

第一章 绪论

2022-02-26 14:45

一、数据结构研究的内容

★ 通常用计算机解题的步骤：

【一、具体问题抽象为数学模型】

实质：分析问题—>提取操作对象—>找出操作对象之间的关系—>用数学语言描述（数据结构）

【二、设计算法】

【三、编程、调试、运行】

★ 应用举例：

例一、学生学籍管理系统

▪ 例1 学生学籍管理系统

表1.1 学生基本信息表

学号	姓名	性别	籍贯	专业
60214201	杨阳	男	安徽	计算机科学与技术
60214205	蒋林	男	福建	计算机科学与技术
60214215	王诗萌	女	吉林	计算机科学与技术
60214216	冯子晗	女	山东	计算机科学与技术

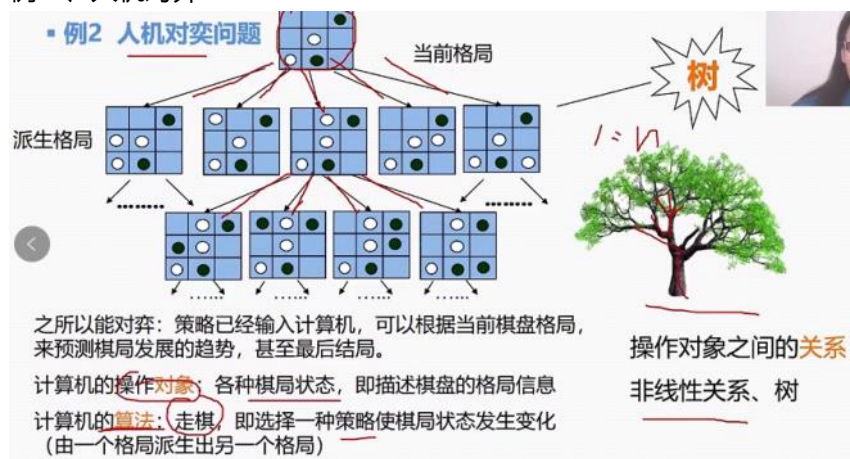
操作对象：每位学生的信息（学号、姓名、性别、籍贯、专业...）。

操作算法：查询、插入、修改、删除等。

操作对象之间的关系：线性关系 数据结构：线性数据结构、线性表。

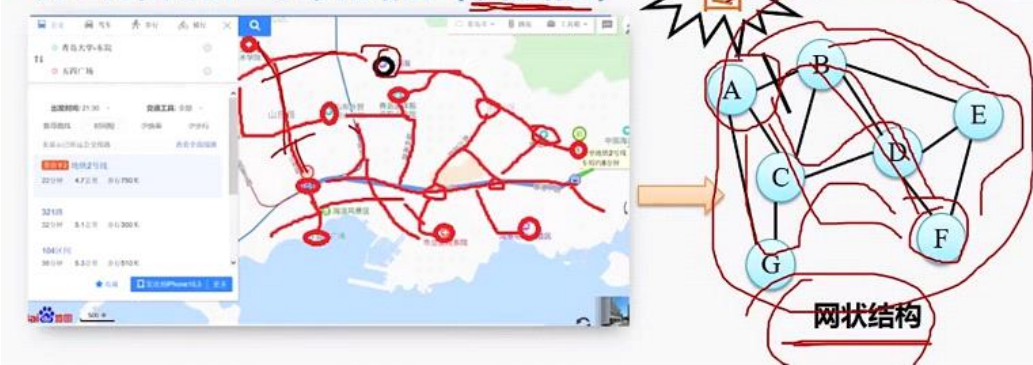
例二、人机对弈

▪ 例2 人机对弈问题



例三、地图导航—求最短路径（最快路径）

例3 地图导航—求最短路径（最快路径）



★ 小结：

- 这些问题的共性都是无法用数学的公式或方程来描述，是一些“非数值计算”的程序设计问题。
- 描述非数值计算问题的数学模型不是数学方程，而是诸如表、树和图之类的具有逻辑关系的数据。
- 数据结构是一门研究非数值计算的程序设计中计算机的操作对象以及它们之间的关系和操作的学科。

二、基本概念和术语

★ 基本概念

数据：

是能输入计算机且能被计算机处理的各种符号的集合

信息的载体

是对客观事物符号化的表示

能够被计算机识别、存储和加工

包括：

数值型的数据：整型、实数等

非数值型的数据：文字、图像、图形、声音等

数据元素：

是数据的基本单位，在计算机程序中通常作为一个整体进行考虑和处理。

也简称为元素，或称为记录、结点或顶点

一个数据元素可由若干个数据项组成

数据项：

构成数据元素的不可分割的最小单位

学籍表				
学号	姓名	性别	出生日期	政治面貌
0001	陆宇	男	986/09/02	团员
0002	李明	男	1985/12/25	党员
0003	汤晓影	女	1986/03/26	团员

数据项

数据元素

数据、数据元素、数据项三者之间的关系：

数据 > 数据元素 > 数据项

例：学生表 > 个人纪录 > 学号、姓名...

数据对象：

是性质相同的数据元素的集合，是数据的一个子集。

例如：

整数数据对象是集合 $N=\{0, \pm 1, \pm 2, \dots\}$

字母字符数据对象是集合 $C=\{ 'A' , 'B' , \dots , 'Z' \}$

学籍表也可看作一个数据对象

数据元素与数据对象

数据元素——组成数据的基本单位

与数据的关系：是集合的个体

数据对象——性质相同的数据元素的集合

与数据的关系：集合的子集

★ 数据结构：

概念：

数据元素不是孤立存在的，它们之间存在着某种关系，数据元素相互之间的关系称为结构

是指相互之间存在一种或多种特定关系的数据元素集合

或者说，数据结构是带结构的数据元素的集合。

数据结构包括以下三个方面的内容：

- 1.数据元素之间的逻辑关系，也称逻辑结构
- 2.数据元素及其关系在计算机内存中的表示（又称为映像），称为数据的物理结构或数据的存储结构。
- 3.数据的运算和实现，即对数据元素可以施加的操作以及这些操作在相应的存储结构上的实现。

数据结构的两个层次

逻辑结构：

描述数据元素之间的逻辑关系

与数据的存储无关，独立于计算机

是从具体问题抽象出来的数学模型

物理结构（存储结构）：

数据元素及其关系在计算机存储器中的结构（存储方式）

是数据结构在计算机中的表示

逻辑结构与存储结构的关系：

存储结构是逻辑关系的映像与元素本身的映像

逻辑结构是数据结构的抽象，存储结构是数据结构的实现

两者综合起来建立了数据元素之间的结构关系。

逻辑结构的种类：

划分方法一：

(1) 线性结构

有且仅有一个开始和一个终端结点，并且所有结点都最多只有一个直接前趋和直接后继

例如：线性表、栈、队列、串

(2) 非线性结构

一个结点可能有多个直接前趋和直接后继

例如：树、图

划分方法二：

- (1) 集合结构：结构中的数据元素之间除了同属于一个集合的关系外，无其他任何关系
- (2) 线性结构：结构中的数据元素之间存在着一对一的线性关系。
- (3) 树形结构：结构中的数据元素之间存在着一对多的层次关系。
- (4) 图形结构或网状结构：结构中的数据元素之间存在着多对多的任意关系。

物理结构的种类：

顺序存储结构

用一组连续的存储单元依次存储数据元素，数据元素之间的逻辑关系由元素的存储位置来表示。

链式存储结构

用一组任意的存储单元存储数据元素，数据元素之间的逻辑关系用指针来表示。

索引存储结构

散列存储结构

三、抽象数据类型的表示与实现

在使用高级程序设计语言编写程序时，必须对程序中出现的每个变量、常量或表达式，明确说明它们所属的数据类型。如C语言中char，float类型，一些最基本数据结构可以用数据类型来实现，如数值、字符串等，二另一些常用的数据结构，如栈、队列、树、图等，不能直接用数据类型来表示。高级语言中的数据类型明显或隐含地规定了在程序执行期间变量和表达的所有可能的取值范围。

由此可以看出数据类型的作用：1.约束变量或常量的取值范围 2.约束变量或常量的操作

★ 定义：

数据类型：是一组性质相同的值的集合以及定义于这个值集合上的一组操作的总称

数据类型=值的集合+值集合上的一组操作

抽象数据类型：是指一个数学模型以及定义在次数学模型上的一组操作。（ADT）

由用户定义，从问题抽象出数据模型（逻辑结构）

还包括定义在数据模型上的一组抽象运算（相关操作）

不考虑计算机内的具体存储结构于运算的具体实现算法

★ 使用：

抽象数据类型的形式定义：

抽象数据类型可用 (D, S, P) 三元组表示

其中：D是数据对象

S是D上的关系集

P是对D的基本操作集

基本操作定义格式说明：

参数表：赋值参数 只为操作提供输入值。

引用参数以&打头，除可提供输入值外，还将返回操作结果。

初始条件：描述操作执行之前数据结构和参数应满足的条件，若不满足则操作失败，并返回相应出错信息。

若初始条件为空，则省略之。

操作结果：说明操作正常完成之后，数据结构的变化状况和应返回的结果。

抽象数据类型定义举例：circle的定义

ADT 抽象数据类型名

```
{
    Data
        数据对象的定义
        数据元素之间逻辑关系的定义
    Operation
        操作 1
            初识条件
            操作结果描述
        操作 2
        ....
        操作 n
        ....
}ADT 抽象数据类型名
```

ADT Circle

```
{
    数据对象:  $D = \{r, x, y \mid r, x, y \text{ 均为实数}\}$ 
    数据关系:  $R = \{ \langle r, x, y \rangle \mid r \text{ 是半径}, \langle x, y \rangle \text{ 是圆心坐标} \}$ 
    基本操作:
        Circle(&C, r, x, y)
            操作结果: 构造一个圆。
        double Area©
            初识条件: 圆已存在。
            操作结果: 计算面积。
        double Circumference©
            初始条件: 圆已存在。
            操作结果: 计算周长。
}ATD Circle
```

★ 抽象数据类型如何实现:

C语言实现抽象数据类型: 用已有数据类型定义描述它的存储结构; 用函数定义描述它的操作
抽象数据类型可以通过固有的数据类型 (如整型、实型、字符型等) 来表示和实现。

```
Void assign (Complex *A, float real, float imag){
    A->realpart = real; /* 实部赋值 */
    A->imagpart = imag; /* 虚部赋值 */
} /* End of assign() */

void add (Complex *c, Complex A, Complex B) { /* c = A + B */
    c->realpart = A.realpart + B.realpart; /* 实部相加 */
    c->imagpart = A.imagpart + B.imagpart; /* 虚部相加 */
} /* End of Add() */
.....
```

★ 概念小结:



四、算法和算法分析

★ 算法:

算法的定义

对特定问题求解方法和步骤的一种描述, 它是指令的有限序列。其中每个指令表示一个或多个操作。
简而言之, 算法就是解决问题的方法和步骤。

算法的描述

自然语言：英语、中文

流程图：传统流程图、NS流程图

伪代码：类语言（类C语言）

程序代码：C语言程序、Java程序

算法与程序

算法是解决问题的一种方法或一个过程，考虑如何将输入转换成输出，一个问题可以用多种算法。

程序是用某种程序设计语言对算法的具体实现。

程序=数据结构+算法

数据结构通过算法实现操作

算法根据数据结构设计程序

算法的特性：

有穷性、确定性、可行性、输入、输出

算法设计的要求：

正确性、可读性、健壮性、高效性

★ 算法分析

算法在满足前面的特性的前提下，主要考虑算法的效率，通过算法的效率高低来评判不同算法的优劣程度。算法效率主要从以下两个方面考虑：1.时间效率（时间复杂度）2.空间效率（空间复杂度）但这两个方面有时候是矛盾的。

算法时间效率的度量

算法时间效率可以用依据该算法编制的程序在计算机上执行所消耗的时间来度量。度量的方法有两种一种为事后统计，一种为事前分析，但事前分析需要编写程序来实现，会耗费更多的时间与精力并且所得结果依赖于计算机的软硬件等环境因素，掩盖算法本身的优劣。

因此我们在判断一个算法的时间效率时，通过事前分析相对来说更好。

事前分析判断算法时间效率：

一个算法的运行时间是指一个算法在计算机上运行所耗费的时间大致可以等于计算机执行一种简单操作（如赋值、比较、移动等）所需要的时间与算法中进行的简单操作次数乘积

算法运行时间=一个简单操作所需要的时间×简单操作次数

= \sum 每条语句的执行次数×该语句执行一次所需要的时间

一般情况下我们假设执行每条语句所需要的时间均为单位时间，此时对算法的运行时间的讨论就可转化为讨论该算法中所有语句的执行次数，即频度之和。

例一：两个 $n \times n$ 矩阵相乘的算法可描述为：

```

1 for(i=1;i<=n;i++)           // n+1次
2   for(j=1;j<=n;j++){        // n(n+1)次
3     c[i][j]=0;               // n*n次
4     for(k=0;k<=n;k++)        // n*n*(n+1)次
5       c[i][j]=c[i][j]+a[i][k]*b[k][j]; // n*n*n次
6   }

```

我们把算法所耗费的时间定义为该算法中每条语句的频度之和，则上述算法的时间消耗 $T(n)$ 为：

$$T(n) = 2n^3 + 3n^2 + 2n + 1$$

这是一个关于 n 的函数

为了便于比较不同算法的时间效率，我们仅比较它们的数量级即可

若有某个辅助函数 $f(n)$ ，使得当 n 趋近于无穷大时， $T(n)/f(n)$ 的极限值为不等于零的常数，则称 $f(n)$ 是 $T(n)$ 的同数量级函数。记作 $T(n) = O(f(n))$ ，称 $O(f(n))$ 为算法的渐近时间复杂度，简称时间复杂度。

对于上述问题， $T(n) = 2n^3 + 3n^2 + 2n + 1$ 。当 $n \rightarrow \infty$ 时， $T(n)/n^3 \rightarrow 2$ ，这表示 n 充分大时， $T(n)$ 与 n^3 是同阶或同数量级，引入大“O”记号，则 $T(n) = O(n^3)$ 。

且一般情况下，不必计算所有操作的执行次数，而只考虑算法中基本操作执行的次数。

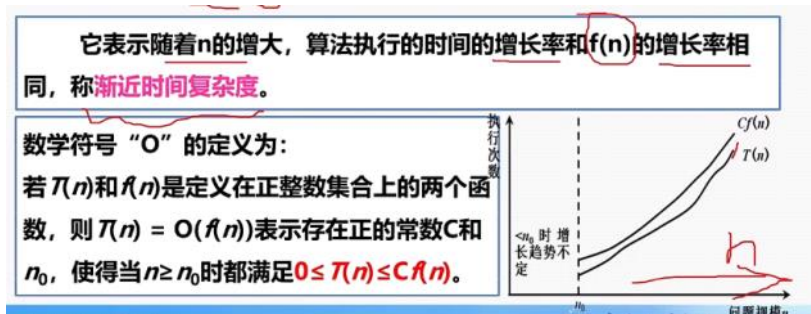
算法时间复杂度定义

算法中基本语句重复执行的次数是问题规模 n 的某个函数 $f(n)$ ，算法时间量度记作： $T(n) = O(f(n))$

- 算法中重复执行的次数和算法的执行时间成正比的语句
- 对算法运行时间的贡献最大
- 执行次数最多

n 越大算法执行时间越长

- 排序： n 为记录数
- 矩阵： n 为矩阵的阶数
- 多项式： n 为多项式的项数
- 集合： n 为元素个数
- 数： n 为树的结点个数
- 图： n 为图的顶点数或边数



分析算法时间复杂度的基本方法

定理1.1

忽略所有低次幂和最高次幂系数，体现出增长率的含义。

方法：

1. 找出语句频度最大的那条语句作为基本语句
2. 计算基本语句的频度得到问题规模 n 的某个函数 $f(n)$
3. 取其数量级用符号“O”表示

例题：

```

x = 0; y = 0;
for ( int k = 0; k < n; k ++ )
    x ++;
for ( int i = 0; i < n; i ++ )
    for ( int j = 0; j < n; j ++ )
        y ++;

```

$T(n) = O(n^2)$

$f(n) = n(n+1)$

```

void exam ( float x[ ][ ], int m, int n ) {
    float sum [ ];
    for ( int i = 0; i < m; i ++ ) {
        sum[i] = 0.0;
        for ( int j = 0; j < n; j ++ )
            sum[i] += x[i][j];
    }
    for ( i = 0; i < m; i ++ )
        cout << i << " : " << sum [i] << endl;
}

```

$T(n) = O(m*n)$

$f(n) = m*n$

时间复杂度是由嵌套最深层语句的频度决定的

例2:

```

for( i=1; i<=n; i++)
    for( j=1; j<=i; j++)
        for( k=1; k<=j; k++)
            x=x+1;

```

语句频度 = $\sum_{i=1}^n \sum_{j=1}^i \sum_{k=1}^j 1 = \sum_{i=1}^n \sum_{j=1}^i j = \sum_{i=1}^n \frac{i(i+1)}{2}$

例3: 分析以下程序段的时间复杂度

```

i=1; ①
while(i<=n)
    i=i*2; ②

```

关键是要找出来执行次数x与...关系, 并表示成n的函数

若循环执行1次: $i = 1 * 2 = 2$
 若循环执行2次: $i = 2 * 2 = 2^2$
 若循环执行3次: $i = 2 * 2 = 2^3$,
 若循环执行x次: $i = 2^x$

设语句②执行次数为x次, 由循环条件 $i \leq n$, $\therefore 2^x \leq n \therefore x \leq \log_2 n$

$2^{f(n)} \leq n$

即 $f(n) \leq \log_2 n$, 取最大值 $f(n) = \log_2 n$

注意:

在有的情况下, 算法中基本操作重复执行的次数还随问题的输入数据集不同而不同。

【例】顺序查找, 在数组a[i]中查找值等于e的元素, 返回其所在位置。

```

for (i=0; i<n; i++)
    if (a[i]==e) return i+1; //找到, 则返回是第几个元素
return 0;

```

- 最好情况: 1次
- 最坏情况: n
- 平均时间复杂度为 $O(n)$

一般法则:

1) for循环——一次for循环运行时间至多是该for循环内语句 (包括测试) 的运行时间乘以迭代次数

- 2) 嵌套的for循环——该语句的运行时间乘以该组所有的for循环的大小的乘积
- 3) 顺序语句——各个运行时间求和即可，其中最大值就是所得的运行时间
- 4) if/else语句——其从不超过判断的时间加上S1和S2中运行时间较长者的总的运行时间
- 5) while语句——对数

几个新概念：

最坏时间复杂度：指在最坏的情况下，算法的时间复杂度。——>一般考虑它

平均时间复杂度：指在所有可能输入实例在等概率出现的情况下，算法的期望运行时间。

最好时间复杂度：指在最好情况下，算法的时间复杂度。

对于复杂的算法，可以将它分成几个容易估算的部分，然后利用大O加法法则和乘法法则。

a) 加法法则： $T(n) = T_1(n) + T_2(n) = O(f(n)) + O(g(n)) = O(\max(f(n), g(n)))$

b) 乘法法则： $T(n) = T_1(n) \times T_2(n) = O(f(n)) \times O(g(n)) = O(f(n) \times g(n))$

渐近空间复杂度：

算法所需存储空间的度量，记作 $S(n) = O(f(n))$

算法要占据的空间：

1. 算法本身要占据的空间，输入/输出，指令，常数，变量等
2. 算法要使用的辅助空间。

★ 总结：

设计好算法的过程

