

73.014  
C146  
874

# 形式语言 与 自动机

陈崇昕 编著

874  
126

北京邮电学院出版社

# 形式语言与自动机

陈崇昕 编著

北京邮电学院出版社

## 内 容 提 要

本书是讨论形式语言、自动机及语法分析等方面的内容，是属专业基础性教材。

书的内容，除第一章为基础知识外，分两方面：一是形式语言与自动机的基本体系，包括右线性文法与有限自动机、上下文无关文法与下推自动机、图灵机以及无限制文法等。二是形式语言与自动机的重要应用，即语法分析，包括翻译和解析方法。

本书可作为高等学校计算机有关专业高年级学生和硕士研究生学习本课程的教材或教学参考书，也可作为计算机应用领域中广大科技工作者的参考书。

### 形式语言与自动机

陈崇昕 编著

责任编辑 王履榕

北京邮电学院出版社出版

新华书店北京发行所发行 各地新华书店经售

北京市怀柔黄坎印刷厂印刷

850×1168毫米 1/32 印张9.25 字数231.28千字

1988年11月第一版 1988年11月第一次印刷

印数：1-3000册

ISBN 7-5635-0021-9/TP·5 定价：2.25元

## 序 言

本世紀中期，人們在研究使“自然語言”形式化的過程中，提出和發展了形式語言與自動機理論。由於它與計算機科學中的程序語言有著密切的關係，所以隨著計算機科學的飛速發展，形式語言與自動機理論和方法的研究也越來越受到人們的重視，當前已成為計算機科學的理論基礎。應用領域已擴展到圖象處理、模式識別、自動控制以及生物工程等方面。

本書共分七章，第一章是基礎知識，涉及到本書所使用的最基本的數學知識。第二、三、四、五章在介紹語言和文法分類後，重點討論右綫性文法與有限自動機、上下文無關文法與下推自動機、圖靈機以及無限制文法等，這些構成了形式語言與自動機的基本體系。第六、七章是形式語言與自動機理論的重要應用部分，主要討論編譯程序的理論基礎、語法分析以及有關實用方法等。

為了保持理論的系統性，給讀者一個較完整的理論體系，書中討論的重點放在有實用意義的右綫性文法與有限自動機、上下文無關文法與下推自動機，而對圖靈機和無限制文法等僅作扼要地介紹。由於形式語言與自動機是一門應用性較強的學科，本書以第六、七章重點聯系在編譯方面的應用。

本書可作為高等學校計算機有關專業本科學生和碩士研究生學習形式語言與自動機理論的教材或教學參考書，亦可作為廣大計算機科技工作者的參考書。

編 者

1988年6月

# 目 录

## 第一章 基础知识

§ 1.1 集合与关系	( 1 )
§ 1.2 逻辑	( 10 )
§ 1.3 图	( 13 )
§ 1.4 归纳证明	( 24 )
习 题	( 25 )

## 第二章 语言及文法

§ 2.1 语言的定义与运算	( 30 )
§ 2.2 文法	( 33 )
§ 2.3 文法的分类	( 36 )
习 题	( 39 )

## 第三章 有限自动机和右线性文法

§ 3.1 有限自动机	( 41 )
§ 3.2 不确定的有限自动机	( 46 )
§ 3.3 <i>DFA</i> 与 <i>NFA</i> 的等效	( 48 )
§ 3.4 有 $\epsilon$ 转换的不确定的有限自动机	( 53 )
§ 3.5 正则集与正则式	( 59 )
§ 3.6 右线性文法和正则集	( 63 )
§ 3.7 右线性语言与有限自动机	( 66 )
§ 3.8 右线性语言的性质	( 71 )
§ 3.9 双向和有输出的有限自动机	( 84 )
习 题	( 90 )

## 第四章 上下文无关文法与下推自动机

- § 4.1 推导树与二义性..... ( 94 )
- § 4.2 上下文无关文法的变换..... ( 100 )
- § 4.3 *Chomsky* 范式和 *Greibach* 范 式..... ( 112 )
- § 4.4 下推自动机..... ( 117 )
- § 4.5 上下文无关文法与下推自动机..... ( 122 )
- § 4.6 上下文无关语言的性质..... ( 129 )
- § 4.7 受限型上下文无关文法..... ( 138 )
- 习 题..... ( 139 )

## 第五章 图灵机

- § 5.1 基本图灵机..... ( 144 )
- § 5.2 图灵机的构造技术..... ( 149 )
- § 5.3 修改型图灵机..... ( 157 )
- § 5.4 图灵机与无限制文法..... ( 166 )
- § 5.5 线性有界自动机与上下文有关文法..... ( 170 )
- 习 题..... ( 171 )

## 第六章 翻译

- § 6.1 翻译式..... ( 173 )
- § 6.2 转换器..... ( 180 )
- § 6.3 词法分析..... ( 192 )
- § 6.4 句法分析..... ( 197 )
- 习 题..... ( 207 )

## 第七章 解析方法

- § 7.1 回溯解析..... ( 210 )

§ 7.2	$LL(k)$ 文法	( 226 )
§ 7.3	$LR(k)$ 文法	( 247 )
§ 7.4	优先文法	( 268 )
§ 7.5	表格解析	( 274 )
	习 题	( 285 )
	参考文献	( 290 )

# 第一章 基础知识

作为阅读本书的一些基础知识，在本章内引入有关集合、图和逻辑等方面的基本概念。

## § 1.1 集合与关系

### 集合

当我们研究某一类对象时，可把这类对象的整体称为**集合**，组成一个集合的对象称为该集合的**元素**。

设 $A$ 是一个集合， $a$ 是集合 $A$ 中的元素，可表示为 $a \in A$ ，读作 $a$ 属于 $A$ 。如果 $a$ 不是集合 $A$ 的元素，则表示为 $a \notin A$ ，读作 $a$ 不属于 $A$ 。

例如，26个小写英文字母，可组成一个字母集合 $A$ ，每个小写字母皆属于 $A$ ，可写为 $a \in A$ ， $b \in A$ ， $\dots$ ， $z \in A$ 。所有阿拉伯数字都不属于 $A$ ，则有 $2 \notin A$ ， $8 \notin A$ 等。

**注意：**为书写方便，今后对 $a \in A$ ， $b \in A$ ， $\dots$ ， $z \in A$ ，可改写为 $a$ ， $b$ ， $\dots$ ， $z \in A$ 。

有限个元素 $x_1$ ， $x_2$ ， $\dots$ ， $x_n$ 组成的集合，称为**有限集合**。无限个元素组成的集合，称为**无限集合**。例如，整数构成的集合是一个无限集合。

我们把不含元素的集合，称为**空集**，记为 $\phi$ 。

集合的表示法，有列举法和描述法。

**列举法：**是把集合的元素一一列举出来。例如26个小写英文



字母组成的集合 $A$ ，可写成 $A = \{a, b, c, \dots, z\}$ ；阿拉伯数字的集合 $D = \{0, 1, 2, \dots, 9\}$ 以及集合 $C = \{a^1, a^2, a^3, \dots\}$ 等。

**描述法：**是描述出集合中元素所符合的规则。

例如， $N = \{n | n \text{ 是自然数}\}$ ，表明 $N$ 是自然数集合。

$A = \{x | x \in \mathbb{Z} \text{ 且 } 0 \leq x \leq 5\}$ ，其中 $\mathbb{Z}$ 是整数，则

$A = \{0, 1, 2, 3, 4, 5\}$ 。

### 集合之间的关系

(1) 设两个集合 $A$ 、 $B$ 包含的元素完全相同，则称集合 $A$ 和 $B$ 相等，表示为 $A = B$ 。

例如，集合 $A = \{a, b, c\}$ ，集合 $B = \{b, a, c\}$ ，则有 $A = B$ 。

应指出，一个集合中元素排列的顺序是无关紧要的。

对有限集合 $A$ 中不同元素的个数称为集合的**基数**，表示为 $\#A$ 或 $|A|$ 。

例如， $B = \{a, b, c, c, 4, 8\}$ ，其基数 $\#B = 5$ 。

(2) 设两个集合 $A$ 、 $B$ ，当 $A$ 的元素都是 $B$ 的元素，则称 $A$ 包含于 $B$ ，或称 $A$ 是 $B$ 的子集，表示为 $A \subseteq B$ 。当 $A \subseteq B$ 且 $A \neq B$ ，称 $A$ 是 $B$ 的真子集，表示为 $A \subset B$ 。

如果所研究的集合皆为某个集合的子集时，称该集合为**全集**，记为 $E$ 。

(3) 根据(1)和(2)，对于任意两个集合 $A$ 、 $B$ ， $A = B$ 的充要条件是 $A \subseteq B$ 且 $B \subseteq A$ 。

### 幂集

设 $A$ 是集合， $A$ 的所有子集组成的集合称为是 $A$ 的**幂集**，表示为 $2^A$ 或 $\rho(A)$ 。

例如， $A = \{a, b, c\}$ ，则有

$$\rho(A) = \{\phi, \{a\}, \{b\}, \{c\}, \{a, b\}, \{b, c\}, \\ \{a, c\}, \{a, b, c\}\}$$

当 $A$ 是有限集,  $\#A = n$ , 则 $\rho(A)$ 的元素数为

$$C_n^0 + C_n^1 + \dots + C_n^n = 2^n$$

应指出, 空集 $\phi$ 是任何集合的一个子集。

### 集合的运算

(1) 设两个集合 $A$ 、 $B$ , 由 $A$ 和 $B$ 的所有共同元素构成的集合, 称为 $A$ 和 $B$ 的交集, 表为 $A \cap B$ 。

例如,  $A = \{a, b, c\}$ ,  $B = \{c, d, e, f\}$ , 则 $A \cap B = \{c\}$ 。

(2) 设两个集合 $A$ 、 $B$ , 所有属于 $A$ 或属于 $B$ 的元素组成的集合, 称为 $A$ 和 $B$ 的并集, 表示为 $A \cup B$ 。

例如,  $A = \{a, b\}$ ,  $B = \{7, 8\}$ , 则 $A \cup B = \{a, b, 7, 8\}$ 。

(3) 设两个集合 $A$ 、 $B$ , 所有属于 $A$ 而不属于 $B$ 的一切元素组成的集合, 称为 $B$ 对 $A$ 的补集, 表示为 $A - B$ 。

例如,  $A = \{a, b, c, d\}$ ,  $B = \{c, d, e\}$ , 则

$$A - B = \{a, b\}, B - A = \{e\}。$$

设全集 $E$ 和集合 $A$ , 则称 $E - A$ 是集合 $A$ 的补集, 表示为 $\overline{A}$ 。

(4) 设两个集合 $A$ 、 $B$ , 所有序偶 $(a, b)$ 组成的集合, 称是 $A$ 、 $B$ 的笛卡儿乘积, 表示为 $A \times B$ 。

$$A \times B = \{(a, b) \mid a \in A \text{ 且 } b \in B\}$$

例如,  $A = \{a, b, c\}$ ,  $B = \{0, 1\}$ , 则

$$A \times B = \{(a, 0), (a, 1), (b, 0), (b, 1), (c, 0), (c, 1)\}$$

序偶的元素排列是有顺序的, 不能任意颠倒,  $(a, b)$ 和 $(b, a)$ 是不相同的两个序偶, 因此两个序偶相等, 应该是对应元素相同, 例如,  $(a, b) = (c, d)$ , 应有 $a = c$ 和 $b = d$ 。

对任意集合 $A, B, C$ 有如下运算律:

$$(1) A \cup A = A, A \cap A = A;$$

$$(2) A \cup B = B \cup A, A \cap B = B \cap A;$$

- (3)  $(A \cup B) \cup C = A \cup (B \cup C)$ ,  
 $(A \cap B) \cap C = A \cap (B \cap C)$ ;
- (4)  $A \cup (B \cap C) = (A \cup B) \cap (A \cup C)$ ,  
 $A \cap (B \cup C) = (A \cap B) \cup (A \cap C)$ ;
- (5)  $A \cup (A \cap B) = A$ ,  $A \cap (A \cup B) = A$ ;
- (6)  $A \cup \bar{A} = E$ ,  $A \cap \bar{A} = \phi$ ;
- (7)  $\overline{A \cup B} = \bar{A} \cap \bar{B}$ ,  $\overline{A \cap B} = \bar{A} \cup \bar{B}$ ;
- (8)  $E \cup A = E$ ,  $E \cap A = A$ ;
- (9)  $A \cup \phi = A$ ,  $A \cap \phi = \phi$ .

## 关系

关系的概念在数学中是常用的, 诸如大于、小于、等于、包含等都属于关系。下面给出关系的形式定义。

**定义1.1.1** 设 $A$ 是一个集合,  $A \times A$ 的一个子集 $R$ , 称为是集合 $A$ 上的一个二元关系, 简称关系。

对于 $a \in A$ ,  $b \in A$ , 如果 $(a, b) \in R$ , 称 $a$ 和 $b$ 存在关系 $R$ , 表示为 $aRb$ ; 如果 $(a, b) \notin R$ , 称 $a$ 和 $b$ 不存在关系 $R$ , 表示为 $a \not R b$ 。

例如, 自然数集合 $N$ 中的大于关系, 可表示为

$$> = \{(a, b) \mid a, b \in N \text{ 且 } a > b\}$$

当有两个集合 $A$ 、 $B$ , 则从 $A$ 到 $B$ 的关系是 $A \times B$ 的一个子集。

**例1** 设 $A = \{x, y, z\}$ ,  $B = \{0, 1\}$

$$A \times B = \{(x, 0), (x, 1), (y, 0), (y, 1), (z, 0), (z, 1)\}$$

则下列子集均为从 $A$ 到 $B$ 的关系:

$$R_1 = \{(x, 0), (y, 0)\}$$

$$R_2 = \{(x, 1), (y, 1), (z, 0)\}$$

$$R_3 = \{(y, 0)\}$$

**定义1.1.2** 设集合  $A$ ,  $R$  是  $A$  上的关系:

对每个  $a \in A$ , 如果有  $aRa$ , 称  $R$  是自反的;

对于  $a, b \in A$ , 如果有  $aRb$ , 又有  $bRa$ , 称  $R$  是对称的;

对于  $a, b \in A$ , 如果有  $aRb$  和  $bRa$ , 则必有  $a = b$ , 称  $R$  是反对称的;

对于  $a, b, c \in A$ , 如果有  $aRb$  和  $bRc$ , 则有  $aRc$ , 称  $R$  是传递的;

对每个  $a \in A$ , 如果  $a \not R a$ , 称  $R$  是反自反的。

例如, 数之间的相等关系, 具有自反性、对称性和传递性, 小于关系和大于关系没有自反性, 但有传递性。

**定义1.1.3** 设  $R$  是非空集合  $A$  上的一个关系, 如果  $R$  有自反性、对称性和传递性, 则称  $R$  是一个等价关系。

由等价关系  $R$  可以把  $A$  分为若干子集, 每个子集称为一个等价类, 同一等价类中的元素互相是等价的。

**例2** 设  $Z$  是整数集合,  $R$  是  $Z$  上模3同余关系, 也是一个等价关系, 即

$$R = \{(x, y) \mid x, y \in Z \text{ 且 } x \equiv y \pmod{3}\}$$

由于  $R$  是等价关系, 则存在三个等价类为

$$[0]_R = \{\dots, -6, -3, 0, 3, 6, \dots\}$$

$$[1]_R = \{\dots, -5, -2, 1, 4, 7, \dots\}$$

$$[2]_R = \{\dots, -4, -1, 2, 5, 8, \dots\}$$

其中  $[0]_R, [1]_R, [2]_R$  是表示等价类的符号。

### 逆关系

设  $R$  是集合  $A$  上的一个关系, 则

$$R^{-1} = \{(y, x) \mid x, y \in A \text{ 且有 } (x, y) \in R\}$$

则称  $R^{-1}$  是关系  $R$  的逆关系

例如，小于关系的逆关系是大于关系，相等关系的逆关系仍然是相等关系。

### 偏序关系

定义1.1.4 设 $R$ 是集合 $A$ 上的一个关系，如果 $R$ 有自反性、反对称性和传递性，则称 $R$ 是偏序关系（或部分序关系）。

例3 设集合 $C = \{2, 3, 6, 8\}$ ， $R$ 是集合 $C$ 上的整除关系，即

$$R = \{(x, y) \mid x, y \in C \text{ 且 } x \text{ 整除 } y\}$$

可得

$$R = \{(2, 2), (3, 3), (6, 6), (8, 8), (2, 6), (2, 8), (3, 6)\}$$

例4 设集合 $A = \{\phi, \{a\}, \{b\}, \{a, b\}\}$ ，幂集 $\rho(A)$ 上的包含关系 $\subseteq$ ，是一个偏序关系。这里

$$\rho(A) = \{\phi, \{a\}, \{b\}, \{a, b\}\}$$

在 $\rho(A)$ 上的包含关系可用图的方法表示，如图1.1.1所示。

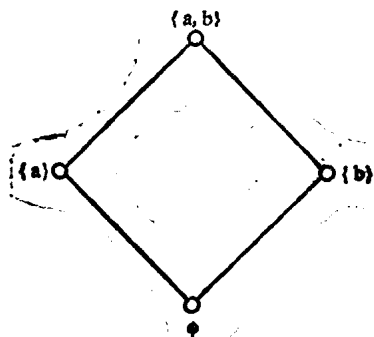


图 1.1.1

描写偏序关系的图称为哈斯图。结合本例说明哈斯图的画法：由于 $\rho(A)$ 中存在 $\{a\} \subseteq \{a, b\}$ ，所以哈斯图中有一条从结点 $\{a\}$ 到结点 $\{a, b\}$ 的边，这条边是自下而上的。又因 $\phi \subseteq \{a\}$ ，故从结点 $\phi$ 到结点 $\{a\}$ 也存在一条自下而上的边。而对于 $\phi \subseteq \{a, b\}$ ，由于以上两条边的存在，靠偏序关系的传递性，所以从结点 $\phi$ 到结点 $\{a, b\}$ 之间的边是不必要的。

性，所以从结点 $\phi$ 到结点 $\{a, b\}$ 之间的边是不必要的。

## 关系的闭包

定义1.1.5 设 $R$ 是集合 $A$ 上的关系, 如果另有关系 $R'$ 满足

(1)  $R'$ 是传递的(自反的, 对称的);

(2)  $R' \supseteq R$ ;

(3) 对任何传递的(自反的, 对称的)关系 $R''$ , 当有 $R'' \supseteq R$ , 就有 $R'' \supseteq R'$ , 则称关系 $R'$ 是 $R$ 的传递(自反, 对称)闭包。

$R$ 的自反闭包表示为 $r(R)$ ,  $R$ 的对称闭包表示为 $s(R)$ 、 $R$ 的传递闭包表示为 $t(R)$ 。

如果给定一个集合 $A$ 上的关系 $R$ , 可用以下方法找出传递闭包 $t(R)$ , 自反闭包 $r(R)$ 和对称闭包 $s(R)$ ;

(1)  $r(R) = R \cup I_A$ , 其中 $I_A = \{(x, x) \mid x \in A\}$ ;

(2)  $s(R) = R \cup R^{-1}$ ;

(3)  $t(R) = R \cup R^2 \cup \dots \cup R^n$ , 其并 $A = n$ 。

例5 设集合 $A = \{a, b, c\}$ ,  $A$ 上的关系

$R = \{(a, b), (b, b), (b, c)\}$ , 则 $R$ 的传递闭包为

$t(R) = \{(a, b), (b, b), (b, c), (a, c)\}$

而 $R$ 的自反传递闭包表示为 $tr(R)$

$tr(R) = \{(a, a), (a, b), (b, b), (b, c), (a, c), (c, c)\}$

今后用 $R^+$ 表示 $R$ 的传递闭包, 用 $R^*$ 表示 $R$ 的自反传递闭包。

## 映射

映射是关系的一个特殊类型, 也称函数。

定义1.1.6 设集合 $A$ 和 $B$ ,  $f$ 是从 $A$ 到 $B$ 的一个关系, 如果对每一个 $a \in A$ , 有唯一的 $b \in B$ , 使得 $(a, b) \in f$ , 称关系 $f$ 是函数, 记为 $f: A \rightarrow B$ 。

如果存在 $(a, b) \in f$ , 则 $a$ 是 $f$ 的自变量,  $b$ 是 $f$ 作用下的象点, 因此 $(a, b) \in f$ 亦可写成 $f(a) = b$ 。

由定义1.1.6可知, 函数有如下特点:

(1) 函数 $f$ 的定义域是 $A$ , 不能是 $A$ 的某个真子集。

(2) 一个 $a \in A$ 只能对应于唯一的一个 $b$ , 或者说 $f(a)$ 是单值的。 $f$ 的值域是 $B$ 的子集, 记为 $R_f$ 。

例6 设集合 $A = \{a, b, c\}$ ,  $B = \{x, y\}$

$$f_1 = \{(a, x), (b, x), (c, y)\}$$

$$f_2 = \{(a, y), (b, y), (c, y)\}$$

$$f_3 = \{(a, y), (b, x), (c, x)\}$$

$$f_4 = \{(a, x), (a, y), (b, x)\}$$

$$f_5 = \{(a, x)\}$$

其中 $f_1, f_2$ 和 $f_3$ 均为函数,  $f_4$ 和 $f_5$ 不是函数, 是关系。

函数的几种特殊类型:

(1) 对于 $f: A \rightarrow B$ 。如果 $f$ 的值域 $R_f = B$ , 即 $B$ 的每一个元素都是 $A$ 中一个或多个元素的象点, 则称 $f$ 是满射的。

例如, 集合 $A = \{a, b, c, d\}$ ,  $B = \{x, y, z\}$ , 如果 $f: A \rightarrow B$ 为:

$$f(a) = x \qquad f(b) = x$$

$$f(c) = y \qquad f(d) = z$$

则 $f$ 是满射的。

(2) 对于 $f: A \rightarrow B$ 。如果 $A$ 中没有两个元素有相同的象点, 则称 $f$ 是入射的。即对于任意 $a_1, a_2 \in A$ :

如果 $a_1 \neq a_2$ , 则有 $f(a_1) \neq f(a_2)$ , 或者

如果 $f(a_1) = f(a_2)$ , 则有 $a_1 = a_2$ 。

例如, 集合 $A = \{a, b\}$ ,  $B = \{x, y, z\}$ , 如果 $f: A \rightarrow B$ 为:  
 $f(a) = x, f(b) = y$ , 则称 $f$ 是入射的。

(3) 对于 $f: A \rightarrow B$ 。如果 $f$ 既是满射的, 又是入射的, 则称 $f$ 是双射的, 或称是一一对应的。

例如, 集合 $A = \{a, b, c\}$ ,  $B = \{1, 2, 3\}$ , 如果 $f: A \rightarrow B$

为

$$f(a) = 3, f(b) = 1, f(c) = 2$$

则称 $f$ 是双射的，或者说是一一对应的。

### 集合的划分

定义1.1.7 设非空集合 $A$ ,  $\Pi = \{\pi_1, \pi_2, \dots, \pi_n\}$ , 其中 $\pi_i \subseteq A$ ,  $\pi_i \neq \phi (i = 1, 2, \dots, n)$ , 如果有 $\bigcup_{i=1}^n \pi_i = A$ 且 $\pi_i \cap \pi_j = \phi (i \neq j)$ , 则称 $\Pi$ 是 $A$ 的划分。其中 $\pi_i$ 是一个划分块。

例如, 集合 $S = \{a, b, c, d\}$ , 考虑下列集合:

$$A = \{\{a, b\}, \{c, d\}\}$$

$$B = \{\{a\}, \{b\}, \{c\}, \{d\}\}$$

$$C = \{\{a\}, \{b, c, d\}\}$$

$$D = \{\{a, b, c, d\}\}$$

$$E = \{\{a, b\}, \{b, c, d\}\}$$

$$F = \{\{a, b\}, \{c\}\}$$

则 $A, B, C$ 和 $D$ 都是 $S$ 的划分,  $E$ 和 $F$ 则不是 $S$ 的划分。

### 集合的基数 (或势)

对于有限集而言, 所谓集合的基数, 即为集合中不同元素的个数。但对于无限集来说, 集合的基数是什么? 两个无限集的基数是否相同呢? 在讨论了函数之后, 我们可以使用一一对应 (双射) 来讨论集合的基数。

定义1.1.8 设有集合 $A, B$ , 如果存在双射函数 $f: A \rightarrow B$ , 则说 $A$ 和 $B$ 有相同的基数, 或者说 $A$ 和 $B$ 等势, 记为 $A \sim B$ 。

显然, 对于有限集合 $A$ 和 $B$ , 称它们有相同的基数, 是指它们的元素个数相同。对于无限集可以看下面的例子。

例7 设偶数集合 $Ne = \{2, 4, 6, \dots\}$ , 定义函数 $f: N \rightarrow Ne$ ,



$N$  为自然数集合。如果对每个  $n \in N$ , 有  $f(n) = 2n$ , 显然  $f$  是从  $N$  到  $N_e$  的双射, 所以存在  $N \sim N_e$ 。

**例7说明**, 一个无限集, 存在着它与其自身的一个真子集有相同的基数。这里  $N_e$  和自然数集合都是无限集。

通常, 我们考虑一个无限集的基数时, 总是看它与自然数集合能否建立一一对应。我们把能与自然数集合建立一一对应的无限集称为**可数集**, 不能与自然数集合建立一一对应的无限集称为**不可数集**。

例如: 整数集合是可数集;

集合  $\{1, 3, 5, 7, \dots\}$  是可数集;

实数集合  $R$  是不可数集;

集合  $\{x | x \in R, 0 < x < 1\}$  是不可数集,

其中  $R$  是实数。

## § 1.2 逻辑

### 命题与连接词

命题是一个能判断真假的语句, 一般可用一个大写英文字母表示一个命题。例如下列语句皆为命题:

$P$ : 3是奇数

$Q$ : 铜是金属

$R$ : 1加4是2

可见, 命题  $P$  和命题  $Q$  的真值均为真 ( $T$ ), 命题  $R$  的真值为假 ( $F$ )。

连接词是用于把命题构成复合的命题, 连接词包括“非”、“与”、“或”和“蕴含”。通常用符号“ $\neg$ ”表示“非”, 符号“ $\wedge$ ”表示“与”、符号“ $\vee$ ”表示“或”和符号“ $\rightarrow$ ”表示“蕴含”。下面用真值表的方法, 给出这些连接词的定义, 见表1.2.1。