

522031910747+李若彬+lab0

522031910747+李若彬+lab0

- 实验环境
 - 哈夫曼树的实现
 - 哈夫曼树类
 - 哈夫曼树的建树过程
 - 功能函数
- 实验设计
 - 三组不同的文本分析

实验环境

处理器	12th Gen Intel(R) Core(TM) i7-12700H 2.70 GHz
机带 RAM	16.0 GB (15.8 GB 可用)
设备 ID	14786D6D-2D9C-42A3-8C80-A819975DA58B
产品 ID	00342-30630-82359-AAOEM
系统类型	64 位操作系统, 基于 x64 的处理器
笔和触控	没有可用于此显示器的笔或触控输入

相关链接 域或工作组 系统保护 高级系统设置

Windows 规格

复制 ^

版本	Windows 11 家庭中文版
版本号	24H2
安装日期	2024/2/20
操作系统版本	26058.1300

哈夫曼树的实现

哈夫曼树类

- 私有成员：
 - 首先需要定义节点类

```
1 struct Node
2 {
3     std::string data; // 储存的字符
4     int frequency; // 字符出现频率
5     Node *left; // 左儿子
6     Node *right; // 右儿子
7     Node( std::string data, int frequency, Node * left, Node * right
8 ) : data( data ), frequency( frequency ), left( left ), right(
9 right ) // 节点构造函数
10
11     std::string getMinData() const { // 获取节点最小字典序
12         if ( this->left == nullptr && this->right == nullptr )
13         {
14             return(data); // 如果是叶子节点, 直接返回节点的数据
15         }
16         std::string minData = this->data; // 初始化最小字典序为当前节点的数据
17
18         if ( this->left != nullptr )
```

```

15     {
16         std::string leftMinData = this->left->getMinData();//递归
        获取左子树的最小字典序
17         if ( leftMinData < minData )
18         {
19             minData = leftMinData;//更新最小字典序
20         }
21     }
22     if ( this->right != nullptr )
23     {
24         std::string rightMinData = this->right->getMinData();//递
        归获取右子树的最小字典序
25         if ( rightMinData < minData )
26         {
27             minData = rightMinData;//更新最小字典序
28         }
29     }
30     return(minData);
31 }
32 };

```

- 比较函数（用于后续哈夫曼树构造的优先级队列）

```

1 static bool compare( const Node* a, const Node* b )
2 {
3     if ( a->frequency == b->frequency )
4     {
5         std::string minDataA = a->getMinData();//获取树a中所有节点
        的最小字典序
6         std::string minDataB = b->getMinData();// 获取树b中所有节
        点的最小字典序
7         return(minDataA < minDataB);//比较最小字典序
8     }
9     return(a->frequency < b->frequency);
10 }

```

- 一个节点类的指针

```

1 Node *ptr;

```

- 公有成员（handout 已给出的 skeleton code）

- Option
- 构造函数
- 编码函数

```

1 enum class Option
2 {
3     SingleChar,
4     Multichar
5 };
6 hfTree( const std::string & text, const Option op );
7 std::map<std::string, std::string> getCodingTable();

```

哈夫曼树的建树过程

- 构造函数

```
1  hfTree::hfTree(const std::string &text, const Option op)
2  {
3      // TODO: Your code here
4      std::deque<Node *> forest; // 创建一个双端队列forest用于存储节点
5      if (op == Option::SingleChar) // 单字符压缩
6      {
7          std::map<char, int> charCount; // 用一个频率表map来存储字符&频率
8          for (char c : text) {
9              charCount[c]++;
10         }
11         for (auto it = charCount.begin(); it != charCount.end(); it++)
12         { // 使用迭代器将map对应的键值对存入队列中
13             Node *node = new Node(std::string(1, it->first), it->second, nullptr, nullptr);
14             forest.push_back(node);
15         }
16         int j = forest.size();
17         for (int i = 0; i < j - 1; i++) { // 构造哈夫曼树
18             std::sort(forest.begin(), forest.end(), compare); // 始终保持队列中头两个节点最小
19             ptr = new Node(forest[0]->data + forest[1]->data,
20                             forest[0]->frequency + forest[1]->frequency, forest[0], forest[1]);
21             forest.pop_front();
22             forest.pop_front();
23             forest.push_back(ptr);
24         }
25         ptr = forest.front(); // 返回哈夫曼树的根节点
26     }
27     else if (op == Option::MultiChar) // 多字符压缩
28     {
29         std::map<std::string, int> charCount;
30         for (int i = 0; i < text.size(); i++) {
31             std::string temp = text.substr(i, 1);
32             charCount[temp] = 0;
33         }
34         std::map<std::string, int> tempCount; // 用于存储两个字符的键值对
35         for (int i = 0; i < text.size() - 1; i++) {
36             std::string temp = text.substr(i, 2);
37             tempCount[temp]++;
38         }
39         std::vector<std::pair<std::string, int>> topCombinations; // 将两个字符的键值对转换为vector类型便于后续比较
40         for (auto it = tempCount.begin(); it != tempCount.end(); it++)
41         {
42             topCombinations.push_back(*it);
43         }
44         // 对两个字符的键值对排序，频率较大的在前面，如果频率一样则字典序较小的在前面
45         std::sort(topCombinations.begin(), topCombinations.end(), [](const auto& a, const auto& b) {
46             if (a.second == b.second) {
```

```

44         return a.first < b.first;
45     }
46     return a.second > b.second;
47 });
48     for (int i = 0; i < 3; i++) { //将频率最大的三个（不足三个的按全部）两
个字符的键值对插入到map
49         if (i >= topCombinations.size()) break;
50         std::string temp = topCombinations[i].first;
51         charCount[temp] = 0;
52     }
53     for (int i = 0; i < text.size(); i++) { //用map中的键扫描text内
容，并更新其值
54         std::string currentTwoChars = text.substr(i, 2); //首先先匹配
两个字符的键
55         if (charCount.find(currentTwoChars) != charCount.end()) {
56             charCount[currentTwoChars]++;
57             i++;
58         }
59         else { //两个字符没有匹配成功，匹配单个字符
60             std::string currentOneChar = text.substr(i, 1);
61             if (charCount.find(currentOneChar) != charCount.end())
{
62                 charCount[currentOneChar]++; //匹配到单个字符的字符串，
计数加一
63             }
64         }
65     }
66     auto it = charCount.begin();
67     while (it != charCount.end()) { //删除频率表map中值为0对应的键值对
68         if (it->second == 0) {
69             it = charCount.erase(it);
70         } else {
71             ++it;
72         }
73     }
74     for (auto it = charCount.begin(); it != charCount.end(); it++)
{
75         Node *node = new Node(it->first, it->second, nullptr,
nullptr);
76         forest.push_back(node);
77     }
78     int j = forest.size();
79     for (int i = 0; i < j - 1; i++) {
80         std::sort(forest.begin(), forest.end(), compare);
81         ptr = new Node(forest[0]->data + forest[1]->data,
forest[0]->frequency + forest[1]->frequency, forest[0], forest[1]);
82         forest.pop_front();
83         forest.pop_front();
84         forest.push_back(ptr);
85     }
86     ptr = forest.front();
87 }
88 }

```

- 编码

```

1  std::map<std::string, std::string> hfTree::getCodingTable()
2  {
3      // TODO: Your code here
4      if (ptr != nullptr) { //哈夫曼树根节点不为空
5          std::map<std::string, std::string> codingTable; //用一个map来储存
键值对
6          std::string code; //霍夫曼编码
7          std::queue<std::pair<Node *, std::string>> q; //用一个双端队列来存
节点和编码，其每个元素都是一个pair
8          q.push(std::make_pair(ptr, "")); //压入根节点
9          while (!q.empty()) {
10             Node *node = q.front().first; //获取队列中的第一个节点
11             code = q.front().second; //获取节点当前编码
12             q.pop();
13             if (node->left == nullptr && node->right == nullptr) { //叶子
节点
14                 codingTable[node->data] = code; //更新map的键值对
15             }
16             if (node->left != nullptr) { //有左孩子，将左孩子压入队列，编码加0
17                 q.push(std::make_pair(node->left, code + "0"));
18             }
19             if (node->right != nullptr) { //有右孩子，将右孩子压入队列，编码加
1              1
20                 q.push(std::make_pair(node->right, code + "1"));
21             }
22         }
23         return codingTable;
24     }
25     return std::map<std::string, std::string>();
26 }

```

功能函数

- 读取待压缩编码的文件内容

```

1  std::string parseText( const std::string &input )
2  {
3      std::string content;
4      // TODO: Your code here
5      std::ifstream file( input );
6      if ( file.is_open() )
7      {
8          char c;
9          while ( file.get( c ) )
10             {
11                 content += c;
12             }
13         file.close();
14     }
15     return content;
16 }

```

- 将 data 中的内容输出到 output 指定的文件

```

1 void output( const std::string &output, const std::string &data )
2 {
3     // TODO: Your code here
4     std::string extension = output.substr( output.find_last_of( "." ) +
1 1 );//查找输出文件的扩展名
5     if ( extension == "huffzip" )//输出到zip文件
6     {
7         std::ofstream outputFile( output, std::ios::binary );//以二进制形
            式打开文件
8         if ( outputFile.is_open() )
9         {
10             std::string byteString;
11             for ( int i = 0; i < data.size(); i += 8 )//将编码数据按照8位
                一组转换为字节并写入文件
12             {
13                 std::string byte      = data.substr( i, 8 );
14                 uint8_t    byteValue  = 0;
15                 for ( int j = 0; j < 8; j++ )//将8位编码转换为一个字节的数值
16                 {
17                     byteValue <<= 1;
18                     if ( byte[j] == '1' )
19                         byteValue |= 1;
20                 }
21                 //byteValue 是一个 uint8_t 类型的变量，它占用一个字节的内存空
                间。reinterpret_cast 将其地址转换为 const char* 类型的指针，以便将字节数据写入
                文件。write 函数根据指针和字节数将数据写入文件
22                 outputFile.write( reinterpret_cast<const char *>
                (&byteValue), sizeof(byteValue) );
23             }
24             outputFile.close();
25         }else {
26             std::cerr << "Failed to open output file" << std::endl;
27         }
28     }else { //正常输出字符串
29         std::ofstream file( output );
30         if ( file.is_open() )
31         {
32             file.write( data.c_str(), data.size() );
33             file.close();
34         }
35     }
36 }
37 }

```

- 将编码表转为字符串

```

1  std::string codingTable2String( const std::map<std::string,
std::string> &codingTable )
2  {
3      std::string result = "";
4      // TODO: Your code here
5      for ( auto it = codingTable.begin(); it != codingTable.end(); it++
6      )
7      {
8          result += it->first + " " + it->second + "\n"; //原符号+空格+对应编
码+换行符
9      }
10     return(result);
11 }

```

- 读取编码表

```

1  void loadCodingTable( const std::string &input, std::map<std::string,
std::string> &codingTable )
2  {
3      // TODO: Your code here
4      std::ifstream file( input );
5      if ( file.is_open() )
6      {
7          char buffer[2]; //一次读取两个字符
8          while ( file.read( buffer, sizeof(buffer) ) )
9          {
10             if ( file.peek() == '0' || file.peek() == '1' ) //下一个字符为
数字，说明为单个字符，多取了一个字符
11             {
12                 std::string key( buffer, 1 ); //只取buffer的第一个字符作为
key
13                 std::string value;
14                 std::string plus( buffer, 2 );
15                 value += plus;
16                 std::getline( file, value ); //读取0或1后的所有字符作为value
17                 codingTable[key] = value; //更新键值对
18             } else if ( file.peek() == ' ' ) //下一个字符为空格，说明正好取两
个字符
19             {
20                 file.ignore( 1 ); //忽略空格
21                 std::string key( buffer, sizeof(buffer) ); //这两个所取得
字符即为键
22                 std::string value;
23                 std::getline( file, value ); //读取空格后的所有字符作为value
24                 codingTable[key] = value; //更新键值对
25             }
26         }
27         file.close();
28     }
29 }

```

- 压缩

```

1  std::string compress( const std::map<std::string, std::string>
&codingTable, const std::string &text )

```

```

2 {
3     std::string result;
4     // TODO: Your code here
5     std::string bitstream; //用于存储有效比特流
6     std::string currentCode;
7     for ( size_t i = 0; i < text.size(); ++i )
8     {
9
10        //优先比较两个字符
11        if ( i + 1 < text.size() )
12        {
13            std::string doubleChar = text.substr( i, 2 );
14            if ( codingTable.count( doubleChar ) > 0 )
15            {
16                bitstream += codingTable.at( doubleChar );
17                i += 1;
18                continue;
19            }
20        }
21        //如果两个字符的匹配失败，那么就匹配一个字符
22        std::string singleChar = text.substr( i, 1 );
23        if ( codingTable.count( singleChar ) > 0 )
24        {
25            bitstream += codingTable.at( singleChar );
26        }
27    }
28    uint64_t effectiveBits = bitstream.size(); //存储比特流的有效位
29    for ( int i = 0; i < 8; ++i )
30    {
31        result += static_cast<char>( (effectiveBits >> (i * 8) ) & 0xFF
32    );
33    }
34    while ( bitstream.size() % 8 != 0 ) //将比特流补0使得成为字节流
35    {
36        bitstream += '0';
37    }
38    std::string ans = "";
39    std::stack<int> binaryEffectiveBits; //将比特流的有效位数转换为二进制
40    while ( effectiveBits > 0 )
41    {
42        binaryEffectiveBits.push( effectiveBits % 2 );
43        effectiveBits /= 2;
44    }
45    std::string binaryStr = "";
46    while ( !binaryEffectiveBits.empty() )
47    {
48        binaryStr += std::to_string( binaryEffectiveBits.top() );
49        binaryEffectiveBits.pop();
50    }
51    while ( binaryStr.size() != 64 ) //将二进制的有效位数补0使得其有64位
52    {
53        binaryStr = "0" + binaryStr;
54    }
55    std::string reverse = "";
56    int jj = binaryStr.size() / 8;
57    for ( int i = 0; i < jj; i++ ) //以小端的方式表示64位的二进制下的有效位数

```



```

57     {
58         reverse += binaryStr.substr( binaryStr.length() - 8 );
59         binaryStr.erase( binaryStr.length() - 8 );
60     }
61     ans = reverse + bitstream;
62     result = ans;
63     return result;
64 }

```

实验设计

- 三组不同的文本（长度有较为显著的差异，但三个文件的大小之和不超过 1 MB）
- 比较两种压缩方式在这些文本上压缩效果的差异

三组不同的文本

1.

文本大小:

大小: 8.83 KB (9,045 字节)

sin.huffzip 大小:

大小: 5.03 KB (5,152 字节)

mul.huffzip 大小:

大小: 5.02 KB (5,144 字节)

sin 压缩率: $5152 \div 9045 = 56.9596\%$

mul 压缩率: $5144 \div 9045 = 56.8712\%$

mul 相较于 sin 方式压缩提升率: $(5152 - 5144) \div 5152 = 0.1553\%$

2.

文本大小:

大小: 328 KB (336,394 字节)

sin.huffzip 大小:

大小: 204 KB (209,038 字节)

mul.huffzip 大小:

大小: 201 KB (206,556 字节)

sin 压缩率: $209038 \div 336394 = 62.1408\%$

mul 压缩率: $206556 \div 336394 = 64.4030\%$

mul 相较于 sin 方式压缩提升率: $(209038 - 206556) \div 209038 = 1.1873\%$

3.

文本大小:

大小: 803 KB (822,364 字节)

sin.huffzip 大小:

大小: 599 KB (614,324 字节)

mul.huffzip 大小:

大小: 559 KB (572,807 字节)

sin 压缩率: $614324 \div 822364 = 74.7022\%$

mul 压缩率: $572807 \div 822364 = 69.6537\%$

mul 相较于 sin 方式压缩提升率: $(614324 - 572807) \div 614324 = 6.7582\%$

分析

- sin 和 mul 两种压缩方式都在一定程度上减少了原文件的大小，压缩率大约在50%以上，且原文件越大，压缩率也相应增加，压缩效果对应降低
- mul 比 sin 方式在各种原文件大小的情况下压缩率都更小，压缩效果都更好，且随着原文件大小的增加，mul 比 sin 在压缩方面的提升越来越大
- 对于更优的压缩策略：当原始文件大小较大时，应当更多地使用连续两个字符参与编码，例如设置参与的数量与原始文本大小正相关，因为在大量的文本中会有大量的重复字符串，这样可以更高效率地压缩并且能达到更低地压缩率，同时还可以适当地引入一些更长的连续字符参与编码，这样会大大降低编码时的资源使用，并压缩后的文件大小也会显著降低。