}$$

where the underscored attributes are keys.

图 2. 从 FD 导出模式

S 的子集（即关系的集合）以及对于每个  $C_i$  的  $C_i$  子集的集合（即每个关系的键的集合）满足三个属性：首先，F “体现”在 C 中（即关系体现给定的 FD）。其次，每个  $C_i$  不能有传递依赖（即它在 3NF 中）。第三，C 的基数最小。这种对问题的处理仍然有些模糊，因为我们还没有讨论组成 FD 的代数规则。为了激发对这些规则的需求，我们提出了一种简单的合成算法。该算法忽略了代数考虑，并且将被证明是不充分的。

### 3.3 简单的合成过程

从给定的 FD 集合获取关系的一种（过于）简单的方法是将功能上依赖于同一组属性的所有属性组合在一起。这表明了以下过程。首先，将给定的 FD 集合划分为组，使得每组中的所有 FD 具有相同的左侧。然后，为每个组构造一个由该组中出现的所有属性组成的关系。每组中的 FD 的左侧是对应关系的键。例如，见图 2。

从示例中可以看出该方法的几个不良特性。首先，综合关系不属于 3NF。例如，在图 2 的关系  $R_i$  中，C 传递地依赖于密钥 A。在  $R_4$  中，B 部分地依赖于密钥 AE。非标准化关系是由于给定 FD 组中的冗余造成的。稍后我们将看到  $f_2$  是多余的，并且 B 是  $f_6$  中的无关属性。

其次，FD 的左侧不一定是关系的键，尽管它们始终是超键。在  $R_4$  中，ABE 是一个超级密钥，但不是密钥，因为 B 是无关的。

第三，这个过程综合了太多的关系。由于  $f_4$  和  $f_5$  彼此相反，因此关系  $R_3$  是无关的。这是由于在构建  $R_2$  以凭借  $f_5$  将  $D$  识别为第二个密钥时的过程失败，而不是将  $f_5$  放入单独的关系中。

为了解决这些问题，我们首先必须形式化冗余FD的概念。然后我们将返回克服上述困难的合成算法的介绍。

## 4 函数依赖的代数和语义

### 4.1 阿姆斯特朗函数依赖性公理化

Armstrong [1] 给出的 FD 的完整公理化为后面章节中讨论的 FD 代数的研究提供了理论背景。阿姆斯特朗表明，如果给定的 FD 集合存在于关系（的扩展）中，则可以使用公理从给定集合导出的任何 FD 也必须存在。阿姆斯特朗提出了几个等价的 FD 公理化。我们将使用的基于 Delobel 和 Casey [10] 证明的 FD 属性。他们是

A1（自反性） $X \rightarrow X$ 。

A2（增广）如果  $X \rightarrow Z$  then  $X + Y \rightarrow Z$ 。

A3（伪及物性）如果  $X \rightarrow Y$  and  $Y + Z \rightarrow W$  then  $X + Z \rightarrow W$ 。其中符号“+”表示“集合并集”（不一定是相交的集合）。

如果  $R(A, B)$  是关系，则公理 A1 可以应用  $X = \{A, B\}$  来证明  $A, B \rightarrow A, B$  或应用  $X = \{A\}$  来证明  $A \rightarrow A$ 。

A2 的含义很简单，如果  $f: X \rightarrow Z$ ，则可以创建另一个 FD  $g$ ，其中  $g$  的域包括  $X$  以及一些其他无关属性  $Y$ ，其值对所选  $Z$  的值没有影响由  $g$ 。因此，知道  $A \rightarrow A$ ，我们就可以得到  $A, B \rightarrow A$  (i.e.  $X = \{A\}$ ,  $Z = \{A\}$ , and  $Y = \{B\}$ )。

A3 是组成 FD 的替换规则。设  $f: X \rightarrow Y$  和  $g: Y + Z \rightarrow W$ 。该公理声称存在  $h: X + Z \rightarrow W$ 。要了解  $h$  从何而来，请考虑将  $h$  应用于给定的  $x \in \text{DOM}(X)$  和  $z \in \text{DOM}(Z)$  分两步进行。首先，将  $f$  应用于  $x$ ，产生唯一的  $y \in \text{DOM}(Y)$ 。其次，将  $g$  应用于  $y$  和  $z$ ，产生唯一的  $w \in \text{DOM}(W)$ ，从而完成  $h$  的应用。象征性地，我们可以说  $h(x, z)$  被定义为  $g(f(x), z)$ 。另外，请注意，在公理 A3 的陈述中，如果  $Z$  是空集，则伪传递性变为简单传递性。

令  $G$  为 FD 的集合。 $G$  的闭包，表示为  $G^+$ ，被定义为在 A1、A2 和 A3 下闭包的  $G$  的最小超集。对于给定的  $G$ ， $G^+$  可以被证明是唯一的。根据阿姆斯特朗理论我们知道，如果  $G$  是关系  $R$  的给定 FD 集合，则  $G^+$  中的每个 FD 也存在于  $R$  中。

如果  $G^+ = (G - \{g\})^+$ ，则 FD  $g \in G$  在  $G$  中是冗余的。如果  $G^+ = H^+$  并且  $H$  不包含冗余 FD，则  $H$  是给定 FD 集合  $G$  的非冗余覆盖。

引理 1 中阐述了 FD 的一个重要属性，稍后将用于证明许多定理。它基于“推导”的概念，我们将非正式地将其视为阿姆斯特朗公理在给定一组 FD。正式的进展出现在附录中。

引理 1. 令  $G$  为 FD 集合，并令  $g: X \rightarrow Y$  为  $G$  中的 FD。如果  $h: V \rightarrow W$  在  $G^+$  中且  $g$  用于从  $G$  导出  $h$ ，则  $V \rightarrow X$  在  $G^+$  中。

证明: 我们在这里使用推导的非正式概念给出一个直观的论证。附录中给出了使用推导的“推导树”模型的正式证明。我们引入符号  $U \Rightarrow Z$  来表示 FD  $U \rightarrow Z$  可以通过在给定的 FD 集合上应用阿姆斯特朗公理之一来导出。符号  $U \Rightarrow^* Z$  意味着  $U \rightarrow Z$  可以通过使用公理的多个应用来推导。现在引理表明存在使用  $g$  的推导  $V \Rightarrow^* W$ 。也就是说，对于某些（可能为空）属性集  $Z$  存在推导  $V \Rightarrow^* ZX \Rightarrow ZY \Rightarrow^* W$ （步骤  $ZX \Rightarrow ZY$  是使用  $g$  的步骤）。但  $V \Rightarrow^* ZX$  意味着  $V \rightarrow ZX$ ，又意味着  $V \rightarrow X$ ，从而证明了引理。

### 4.2 函数依赖的语义

本文对 FD 的处理是基于阿姆斯特朗公理的严格句法处理。为了使用这种方法，我们必须做出以下唯一性假设：对于给定的一组 FD  $G$  和 FD  $X \rightarrow Y$ ，要么  $X \rightarrow Y$  不在  $G^+$  中，要么在  $G^+$  中存在唯一的 FD  $X \rightarrow Y$ 。也就是说，如果同一组属性上有两个 FD，那么它们就是同一个 FD；如果  $f: X \rightarrow Y$  且  $g: X \rightarrow Y$ ，则  $f$  与  $g$  相同。因此，被接受为综合算法输入的 FD 集合被假定为不仅满足阿姆斯特朗公理，而且满足唯一性假设。（这两个假设也是所有以前的 3NF 句法方法所必需的（例如 [10, 11, 2]））。从几个例子可以看出这个假设是相当有力的。

Let  $f_1: \text{DEPT} \# \rightarrow \text{MGR} \#$  和  $\text{MGR} \#, \text{FLOOR} \rightarrow \text{NUMBER OF EMPLOYEES}$ 。j) 和  $f_2$  的一种解释是  $f_1$  确定每个部门的经理， $f_2$  确定在特定楼层为特定经理工作的员工数量。通过对  $f_1$  和  $f_2$  应用伪及物性，我们得到  $f_3: \text{DEPT} \#, \text{FLOOR} \rightarrow \text{NUMBER\_OF\_EMPLOYEES}$ ，它确定特定楼层的特定部门经理的员工数量。如果一名经理可以管理多个部门，则  $f_3$  与语法上相同的 FD  $g_1$  不同。  $\text{DEPT} \#, \text{FLOOR} \rightarrow \text{NUMBER\_OF\_EMPLOYEES}$ ，确定特定楼层特定部门的员工数量。为了使  $g_1$  与  $f_3$  不同，必须更改属性名称以使 FD 在语法上不同。例如，可以改变  $f_2$  和  $p_1$  使得  $f_1: \text{MGR} \#, \text{FLOOR} \rightarrow \text{NUMBER\_OF\_EMPLOYEES\_OF\_}$

经理和gi: 部门。现在gi与fi和f2的组合不同。

作为第二个示例, 令  $f_4: EMP \rightarrow MGR$  和  $f_5: MGR \rightarrow EMP$ 。这里的情况一定是  $f_4$  是  $f_5$  的倒数。因为如果我们组合  $f_i$  和  $f_6$ , 我们就得到  $g_i: EMP \rightarrow EMP$ 。由于只有一个 FD 将  $EMP$  连接到  $EMP$  (根据我们的假设), 并且根据阿姆斯特朗公理, 恒等函数必须存在, 因此  $g_2$  必定是恒等映射。这意味着  $f_i = U$  如果我们采取这样的解释:  $f_i$  将雇员映射到他的经理, 并且  $f_3$  将经理的  $MGR$  映射到他相应的  $EMP$ , 那么当然  $f_4 f_5$ 。所以要把这个纳入解释时, 必须使  $f_i$  和  $f_5$  在语法上不同 (例如  $f_5: MGR \rightarrow EMP \text{ OF } MGR$ )。

作为第三个示例, 令  $STOCK \rightarrow STORE$  和  $f_i: STOCK, STORE \rightarrow QTY$ 。由于  $f_6$  和  $f_i$  的组合是  $g_3: STOCK \rightarrow QTY$ , 所以 (根据我们的假设)  $f_i$  中的属性  $STORE$  一定是不需要的。但是假设  $f_a$  将  $STOCK$  映射到负责订购该商品的商店的  $STORE$ , 并且  $f_7$  将商品的  $STOCK$  和出售该商品的商店的  $STORE$  映射到现有数量。在这种情况下,  $g_3$  并不意味着  $STORE$  在  $f_7$  中是无关的。为了防止发生这种语法推断, 我们必须更改属性名称 (例如  $f_6: STOCK \rightarrow ORDERING\_STORE$ )。

在每个例子中, 句法推断要么是错误的, 要么是误导性的。在每种情况下, 我们都通过重命名属性以将其与另一个属性区分开来解决问题。这种重命名实质上将我们拥有的关于 FD 的一些语义知识转移到了句法层面, 可供 FD 的代数使用。

指定一组不会导致无效句法推理的 FD 显然是一个难题。因为没有任何语法检查仅基于 FD 的代数可以确定给定的 FD 集合是否满足唯一性假设。然而, 如果我们要使用 FD 的形式代数, 我们必须假设所有句法推论都是有效的。如果我们有一个自动语义分析器可以判断每个句法推理的有效性, 那么我们就可以用它作为筛子来剔除无效的推理。不幸的是, 这样的语义分析器远远超出了现有技术的水平。因此, 我们将在句法推论有效性的假设中添加一个条件, 即所有句法推论都 (或至少可以) 检查语义有效性。如果推论无效, 则可能会导致某些属性的重命名或简单地被拒绝。

第三范式是受 FD 代数控制的严格句法属性。在本文中, 我们给出了从 FD 到 3NF 模式的映射的完整说明, 假设阿姆斯特朗公理和唯一性假设被接受。鉴于阿姆斯特朗的完整性证明, 我们相信这些假设在关系数据库建模中是相当合理的。我们并不是要解决如何判断句法推论的语义有效性的问题。这类语义问题不太容易理解, 似乎比确定 3NT 的句法问题更困难。他们的解决方案仍有待进一步研究。

## 5 更复杂的合成过程

### 5.1 算法描述

3.3 节的简单综合过程导致了问题, 因为忽略了组成 FD 的规则。主要困难在于, 过滤到合成模式中的冗余 FD 会创建额外的属性, 并导致属性之间的非标准化连接。通过首先对给定的 FD 集进行非冗余覆盖, 可以缓解规范化问题。例如, 在图 2 中,  $f_a$  是冗余的, 因此不会出现在给定 FD 的非冗余覆盖中, 从而避免了  $R_i$  的 3NF 违规。

找到非冗余覆盖不足以避免诸如图 2 中的  $f_e$  之类的问题 FD。可以通过从 FD 左侧删除无关属性来消除这个问题。属性  $X_i$  在  $FD \ g \in G, g: X_i, \dots$  中是无关的, 如果  $X_p \rightarrow Y$ , 如果  $X_i, \dots, X_{i-1}, X_{i+1}, \dots, X_p \rightarrow Y$  在  $G$  中。消除无关属性有助于避免部分依赖关系和非键的超键, 如图 2 的  $R_4$  所示。

如果两个关系的键在功能上相互依赖 (即等效), 则这两个关系可以合并在一起。如果左侧功能相同, 则可以通过将两组 FD 合并在一起来完成合成过程。例如, 图 2 中的  $g_2$  和  $g_3$  可以合并为一个组。

以下过程包括上述改进:

算法 1

1. (消除无关属性。) 令  $F$  为给定的 FD 集合。从  $F$  中每个 FD 的左侧消除无关属性, 产生集合  $G$ 。如果某个属性的消除不改变 FD 集合的闭包, 则该属性是无关的。
2.  $G$  的非冗余覆盖  $H$ 。
3. (划分) 将  $H$  划分为组, 使得每组中的所有 FD 具有相同的左侧。
4. (合并等效键。) 对于每对组, 假设  $H_i$  和  $H_j$ , 左侧分别为  $X$  和  $Y$ , 如果  $H$  中存在双射  $X \leftrightarrow Y$ , 则将  $H_i$  和  $H_j$  合并在一起。
5. (构造关系。) 对于每个组, 构造一个由该组中出现的所有属性组成的关系。组中任何 FD 左侧出现的每组属性都是关系的键。(第 1 步保证此类集合不包含任何额外属性。) 该算法找到的所有键都将被称为合成键。构造的关系集构成了给定的 FD 集的模式。

接下来, 我们将删除步骤 4 的算法 1 称为算法 1(a)

[2] 中提出了一种用于测试 FD 集合闭包中的成员资格的线性时间算法。使用此过程，可以在时间界限为  $O(L^2)$  的情况下实现算法 1，其中  $L$  是编码给定 FD 集合的字符串的长度。

## 5.2 综合模式的完整性

如果  $S$  中体现的 FD 的闭包等于  $F^+$ ，则模式  $S$  完全表征了一组 FD  $F$ 。为了表明算法 1 合成了一个完整表征给定 FD 的模式，请考虑作为算法 1 的输入给出的一组 FD  $F$ 。令  $H$  为通过消除无关属性和冗余 FD 产生的 FD 集合。显然， $H^+$  仍然等于  $F^+$ 。设  $S$  是从  $F$  合成的模式。由于  $H$  正是  $S$  中体现的 FD 的集合，并且  $H^+ = F^+$ ，所以  $F$  中的每个 FD 都可以从  $S$  中体现的 FD 的子集导出。因此， $S$  完全表征了  $F$ 。

我们希望确定给定集合  $F$  中每个 FD 的外延都可以使用关系代数从合成模式  $S$  的外延中检索出来。我们认为这是根据  $S$  完全表征  $F$  的事实得出的。

考虑有  $f: X \rightarrow A \in F$ 。我们首先注意到  $X$  中的无关属性可以被忽略。也就是说，如果扩展名为 FD  $f: X \rightarrow A$  其中  $X \supset c X$  可以从  $S$  的扩展中检索，那么由于  $f$  可以简单地通过增广从  $f$  获得，所以我们可以将  $f$  视为与  $f$  相同的 FD。现在，由于  $S$  完全表征了  $F$ ，因此基于  $S$  中体现的 FD  $H$  可以推导  $f$ 。在附录中，我们表明，如果  $f$  没有无关属性，则可以仅使用伪及物性从  $H$  推导它公理。由于伪及物性的应用完全对应于关系代数中的连接，因此  $H$  的  $f$  的推导可以通过  $S$  关系的外延上的一系列连接来模拟。这样，每个  $f \in F$  的外延可以是使用关系代数从  $S$  的扩展中检索。也就是说，我们的“完整表征”概念满足了这样的直觉：给定 FD 集中指定的所有关系实际上都可以从综合模式的扩展中检索。

## 5.3 非功能关系

我们引入了一种特殊的符号来表示输入 FD 中的非函数关系。我们现在必须确保这些 FD 以预期的方式运行。

如果  $X \rightarrow 0$  位于给定算法 1 的 FD 集合中，则  $X \rightarrow 0$  或  $Y \rightarrow 0$ （其中  $Y \leftrightarrow X$ ）出现在算法合成的模式中。这是以下引理的结果，在附录中得到了证明。

引理 2. 如果  $X \rightarrow 0$  在一组 FD  $G$  中，则对于覆盖  $G$  的  $H$  的任何非冗余， $X \rightarrow 0$  在  $H$  中或  $Y \rightarrow 0$  在  $H$  中，其中  $X \rightarrow Y$  和  $Y \rightarrow X$ 。因此，非功能关系以几乎与在给定 FD 集中指定的形式相同的形式出现在模式中。

根据上述引理，为了允许表示非功能关系而发明的  $0$  属性总是出现在合成模式中。它们是如何解释的？要了解这一点，请考虑以下示例。假设在给定的 FD 集中指定了两个非函数关系  $f_1: AB \rightarrow 0_1$  和  $f_2: AB \rightarrow 0_2$ 。（再次注意，FD 的唯一性假设并不强制  $A$  和  $B$  之间非函数关系的唯一性。）Alg 的步骤 4.1 将这两个 FD 合并为一个组，产生关系  $R(A, B, 0_1, 0_2)$ 。为了区分  $A$  和  $B$  的给定值对是否满足  $f_1$ 、 $f_2$  或同时满足  $f_1$  和  $f_2$ ，必须保留  $0_1$  和  $0_2$  属性。例如  $(a, b, 0_1, 1) \in R$  表示  $a, b$  满足  $f_2$  但不满足  $f_1$ 。请注意，如果一组属性之间只有一个非功能关系，则通常可以删除  $0$  属性，因为区分关系的问题消失了。例如，如果只有  $f_1$  存在，则通常只包含在  $f_1$  下相关的  $(a, b)$  对；元组  $(a, b, 0)$  通常不会包含在扩展中。因此在这种情况下， $0$  属性将被完全删除。

# 6 第三范式模式

## 6.1 简介

在本节中，我们将展示各种综合关系在什么条件下属于 3NF。我们首先证明算法 1(a)（即没有步骤 4 的算法 1）总是产生 3NF 模式。然后我们检查算法 1。引入非素数属性的推导性质，并证明它是算法 1 生成 3NF 模式的充分条件。不幸的是，在某些情况下，FD 不满足此属性，因此可能导致具有传递依赖性的关系。一个这样的例子被提出并被证明是德洛贝尔和凯西给出的定理的反例。

## 6.2 算法 1(a) 模式

为了证明算法 1(a) 合成的每个关系都属于 3NF，我们证明传递依赖意味着非冗余覆盖中存在冗余 FD。我们将使用引理 1 来证明产生矛盾的 FD 的存在性。引理 1 将在所有后续的 3NF 证明中以这种方式使用。

定理 1. 令  $R(A_1, \dots, A_n)$  为使用算法 1(a) 从 FD  $F$  集合合成的关系。那么  $R$  的任何非素数属性都不会传递依赖于  $R$  的任何键。也就是说， $R$  在 3NF 中。

证明。假设  $A_i$  是非素数并且传递依赖于  $R$  的密钥  $K$ 。（ $K$  不需要合成。）也就是说，存在一个  $X \subseteq \{A_1, \dots, A_n\}$  使得  $K \rightarrow X$ 、 $X \rightarrow A_i$  和  $X \rightarrow A_i$  位于  $F^+$  中，且  $A_i$  不在  $X$  中。

我们首先观察到  $A_i$  传递依赖于  $R$  的合成密钥。设  $Z$  为出现在用于合成  $R$  的 FD 左侧的  $R$  密钥。 $Z \rightarrow X$  在  $F^+$  中，因为  $Z$  是此外， $X \rightarrow A_i$ ，因为如果  $X \rightarrow Z$ ，则  $X \rightarrow Z$  和  $Z \rightarrow K$  将意味着  $X \rightarrow K$ ，与原始传递依赖中的  $XK$  相矛盾。因

此  $Z \rightarrow X$ 、 $XZ$  和  $X \rightarrow A_i$  也是传递依赖。

令  $H$  为算法 1(a) 中计算的  $G$  的非冗余覆盖。现在我们将证明  $H$  中出现的  $Z \rightarrow A_i$  是多余的。为此，只需证明  $Z \rightarrow X$  和  $X \rightarrow A$  即可；两者都可以从  $H - \{Z \rightarrow A_i\}$  导出。

由于合成  $R$  时使用的唯一  $FD$  的形式为  $Z \rightarrow A_j$ ，因此对于所有  $X_k \in X$ ， $Z \rightarrow X_k$  必定在  $H$  中。由于  $A_i$  不在  $X$  中，因此  $Z \rightarrow X_k$  在  $H - \{Z \rightarrow A_i\}$ 。

假设存在  $X \rightarrow A$  的推导；在  $H$  中使用  $Z \rightarrow A_i$ 。然后根据引理 1，我们有  $X \rightarrow Z$ 。但这违反了  $X \twoheadrightarrow Z$  的传递依赖关系。所以  $X \rightarrow A$ ；必须可以在不使用  $Z \rightarrow A$  的情况下导出；。

由于  $Z \rightarrow X$  且  $X \rightarrow A$ ；都可以从  $H - \{Z \rightarrow A_i\}$  导出，则  $Z \rightarrow A_i$  在  $H$  中必定是冗余的，与  $H$  是非冗余的相矛盾。但这反过来必然意味着传递依赖不存在。

上述定理由 Wang 和 Wedekind 首先提出 [12]；然而他们的证明是不正确的 [4]，在他们的证明中，他们只认为传递依赖可以在  $H$  中导出，而不是  $H - \{Z \rightarrow A_i\}$ 。根据上述证明，他们声称如果  $K \rightarrow X$  和  $X \rightarrow A$ ；是传递依赖，则  $I_v \rightarrow A$ ；可以通过伪传递性导出。他们断言，这违反了  $K \rightarrow A_i$  处于非冗余覆盖层中的事实。然而，只有当我们能够证明  $K \rightarrow X$  和  $X \rightarrow A$  时，后者才是正确的；可以从闭包中导出，无需使用  $K \rightarrow A_i$ 。例如， $G = [A \rightarrow B, B \rightarrow A, A \rightarrow C]$  是一组  $FD$ ，其中  $A \rightarrow B$  和  $B \rightarrow C$  位于闭包中，但  $A \rightarrow C$  不是冗余的，因为  $B \rightarrow C$  不能在只有  $A \rightarrow C$  的情况下从  $G$  导出。无论如何，他们的定理如所述是正确的，并且上述论证使用引理 1 和  $X \twoheadrightarrow K$  中的重要事实纠正了他们的证明。传递依赖。

有趣的是，他们并没有消除算法 1(a) 版本中的超级密钥。这不是一个错误，因为他们明确假设  $FD$  左侧不存在无关属性。然而，不需要做出这一一般假设，因为可以通过算法消除一些无关的属性。事实上，为了与  $FD$  代数完全一致，必须消除这些无关的属性。当然，并不是所有这些无关的属性都可以通过这种方式消除；由于第 4.2 节中讨论的原因，许多语义错误必须由用户负责。

人们可能期望定理 1 的证明能够推广到算法 1 合成的模式。不幸的是，情况并非如此。不属于 3NF 的模式可以通过算法 1 来合成，如图 3(b) 所示。在下一节中我们

将添加一个足以保证算法 1 生成的模式的 3NF 的进一步前提条件。

### 6.3 3NF的充分条件

为了表明没有非素数属性传递依赖于  $R$  的任何键，我们将使用以下属性：

如果以下命题成立，则称属性  $A$  满足关系  $R$  中的属性  $P$ ：令  $H$  为算法 1 的步骤 2 生成的非冗余覆盖。如果  $K \rightarrow A$  在  $H$  中，并且  $K \rightarrow A$  用于合成  $R$ ，那么对于  $R$  的任何素数属性  $B$ ，可以在不使用  $K \rightarrow A$  的情况下导出  $FD K \rightarrow B$ （即可以在  $H - \{K \rightarrow A\}$  中导出）。

属性  $P$  严格来说是  $FD$  派生的句法属性，据我们所知，在现实世界关系方面没有语义解释。这是我们所知道的足以保证算法 1 产生 3NF 关系的最弱性质。性质  $P$  足以保证 3NF 的证明与定理 1 的证明相同。

定理 2. 设  $R(A_1, \dots, A_n)$  是通过使用算法 1 从一组  $FD F$  合成的关系之一。如果  $R$  的所有非素数属性满足属性  $P$ ，则  $R$  在 3NF 中。

证明。设  $A_i$  是  $R$  的非素数属性，它传递依赖于  $R$  的某个键  $Y$ 。也就是说，存在  $Z \subset \{A_1, \dots, A_n\}$  使得  $Y \twoheadrightarrow Z$ ， $Z \twoheadrightarrow A_i$ ，和  $Z \not\rightarrow A_i$ ；与  $A_i$ ；不在  $Z$ 。

令  $H$  为算法 1 中计算的非冗余覆盖。令  $K$  为一个键，使得  $h: K \rightarrow A_i$  位于  $H$  中。也就是说， $h$  是一个  $FD$ ，它带来了  $A_i$ ；通过算法 1 转化为  $R$ 。

由于  $K$  是键，所以  $K \rightarrow Z$ 。此外， $Z \twoheadrightarrow K$ ，因为如果  $Z \rightarrow K$ ，则  $Z \rightarrow K$  且  $K \rightarrow Y$  意味着  $Z \rightarrow Y$ ，这是一个矛盾。所以我们有一个新的传递依赖： $K \rightarrow Z$ 、 $Z \twoheadrightarrow K$  和  $Z \rightarrow A_i$ 。我们要证明  $K \rightarrow Z$  和  $Z \rightarrow A_i$  在  $(H - \{K \rightarrow A_i\})^+$  中，以建立  $H$  是冗余的矛盾。



FDs	Schema synthesized by Algorithm 1
$f_1: X_1, X_2 \rightarrow A$ $f_2: C \rightarrow X_1, X_2$ $f_3: A, X_1 \rightarrow B$ $f_4: B, X_2 \rightarrow C$ A does not satisfy property P, yet $R_1$ is in 3NF. (a)	$R_1(\underline{X_1, X_2}, \underline{C}, A)$ (from $f_1$ and $f_2$ ) $R_2(\underline{A}, \underline{X_1}, B)$ $R_3(\underline{B}, \underline{X_2}, C)$
FDs	Schema synthesized by Algorithm 1
$g_1: X_1, X_2 \rightarrow A, D$ $g_2: C, D \rightarrow X_1, X_2$ $g_3: A, X_1 \rightarrow B$ $g_4: B, X_2 \rightarrow C$ $g_5: C \rightarrow A$ $S_1$ is not in 2NF since A is partially dependent upon the key CD. (b)	$S_1(\underline{X_1, X_2}, \underline{C}, D, A)$ (from $g_1$ and $g_2$ ) $S_2(\underline{A}, \underline{X_1}, B)$ $S_3(\underline{B}, \underline{X_2}, C)$ $S_4(\underline{C}, A)$

Figure 3. A schema violating property P and its decomposition R

图 3. 侵犯财产 P 的示例

令  $Z = \{B_1, B_2, \dots, B_m\}$ 。我们区分两种情况。如果  $B_j$  是素数，则属性 P 保证  $K \rightarrow B_j$  可以在不使用  $K \rightarrow A_j$  的情况下导出。如果  $B_j$  不是素数，则存在 FD  $K' \rightarrow B_j$  将  $B_j$  带入 R。因为  $K \rightarrow K'$  可从  $H - \{K \rightarrow A_j\}$  和 K7 导出（通过属性 P） $\rightarrow B_j$  在  $H - \{K \rightarrow A_j\}$  中，我们得到  $K \rightarrow B_j$  可以在不使用  $K \rightarrow A$  的情况下导出；。因此  $K \rightarrow Z$  在  $(H - \{K \rightarrow A_j\})^+$  中。

现在假设  $Z \rightarrow A_j$  在其推导中使用  $K \rightarrow A_j$ 。然后根据引理 1， $Z \rightarrow K$ ，与传递依赖相矛盾。因此  $Z \rightarrow A_j$  位于  $(H - \{K \rightarrow A_j\})^+$  中。

FD  $K \rightarrow Z$  和  $Z \rightarrow A_j$  在  $(H - \{K \rightarrow A_j\})^+$  中，证明 H 是冗余的，这是一个矛盾。因此传递依赖不可能存在。□

对属性 P 的需求源于算法 1 步骤 4 中等效键的合并。假设  $K_i$  和  $K_2$  在步骤 4 中合并，因为  $K_i \leftrightarrow K_2$ ，并且假设  $K_i \rightarrow Z, Z \rightarrow K_i$ ， $Z \rightarrow A$  是综合关系中的传递依赖。如果  $K_i$  和  $K_2$  是两个独立关系的键，则这种传递依赖性将不存在，就像使用算法 1(a) 的情况一样。产生传递依赖的原因之一是存在一个  $Z_j \in Z$ ，其中  $K_2 \rightarrow Z_j$ ；（和  $K_i \rightarrow A$ ）在非冗余覆盖中，但  $K_i \rightarrow Z_j$ （和  $K_2 \rightarrow A$ ）不在覆盖中。因此，关系必须同时包含  $K_i$  和  $K_2$  以体现传递依赖。传递依赖中的 FD  $K_j \rightarrow Z$  是  $K_i \rightarrow K_2 \rightarrow Z$  的组合。如果 A 不具有属性 P，则可能需要  $K_i \rightarrow A$  才能获得  $K_i \rightarrow K_2$ ，在这种情况下  $K_i \rightarrow A$  不必是冗余的（就像算法 1(a) 中那样）。然而，如果 A 确实具有属性 P，则  $K_i \rightarrow K_2$  不需要  $K_i \rightarrow A$ ，因此  $K_i \rightarrow A$  是多余的，我们有定理。

考虑图 3(a) 中的 FD 集合，它通过算法 1 生成关系  $R_1(X_1, X_2, C, A)$ 。（读者可以检查  $X_1, X_2 \rightarrow C$  处于闭包中）属性 A 在 R 中是非素数，并且不满足属性 P，因为导出  $X_1, X_2 \rightarrow C$  的唯一方法是使用  $X_1, X_2 \rightarrow A$ 。然而，尽管违反了性质 P，关系  $R_1$  仍处于 3NF 状态。因此性质 P 不是 3NF 的必要条件。

图 3(b) 提供了一个 FD 示例，它表现出与图 3(a) 相同的属性 P 违规，但引起部分（因此，传递）依赖性。根据上述关于传递依赖的讨论，我们有  $X_1, X_2 \rightarrow C, C \rightarrow X_1, X_2$  和  $C \rightarrow A$ ，但是  $X_1, X_2 \rightarrow A$  不是冗余的，因为  $X_j, X_2 \rightarrow C$  在其推导中需要  $X_1, X_2 \rightarrow A$ 。

属性 P 影响另一已发布的从 FD 合成关系的过程。Delobel 和 Casey [10] 声称他们的分解过程在某种意义上与我们的算法 1 相当，产生 3NF 关系。然而，他们的主张是不正确的，因为图 3(b) 中的示例歪曲了他们的定理 [6]。

## 6.4 将关系放入3NF

对属性  $P$  的违反可能会导致 3NF 违反。一旦发现特定的违反 3NF 的情况，为了将关系放入 3NF，必须删除违规的依赖关系。足够方便的是，如果非主属性传递依赖于关系的键，则可以简单地从关系中删除该属性，并且生成的模式仍将体现相同的 FD。

定理3.令  $R_i(A_1, \dots, A_n)$  为使用算法1合成的模式  $S = \{R_1, \dots, R_k\}$  中的关系。令  $H$  为 FD 的集合，具体体现在  $S$  中。设  $A_i$  是  $R$  的一个属性，它不会出现在  $R_k$  的任何合成键中，并让  $A_i$  传递依赖于  $R$  的某个键  $K$ 。假设从  $R_k$  中删除  $A_i$ ，产生一个新关系  $R_r$ ，从而产生一个新模式  $S' = \{R_1, \dots, R_r, \dots, R_k\}$ 。那么  $S'$  中体现的 FD 集合的闭包等于  $H^+$ 。

证明。假设  $A_i$  从  $R_k$  中删除。由于  $A_i$  没有出现在算法 1 合成的任何密钥中，因此它的删除只能影响形式为  $f: X \rightarrow A_i$  的具体 FD，其中  $X$  是  $R_k$  的合成密钥。令  $H'$  为  $S'$  中包含的 FD 集合。如果我们证明所有这样的  $f$  都在  $(H')^+$  中，则  $H^+ = (H')^+$ 。

通过与定理 1 和 2 的证明中使用的相同论证，如果  $A_i$  传递依赖于  $R_k$  的任何键，那么它传递依赖于所有键。对于上述形式的每个  $f$ ， $X$  是一个键。因此，对于每个这样的  $X$ ， $X \rightarrow V$ ， $V \neq X$ ，并且  $V \rightarrow A_i$ ， $A_i$  不在  $V$  中，是传递依赖。由于  $A_i$  不在  $V$  中，因此对于每个  $X$ ，即使在  $A_i$  之后，FD  $X \rightarrow V$  仍然体现在  $R_k$  中；已移除。因此每个  $X \rightarrow V$  都在  $H^+$  中。

为了证明  $Y \rightarrow A_i \in (H')^+$ ，我们必须证明  $V \rightarrow A_i$ ；不能使用任何 FD  $X \rightarrow A_i$ ；在其推导过程中。但这是直接遵循的，因为如果  $X \rightarrow A_i$ ；用于导出  $V \rightarrow A_i$ ，然后通过引理 1  $V \rightarrow X$ ，与传递依赖之一中的  $V \rightarrow X$  相矛盾。

由于  $X \rightarrow V$  和  $V \rightarrow A_i$  在  $(H')^+$  中， $f: X \rightarrow A_i$ ；对于所有这样的  $f$ ，都在  $(H')^+$  中。因此  $(H')^+ = H^+$ 。□

定理 3 为我们提供了一种删除不需要的传递依赖的简单方法：即从关系中删除有问题的属性。该定理保证生成的模式仍然体现给定的 FD 集。

传递依赖如此容易地被删除是相当令人惊讶的，因为形成模式的 FD 是非冗余的。看起来删除一个属性应该会导致 FD 的丢失。然而，在从 FD 的非冗余覆盖合成模式时，我们已经隐式地将新的 FD 添加到算法 1 的步骤 4 中的覆盖中，并且这些 FD 显式地体现在模式中。如果  $X$  和  $Y$  是  $H$  中  $X \rightarrow Y$  和  $Y \rightarrow X$  的非冗余覆盖  $H$  中两个不同 FD 的左侧，则  $X$  和  $Y$  被合并并放入单个关系中。这将  $X \rightarrow Y$  和  $Y \rightarrow X$  添加为两个新的 FD，它们明确体现在模式中，即使它们可能没有出现在覆盖中。例如，在图 3(b) 中，FD  $X_1, X_2 \rightarrow C, D$  和  $C, D \rightarrow X_1, X_2$  明确体现在  $S_i$  中，即使前一个 FD 不在覆盖范围内。额外的 FD 的添加使我们能够在不影响具体 FD 的闭包的情况下消除传递依赖。

从不同的角度来看定理 3，我们现在可以修改算法 1 以合成保证符合 3NF 的模式。令  $H$  为算法 1 步骤 2 产生的非冗余覆盖。令  $J$  为所有 FD  $X \rightarrow Y$  的集合，使得  $X$  和  $Y$  是算法 1 步骤 4 中发现的等效键。令  $h: Z \rightarrow A$ ， $h \notin H$ ，以这样的方式体现在  $R_k$  中： $A$ ；出现在  $R_k$  和  $A$  的合成键中，传递依赖于  $R_k$  的键。那么定理 4 说明  $h$  是冗余的；也就是说， $h \in (H + J - \{h\})^+$ 。因此，如果我们消除右侧不在任何合成键中且  $h \in (H + J - \{A\})^+$  的每个 FD  $h: H$ ，那么我们就消除了所有传递依赖。如果存在非质数  $A$ ；如果它传递地依赖于  $R_k$  的密钥，那么定理 4 将保证我们的额外冗余检查将消除它。这引出了我们的主要结果，算法 2，它综合了一个可证明的 3NF 模式。

算法2

1. (消除无关属性。) 令  $F$  为给定的 FD 集合。从  $F$  中每个 FD 的左侧消除无关属性，生成集合  $O$ 。如果属性的消除不会改变 FD 集合的闭包，则该属性是无关的。
2.  $G$  的非冗余覆盖  $H$ 。
3. (划分) 将  $H$  划分为组，使得每组中的所有 FD 具有相同的左侧。
4. (合并等效键。) 让  $f = \emptyset$ 。对于每对组，说  $H_i$  和  $H_j$ ；如果  $H$  中存在平分  $X \leftrightarrow Y$ ，则分别用左侧  $X$  和  $Y$  将  $H_i$  和  $H_j$  合并在一起。对于每个这样的二等分，将  $X \rightarrow Y$  和  $Y \rightarrow X$  添加到  $f$ 。对于每个  $A \in Y$ ，如果  $X \rightarrow A$  在  $H$  中，则将其从  $H$  中删除。对  $H$  中的每个  $Y \rightarrow B$  与  $B \in X$  执行相同操作。
5. (消除传递依赖性。) 找到一个  $H' \subset H$  使得  $(H' + J)^+ = (H + J)^+$  并且  $H'$  的真子集不具有此属性。将  $J$  的每个 FD 添加到其对应的  $H'$  组中。
6. (构造关系。) 对于每个组，构造一个由该组中出现的所有属性组成的关系。组中任何 FD 左侧出现的每组属性都是关系的键。(第 1 步保证此类集合不包含任何额外属性。) 该算法找到的所有键都将被称为合成键。构造的关系集构成了给定的 FD 集的模式。

步骤 1-4 的有效实现与算法 1 相同。步骤 5 可以通过使用 [5] 的隶属度算法有效地实现，然后算法 2 可以在与算法 1 相同的  $O(L^2)$  时间范围内实现。

## 7 极小性证明

本节目的是检查算法 2（或 1）针对给定的 FD 集合成的关系数量，并与体现这些 FD 的任何其他关系模式进行比较。我们将通过证明合成键的等价类数量在给定 FD 集的所有非冗余覆盖中是相同的来证明所有非冗余覆盖生成相同数量的关系。这将意味着算法 2 合成的模式在合成关系的数量上是最少的。

引理 3. 设  $G_1$  和  $G_2$  是两个非冗余 FD 集合，其中  $G_1 = G_2$ 。如果  $g: X \rightarrow A$  在  $G_1$  中，则存在  $h: Y \rightarrow B$ ，其中  $h \in G_1$  且  $Y \rightarrow X$  和  $X \wedge Y \in G_1$ 。

普库夫。如果  $g \in G_1$  且  $G_1 = G_2$ ，则  $g \in G_2$ ，因此  $G_2$  中存在  $g$  的推导。用于导出  $g$  的每个 FD  $h$  都在  $G_2$  中，因此每个这样的  $h$  都在  $G_1$  中。

如果  $h$  在  $G_1$  中的推导中不需要  $g$ ，那么我们可以构造  $g$  在  $G_1$  中的推导。该推导是通过模仿  $G_2$  中  $g$  的推导来构造的，并将该推导中的每个  $h$  替换为  $G_1$  中推导的  $h$ 。由于这个推导没有使用  $g$ ，所以  $g$  在  $G_1$  中必然是多余的，这是一个矛盾。因此，至少有一个  $h$  必须在  $G_1$  中的推导中需要  $g$ 。

假设  $h: Y \rightarrow B$  在  $G_1$  中的推导中需要  $g: X \rightarrow A$ 。然后由引理 1,  $X \rightarrow Y$ 。此外，由于  $h$  出现在  $G_1$  中  $g$  的推导中，根据引理 1  $Y \rightarrow X$ 。因此  $h: Y \rightarrow B$  在  $G_2$  中，其中  $X \rightarrow Y$  和  $Y \rightarrow X$ ，完成证明。□

引理 2 不能被强化，使得  $X = Y$ 。也就是说，可以有两个具有等效闭包的非冗余覆盖，使得这两个覆盖具有代表关键等价类的不同左侧。例如，在图 4 中， $g_3$  和  $h_3$  具有功能相同的左侧，因为  $CE \leftrightarrow DE$ ，但  $CE \wedge DE$ 。

使用引理 3 并认识到算法 2 从具有功能等效左侧的每个最大 FD 组合成一个关系，我们现在可以看到给定 FD 集的所有非冗余覆盖通过算法 2 产生相同数量的关系。

定理 4. 令  $F$  为 FD 的集合。 $F$  的任何两个非冗余覆盖将通过算法 2 产生相同数量的关系。

证明。设  $G_1$  和  $G_2$  是  $F$  的两个非冗余覆盖。根据引理 3，如果 FD  $g: X \rightarrow A$  在  $G_1$  中，则在  $G_2$  中存在  $f: Y \rightarrow B$ ，其中  $X \rightarrow Y$  且  $Y \rightarrow X$ 。因此，对于  $G_1$  中具有功能等效左侧的任何 FD 组， $G_2$  中一定存在一个这样的组，即具有相同功能等效左侧的组。由于每个这样的组生成一个关系，因此  $G_1$  和  $G_2$  必须生成相同数量的关系。□

定理 4 指出，就合成关系的数量而言，所有非冗余覆盖的选择都同样好。这有点令人惊讶，因为它与直觉相矛盾，即最小尺寸的非冗余覆盖可能会比其他较大的非冗余覆盖产生更少的关系。

该定理还表明，在逻辑层面上，对于如何选择覆盖给定 FD 集合的关系，没有太多选择。一些分解方法（例如 [10,11,12]）声称允许系统在一类可能的模式中进行选择，通过效率考虑来指导选择。由于所有覆盖都有相同的键等价类集，因此可能的模式类实际上非常小。因此，如果一个人在逻辑层面上受到规范化考虑而不是效率考虑的引导，那么就会得到一组几乎相同的可能模式。

从定理 4 可以看出，算法 2 生成的关系数量在所有包含相同给定 FD 集合的关系中是最少的。这给出了 [8] 中讨论的最优 3NF 模式的完整特征。

推论 1. 令  $S$  为通过使用算法 2 从一组 FD  $F$  合成的模式。令  $S'$  为包含覆盖  $F$  的一组 FD  $G$  的任何模式。然后  $|S'| \geq |S|$ 。

证明。令  $H \subseteq G$  为  $F$  的非冗余覆盖。当然，通过算法 2， $H$  不会生成比  $S'$  中更多的关系。此外，根据定理 4，算法 2 将从  $G$  生成与从  $F$  生成的关系数量相同的模式。 $|S'| \geq |S|$ 。

$$G = \{g_1: C \rightarrow D, \quad g_2: D \rightarrow C, \quad g_3: CE \rightarrow F\} \quad H = \{h_1: C \rightarrow D, \quad h_2: D \rightarrow C, \quad h_3: DE \rightarrow F\}$$

$G$  and  $H$  are nonredundant and  $G^+ = H^+$ . However,  $g_3$  and  $h_3$  generate different relations. This is an example of Lemma 3 where  $X \neq Y$ .

Fig. 4. Two equivalent coverings with different keys

## 8 结论

本文的目的是开发一种算法，用于从给定的 FD 集合合成 3NF 模式，并检查此类模式的一些属性。主要结果是：

- (1) 某些用于合成模式的简单算法要么产生太多关系，要么违反 3NF。
- (2) 提出了一种合成可证明的 3NF 模式的算法。该算法的本质在于它从给定的 FD 集中消除了尽可能多的冗余。
- (3) 当使用后一种方法时，所有非冗余覆盖都会产生相同数量的关系。因此，综合模式包含最少数量的关系。

据我们所知，这是首次成功尝试可证明且有效地实施 Codd 标准化程序 [8]。（早期两次类似尝试中的错误已被隔离。）此外，根据推论 1，综合关系满足 Codd 的最优性标准——覆盖相同 FD 的其他模式没有更少的关系。

## 9 致谢

我非常感谢 C. Beeri 帮助我重新思考和扩展本文 [3] 的结果。特别是算法 1 的模式和实施例以及步骤 1 的定义均归功于他。

我要感谢 JM Smith 发现了我最初对非功能性关系的处理中的错误，并指出了早期草案中许多不透明的部分。我还要感谢 HA Schmid、K. Sevcik、JR Swenson 和 D. Tsichritzis 对于本研究的许多结果的方法和意义进行了许多有益的讨论。最后，我要感谢审稿人为提高终稿的可理解性提出的许多有益的建议。