

# 数据结构与算法 II 上机实验 (11.3)

中国人民大学 信息学院 崔冠宇 2018202147

注：请使用 C++17 或以上编译器！

上机题 1 实现最长公共子序列 (LCS) 算法。

## 一、问题描述

实现最长公共子序列算法。

1. 输入：字符串  $X = \langle x_1, x_2, \dots, x_m \rangle$  与字符串  $Y = \langle y_1, y_2, \dots, y_n \rangle$ 。
2. 输出：它们的最长公共子序列  $Z = \langle z_1, z_2, \dots, z_l \rangle$ 。

## 二、算法基本思路

记  $count[i][j]$  为  $X[1\dots i]$  与  $Y[1\dots j]$  的最长公共子序列的长度。根据课上的分析（定理 15-1），总共有下列几种情况：

1. 若两字符串的当前字符相同，即  $X[i] = Y[j]$ ，则两字符串截至此处的最长公共子序列为“两字符串往前缩一个字符的最长公共子序列加上当前字符”，即

$$count[i][j] = count[i-1][j-1] + 1$$

2. 若两字符串的当前字符不同，即  $X[i] \neq Y[j]$ ，则两字符串截至此处的最长公共子序列为“ $X$  往前缩一个字符， $Y$  不动的最长公共子序列”和“ $Y$  往前缩一个字符， $X$  不动的最长公共子序列”中较大的一个，即

$$count[i][j] = \max\{count[i-1][j], count[i][j-1]\}$$

为了重构解，还需要定义  $dir[i][j]$  记录每一次计算  $count[i][j]$  是上述何种情况：

1. 如果是上述第一种情况，则  $dir[i][j] = LEFTUP$ ；
2. 如果是上述第二种情况，且“ $X$  往前缩一个字符， $Y$  不动的最长公共子序列”较长，则  $dir[i][j] = UP$ ；
3. 如果是上述第二种情况，且“ $Y$  往前缩一个字符， $X$  不动的最长公共子序列”较长，则  $dir[i][j] = LEFT$ 。

打印解的时候，首先根据  $dir[i][j]$  的方向递归打印两字符串前缀对应的解，若  $dir[i][j] = LEFTUP$ ，则还要打印  $X[i](Y[j])$ 。

综上，得到状态转移方程

$$count[i][j] = \begin{cases} count[i-1][j-1] + 1, & X[i] = Y[j] \\ \max\{count[i-1][j], count[i][j-1]\}, & X[i] \neq Y[j] \end{cases}$$

$$dir[i][j] = \begin{cases} LEFTUP, & X[i] = Y[j] \\ UP, & X[i] \neq Y[j] \text{ 且 } count[i-1][j] \geq count[i][j-1] \\ LEFT, & X[i] \neq Y[j] \text{ 且 } count[i-1][j] < count[i][j-1] \end{cases}$$

下面给出 LCS 的伪代码：

```

1 // 最长公共子序列
2 // 参数：
3 //      X, Y: 字符串
4 // 返回值：
5 //      count[i][j]: X[1...i] 与 Y[1...j] 的最长公共子序列的长度
6 //      dir[i][j]: 上述最长公共子序列由哪个前缀确定
7 LCS(X, Y):
8     m = X.length()
9     n = Y.length()
10    count[0...m, 0...n], dir[1...m, 1...n]
11    // 初始化
12    for i = 0 to m:
13        count[i, 0] = 0
14    for i = 0 to n:
15        count[0, i] = 0
16    // 填表
17    for i = 1 to m:
18        for j = 1 to n:
19            // 情况一
20            if X[i] == Y[j]:
21                count[i, j] = count[i - 1, j - 1] + 1
22                dir[i, j] = LEFTUP
23            // 情况二
24            else if count[i - 1, j] >= count[i, j - 1]:
25                count[i, j] = count[i - 1, j]
26                dir[i, j] = UP
27            else:
28                count[i, j] = count[i, j - 1]
29                dir[i, j] = LEFT
30    return count and dir

```

以及重构解的伪代码：

```

1 // 打印解
2 // 参数:
3 //     X: 字符串
4 //     dir: LCS 使用的记录方向的矩阵
5 //     i, j: 想求的是 X[1...i], Y[1...j] 的 LCS
6 // 输出:
7 //     X[1...i], Y[1...j] 的 LCS
8 PrintSolution(X, dir, i, j):
9     // 递归出口
10    if i == 0 || j == 0:
11        return
12    // 情况一
13    if dir[i, j] == LEFTUP:
14        PrintSolution(X, dir, i - 1, j - 1)
15        print(X[i])
16    // 情况二
17    else if dir[i, j] == UP:
18        PrintSolution(X, dir, i - 1, j)
19    else:
20        PrintSolution(X, dir, i, j - 1)

```

### 三、算法复杂性分析

先分析 LCS 的复杂度:

#### 1. 时间复杂度:

设  $|X| = m$ ,  $|Y| = n$ , 算法运行时间为  $T(m, n)$ 。根据算法伪代码, 主要的计算步骤在于双层循环, 共循环  $mn$  次, 每次循环内部是常数时间, 所以  $T(m, n) = O(mn)$ 。

#### 2. 空间复杂度:

算法内部的辅助空间是  $count[0...m, 0...n]$  和  $dir[1...m, 1...n]$ , 共有  $(m+1)(n+1) + mn = O(mn)$  个单元, 所以  $S(m, n) = O(mn)$ 。

再分析 PrintSolution 的时间复杂度:

容易看出, 在打印解的时候, 程序递归走过的路径至多是从  $dir$  右下角走到左上角的曼哈顿距离  $m+n$ , 所以  $T(m, n) = O(m+n)$ 。

### 四、程序源代码

## 最长公共子序列: **lcs.cpp**:

```
1 #include <iostream>
2 #include <vector>
3 #include <string>
4 #include <utility>
5
6 // 最长公共子序列问题
7 // （动态规划法）
8
9 // 方向枚举，重构解的时候使用
10 enum class Direction
11 {
12     // 空方向
13     NULLDIR,
14     // 左上
15     LEFTUP,
16     // 左
17     LEFT,
18     // 上
19     UP
20 };
21
22 // 最长公共子序列
23 // 参数：
24 //     X[]: 字符串1
25 //     Y[]: 字符串2
26 // 返回值：
27 //     count[i][j]: 保存 X[1...i] Y[1...j] 的最长公共子序列长度
28 //     dir[i][j]: 保存 X[1...i] Y[1...j] 的最长公共子序列由哪个前缀确定
29 // 时间复杂度：
30 //     O(mn)
31 // 空间复杂度：
32 //     O(mn)
33 //
34 auto LCS(const std::string & X, const std::string & Y)
35 {
36     // 长度
```

```

37     size_t m = X.length();
38     size_t n = Y.length();
39
40     // 辅助空间: O(mn)
41     // count[0...m][0...n]
42     // 用于记录最长公共子序列长度
43     std::vector<std::vector<size_t>> count;
44     // dir[0...m][0...n] (0 忽略)
45     // 用于重构解
46     std::vector<std::vector<Direction>> dir;
47
48     // 初始化为 (m + 1) * (n + 1) 的全零矩阵
49     // 零向量 1 * (n + 1)
50     std::vector<size_t> zeroSizeVector;
51     std::vector<Direction> zeroDirVector;
52     // 构建零向量
53     for(size_t i = 0; i < n + 1; i++)
54     {
55         zeroSizeVector.push_back(0);
56         zeroDirVector.push_back(Direction::NULLDIR);
57     }
58     // 构建零矩阵
59     for(size_t i = 0; i < m + 1; i++)
60     {
61         count.push_back(zeroSizeVector);
62         dir.push_back(zeroDirVector);
63     }
64     // 填表
65     // 时间复杂度: O(mn)
66     for(size_t i = 1; i <= m; i++)
67     {
68         for(size_t j = 1; j <= n; j++)
69         {
70             // 最后字符相同
71             if(X[i - 1] == Y[j - 1])
72             {
73                 count[i][j] = count[i - 1][j - 1] + 1;

```

```

74         dir[i][j] = Direction::LEFTUP;
75     }
76     // max{ LCS_len(X_{i-1}, Y_j), LCS_len(X_i, Y_{j-1}) }
77     else if(count[i - 1][j] >= count[i][j - 1])
78     {
79         count[i][j] = count[i - 1][j];
80         dir[i][j] = Direction::UP;
81     }
82     else
83     {
84         count[i][j] = count[i][j - 1];
85         dir[i][j] = Direction::LEFT;
86     }
87 }
88 }
89 return std::make_pair(count, dir);
90 }
91
92 // 打印最长公共子序列
93 // 参数:
94 //     X[]: 字符串1
95 //     dir[][]: 表示最长公共子序列由谁决定的矩阵
96 // 输出:
97 //     重构出的最长公共子序列
98 // 时间复杂度:
99 //     O(m+n)
100 void PrintSolution(const std::string & X,
101     const std::vector<std::vector<Direction>> & dir,
102     size_t i, size_t j)
103 {
104     // 打印解的时候, 走过的路径至多就是
105     // 从右下角到左下角的曼哈顿距离 O(m+n)
106     if(i == 0 || j == 0)
107         return;
108     // 左上方, 递归打印前边的最长公共子序列
109     // 并打印该字符
110     if(dir[i][j] == Direction::LEFTUP)

```

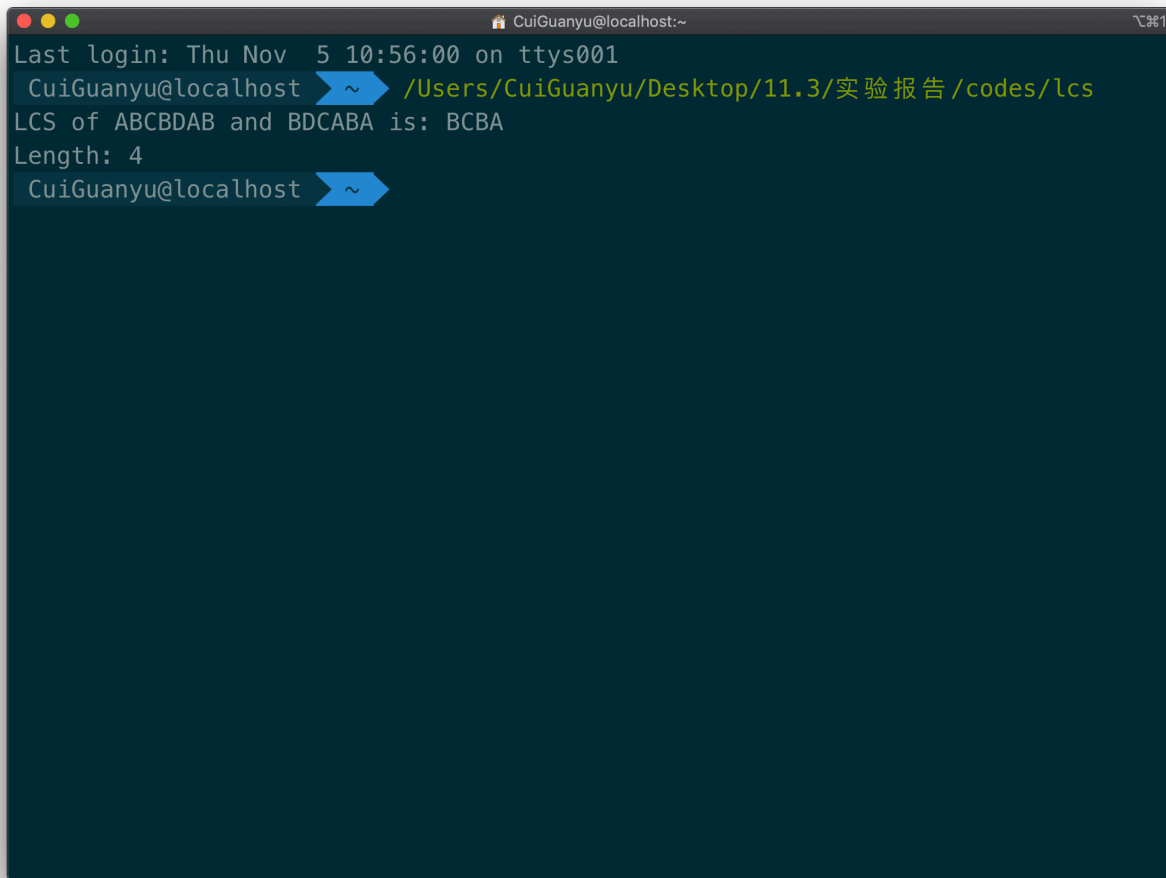
```

111     {
112         PrintSolution(X, dir, i - 1, j - 1);
113         std::cout << X[i - 1];
114     }
115     // 上边或左边, 说明最后字符不同, 需要递归打印
116     else if (dir[i][j] == Direction::UP)
117     {
118         PrintSolution(X, dir, i - 1, j);
119     }
120     else if (dir[i][j] == Direction::LEFT)
121     {
122         PrintSolution(X, dir, i, j - 1);
123     }
124     else
125     {
126         std::cout << "Error";
127     }
128 }
129
130 int main()
131 {
132     // 测试字符串
133     std::string X = "ABCBADAB";
134     std::string Y = "BDCABA";
135     std::cout << "LCS of " << X << " and "
136         << Y << " is: ";
137     // 求解矩阵
138     auto result = LCS(X, Y);
139     // 打印解
140     PrintSolution(X, result.second,
141         X.length(), Y.length());
142     std::cout << std::endl;
143     std::cout << "Length: "
144         << result.first[X.length()][Y.length()] << std::endl;
145     return 0;
146 }

```

## 五、运行结果截图

运行 lcs.cpp，程序以课件上的例子做测试



```
CuiGuanyu@localhost:~  
Last login: Thu Nov  5 10:56:00 on ttys001  
CuiGuanyu@localhost ~ ➤ /Users/CuiGuanyu/Desktop/11.3/实验报告/codes/lcs  
LCS of ABCBDAB and BDCABA is: BCBA  
Length: 4  
CuiGuanyu@localhost ~ ➤
```

图 1: 最长公共子序列测试

可见程序运行正确。



## 上机题 2 实现最优二叉搜索树 (OBST) 算法。

### 一、问题描述

实现最优二叉搜索树算法。

1. 输入：排好序的关键词序列  $K = \langle k_1, k_2, \dots, k_n \rangle$ ，各关键词被查询到的概率数组（可按比例扩大为整数） $P = \langle p_1, p_2, \dots, p_n \rangle$ ，各关键词“之间”<sup>1</sup>被查询的概率数组（可按比例扩大为整数） $Q = \langle q_0, q_1, q_2, \dots, q_n \rangle$ 。
2. 输出：一棵最优二叉搜索树（按先根序打印）。

### 二、算法基本思路

使用动态规划的方法，多阶段决策当前子树的根节点。

首先，定义  $w[i, j] (i = 1 \dots n, j = i - 1 \dots n)^2$  为增加一层时，由  $K[i \dots j]$  构成的子树增加的代价，根据二叉搜索树外节点比内节点多一个的特点，有

$$w[i, j] = Q[i - 1] + \sum_{l=i}^j (P[l] + Q[l])$$

接下来定义  $e[i, j]$  为  $K[i \dots j]$  构成最优二叉搜索树时的最小代价。假若  $K[r]$  为  $K[i \dots j]$  的根时，

$$\begin{aligned} e[i, j] &= P[r] + (e[i, r - 1] + w[i, r - 1]) + (e[r + 1, j] + w[r + 1, j]) \\ &= e[i, r - 1] + e[r + 1, j] + w[i, j] \end{aligned}$$

因为我们要使树代价  $e[i, j]$  最小，所以有

$$e[i, j] = \begin{cases} Q[i - 1], & j = i - 1 \\ \min_{i \leq r \leq j} \{e[i, r - 1] + e[r + 1, j] + w[i, j]\}, & i \leq j \end{cases}$$

为了重构解，我们需要保存  $K[i \dots j]$  的根的位置为  $root[i, j]$ ，即：

$$root[i, j] = \arg \min_{i \leq r \leq j} \{e[i, r - 1] + e[r + 1, j] + w[i, j]\}$$

重构解时，只需要打印出  $root[i, j]$  对应的关键字，并分别打印左子树和右子树即可。

下面给出 OBST 算法的伪代码：

```
1 // 最优二叉搜索树
2 // 参数：
```

---

<sup>1</sup>“伪关键字”，即  $q_i$  为查询词落到  $(w_i, w_{i+1})$  的概率。

<sup>2</sup>当  $j = i - 1$  时，表示子树为空，也即只有外节点。

```

3 //      P[1...n]: 各关键词被查询的概率数组
4 //      Q[0...n]: 各伪关键词被查询的概率数组
5 //      n: 关键词数量
6 // 返回值:
7 //      e[i][j]: 表示关键词i-关键词j的OBST的搜索代价
8 //      root[i][j]: 表示关键词i-关键词j的OBST的根
9 OBST(P, Q, n):
10     e[1...n + 1, 0...n], w[1...n + 1, 0...n], root[1...n, 1...n]
11     // 初始化
12     for i = 1 to n:
13         e[i][i - 1] = Q[i - 1]
14         w[i][i - 1] = Q[i - 1]
15     // 1 长的数组, 类似矩阵链一样斜着填表
16     // n 次
17     for l = 1 to n:
18         // n - l + 1 次
19         for i = 1 to n - l + 1:
20             // 终点
21             j = i + l - 1
22             e[i, j] = +INFINITY
23             w[i, j] = w[i, j - 1] + P[j] + Q[j]
24             // 计算根
25             for r = i to j:
26                 t = e[i, r - 1] + e[r + 1, j] + w[i, j]
27                 if t < e[i, j]:
28                     e[i, j] = t;
29                     root[i, j] = r;
30     return e and root

```

再给出 PrintSolution 的伪代码:

```

1 // 打印解
2 // 参数:
3 //      words: 关键词数组
4 //      root: OBST 算出的根的位置
5 //      i, j: 关键词 i-j 构成的树
6 // 输出:
7 //      最优二叉树, 以先根序给出

```

```

8 PrintSolution(words, root, i, j):
9     // i > j 表示空节点
10    if i > j:
11        print("NULL")
12    // 打印根
13    print(root[i, j])
14    // 左子树
15    PrintSolution(words, root, i, root[i, j] - 1)
16    // 右子树
17    PrintSolution(words, root, root[i, j] + 1, j)

```

### 三、算法复杂性分析

先分析 OBST 的复杂度：

#### 1. 时间复杂度：

类似于矩阵链乘法，OBST 也是在两个循环中填写一个  $O(n) \times O(n)$  的矩阵，其中每个单元确定树根时都要从  $i$  迭代至  $j$ ，共有  $O(n)$  次计算，所以平凡方法的总时间复杂度为  $O(n^3)$ <sup>3</sup>。

#### 2. 空间复杂度：

内部辅助空间  $e$ 、 $w$  与  $root$  都是  $O(n) \times O(n)$  的矩阵，所以空间复杂度  $S(n) = O(n^2)$ 。

再分析 PrintSolution 的时间复杂度：

设打印节点为基本操作。因为二叉搜索树有  $n$  个内节点和  $n+1$  个外节点，所以共会打印  $O(n)$  次，即时间复杂度为  $T(n) = O(n)$ 。

### 四、程序源代码

最优二叉搜索树: **OBST.cpp**:

```

1 #include <iostream>
2 #include <string>
3 #include <vector>
4 #include <utility>
5 #include <limits>
6
7 // 最优二叉搜索树
8 // 参数:
9 //     P[]: 各关键字被检索的概率

```

---

<sup>3</sup>根据课本 15.5-4, Knuth 改进了算法，使得时间复杂度降低到了  $\Theta(n^2)$ 。

```

10 //      Q[]: 落在各关键字“之间”的概率
11 //      n: 关键字个数
12 // 返回值:
13 //      e[i][j]: a_i...a_j 构成的最优二叉树的COST
14 //      root[i][j]: a_i...a_j 构成的最优二叉树的根的位置
15 // 时间复杂度: (认为填表是基本操作)
16 //      T(n)=O(n^3), 若确定根时使用平凡方法
17 //      T(n)=O(n^2), 若确定根时使用 Knuth 改进的方法
18 // 空间复杂度:
19 //      S(n)=O(n^2)
20 //
21 auto OBST(const std::vector<double> & P,
22     const std::vector<double> & Q, size_t n)
23 {
24     // 辅助空间: O(n^2)
25     // 补偿 w[0...n+1, 0...n]
26     std::vector<std::vector<double>> w;
27     // 代价 e[0...n+1, 0...n]
28     std::vector<std::vector<double>> e;
29     // 为了重构解表示 a_i...a_j 的根的位置
30     // root[0...n, 0...n]
31     std::vector<std::vector<size_t>> root;
32
33     // 初始化
34     // 零向量 1 * (n + 1)
35     std::vector<double> zeroDoubleVector;
36     std::vector<size_t> zeroSizeVector;
37     for(size_t i = 0; i < n + 1; i++)
38     {
39         zeroDoubleVector.push_back(0.0);
40         zeroSizeVector.push_back(0);
41     }
42     // w、e、root 赋初值
43     // (n + 2) * (n + 1)
44     for(size_t i = 0; i < n + 2; i++)
45     {
46         w.push_back(zeroDoubleVector);

```

```

47         e.push_back(zeroDoubleVector);
48         root.push_back(zeroSizeVector);
49     }
50
51     // 动态规划填表
52     // 初始化
53     for(size_t i = 1; i <= n + 1; i++)
54     {
55         e[i][i - 1] = Q[i - 1];
56         w[i][i - 1] = Q[i - 1];
57     }
58     // 1 节点的树
59     // 循环 1 次
60     for(size_t l = 1; l <= n; l++)
61     {
62         // 迭代树的起点
63         // 循环 n - l + 1 次
64         for(size_t i = 1; i <= n - l + 1; i++)
65         {
66             // 终点
67             size_t j = i + l - 1;
68             e[i][j] = std::numeric_limits<double>::infinity();
69             w[i][j] = w[i][j - 1] + P[j] + Q[j];
70             // 找根的位置
71             for(size_t r = i; r <= j; r++)
72             {
73                 double t = e[i][r - 1] + e[r + 1][j] + w[i][j];
74                 if(t < e[i][j])
75                 {
76                     e[i][j] = t;
77                     root[i][j] = r;
78                 }
79             }
80         }
81     }
82     return std::make_pair(e, root);
83 }

```

```

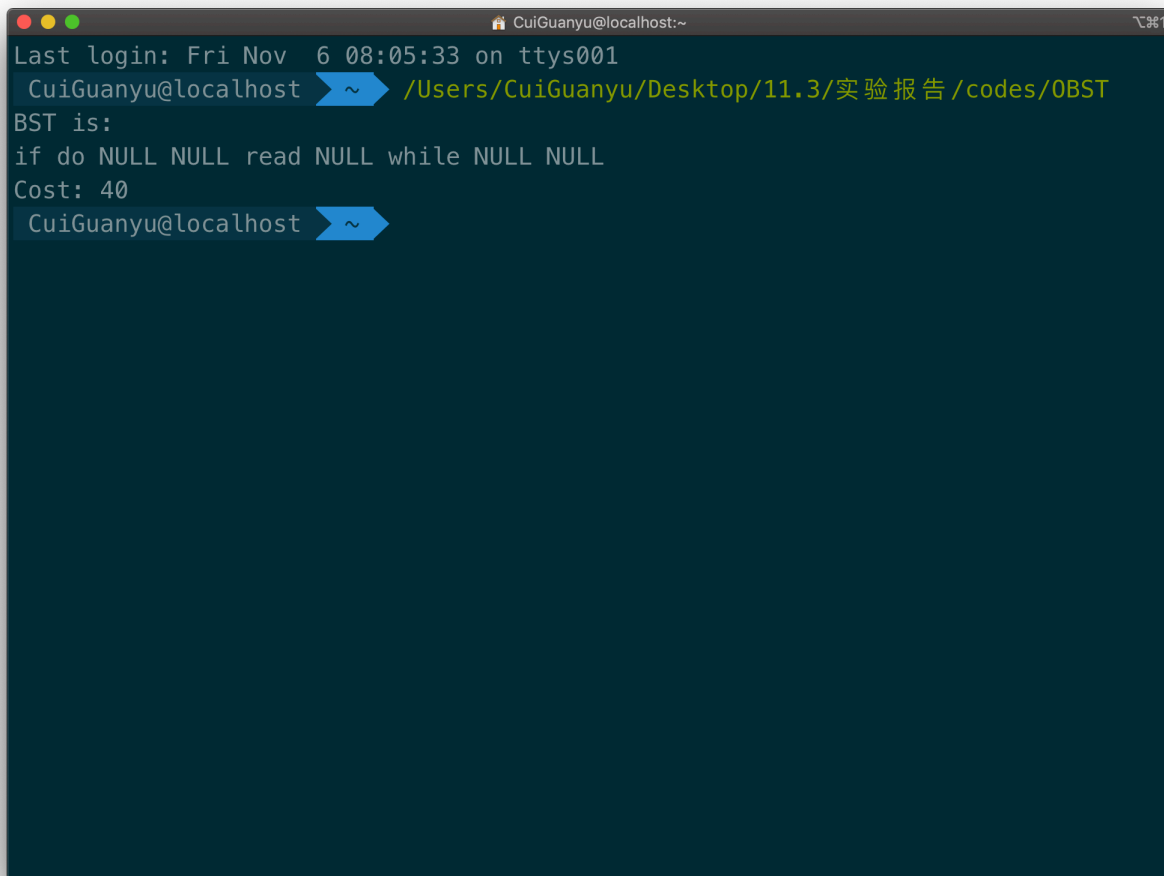
84
85 // 打印解
86 // 参数:
87 //     words[]: 关键字
88 //     root[i][j]: OBST算出的根位置
89 //     i, j: 从 i 开始到 j 结束的树
90 // 输出:
91 //     最优二叉树(按数组的形式先根打印)
92 // 时间复杂度: (认为打印节点是基本操作)
93 //      $T(n) = O(n)$ 
94 void PrintSolution(const std::vector<std::string> & words,
95     const std::vector<std::vector<size_t>> & root,
96     size_t i, size_t j)
97 {
98     // i > j, 表示空节点
99     if(i > j)
100     {
101         std::cout << "NULL ";
102         return;
103     }
104     // 因为会打印 n 个内节点和 n+1 个外节点
105     // 所以打印次数为  $O(n)$ 
106     // 打印根
107     std::cout << words[root[i][j] - 1] << " ";
108     // 打印左子树
109     PrintSolution(words, root, i, root[i][j] - 1);
110     // 打印右子树
111     PrintSolution(words, root, root[i][j] + 1, j);
112 }
113
114 int main()
115 {
116     // 关键字
117     std::vector<std::string> words = {"do", "if", "read", "while"};
118     // P[0] 补0
119     std::vector<double> P = {0, 3, 3, 1, 1};
120     // Q

```

```
121     std::vector<double> Q = {2, 3, 1, 1, 1};
122     auto result = OBST(P, Q, words.size());
123     std::cout << "BST is: " << std::endl;
124     PrintSolution(words, result.second, 1, words.size());
125     std::cout << std::endl;
126     std::cout << "Cost: " << result.first[1][words.size()] << std::endl;
127     return 0;
128 }
```

## 五、运行结果截图

运行 OBST.cpp，程序以课件上的例子做测试（注：老师给的代码与书上的伪代码略有不同。如果按照书上的计算方式，则代价应该为 40；而老师的方法计算结果是 32，没有考虑最底层的外节点。两种方法得到的树是一致的。）



```
CuiGuanyu@localhost:~
Last login: Fri Nov 6 08:05:33 on ttys001
CuiGuanyu@localhost ~ ➜ /Users/CuiGuanyu/Desktop/11.3/实验报告/codes/OBST
BST is:
if do NULL NULL read NULL while NULL NULL
Cost: 40
CuiGuanyu@localhost ~ ➜
```

图 2: 最优二叉搜索树测试

可见程序运行正确。