

数据结构与算法 - 综合训练 1

用树解决问题

吴天阳 2204210460 强基数学 002

1 实验目的

实验 1: 实现二叉检索树 (Binary Search Tree, BST), 需要包含以下 8 种结构

```
1 class BSTBase(): # BST 基类
2     def insert(self, key, value): pass # 插入 (key,value) 键值对
3     def remove(self, key): pass # 删除 key 节点, 若找到并删除节点, 则返回对应
    ↪ 的 value, 若无该节点则返回 None
4     def search(self, key): pass # 查询 key 对应的 value, 若无 key 节点, 则返
    ↪ 回 None
5     def update(self, key, value): pass # 更新 key 对应的 value, 若无 key 节
    ↪ 点, 返回 False, 否则返回 True
6     def isEmpty(self): pass # 判断二叉搜索树是否为空
7     def clear(self): pass # 重新初始化 BST
8     def showStructure(self, file): pass # 输出当前二叉树的节点总数和高度到文
    ↪ 件 file 中
9     def printInorder(self, file): pass # 输出二叉树的中序遍历到文件 file 中
```

读入命令文件 BST_testcases.txt 并处理执行, 将输出运行结果到文件中, 并与 BST_result.txt 进行比较, 判断程序是否正确.

实验 2: 使用 BST 为文稿建立单词索引表, 基于给出的英文文本 article.txt, 记录每个单词在文本中出现的行号 (编号从 1 开始).

2 实验原理

二叉搜索树是一种二叉树形式的数据结构, 支持插入、查找、删除键值对的功能, 对于一个有 n 个节点的二叉搜索树, 每次操作的最优复杂度为 $\mathcal{O}(\log n)$, 最坏为 $\mathcal{O}(n)$, 随机构成一颗二叉树的期望高度为 $\mathcal{O}(\log n)$.

注: 平衡树是在二叉搜索树基础上进行的改进, 可以保证每次查找的复杂度为 $\mathcal{O}(\log n)$.

二叉搜索树定义如下:

1. 空树是二叉搜索树.
2. 若二叉搜索树的左子树非空, 则左子树上所有点的 key 值均小于根节点的 key 值.
3. 若二叉搜索树的右子树非空, 则右子树上所有点的 key 值均大于根节点的 key 值.
4. 二叉搜索树的左右子树均为二叉搜索树.

实现以上操作基本均通过递归即可实现, 节点之间的关联可用指针实现, 具体实现请见下文.

3 实验步骤与结果分析

使用 Python 3.9.12 进行实现. 在 Python 中数值类型数据 `int`, `float`, `str` 无法直接传入实参, 但是使用 `class`, `list`, `dict` 则默认传入实参, 所以需要利用该性质实现指针操作.

代码文件总共有三个: `my_bst.py` 包含 BST 核心类, `main1.py`, `main2.py` 分别为实验 1 和实验 2 的代码, 对应的输出文件分别为 `my_result1.txt`, `my_result2.txt`.

3.1 实现 BST 类

3.1.1 初始化

首先创建节点类 `Node`, 具体有以下属性:

- `key(int)`: 存储 key 值.
- `val0`: 用于初始化 val 值. (用于实验 2, 可以为 list, 存储多个 value 值)
- `val`: 存储 val 值.
- `child(list)`: 长度为 2 的 list, `child[0]`, `child[1]` 分别表示左右子节点.

```

1 class Node(): # 节点子类
2     val0 = 0
3     def __init__(self):
4         self.key = None # 初始化键值
5         self.child = [None, None] # 初始化左右孩子节点
6         if isinstance(self.val0, list):
7             self.val = [] # 由于 list 按照实参赋值, 必须重新创建空 list
8         else: self.val = self.val0

```

初始化 BST 类

```

1 def __init__(self, val0=0):
2     self.Node.val0 = val0 # val0 为每个节点值的初值
3     self.file = None # 将要写入的文件
4     self.root = None # 创建根节点
5     self.height = 0 # 树的高度
6     self.size = 0 # 树的节点总数

```

判断二叉搜索树是否为空, 只需判断 `BST.root` 是否为 `None` 即可.

```

1 def isEmpty(self):
2     return self.root == None

```

重新初始化整个 BST, 由于需要删除掉全部节点, 所以需要递归删除, 这里使用中序遍历中的递归顺便完成该任务:

```

1 def dfs(self, p, delete_node=False): # 输出中序遍历结果
2     if p.child[0]:
3         self.dfs(p.child[0], delete_node)
4     if not delete_node:
5         self.file.write('[{ } ---- < { } >]\n'.format(p.key,
        ↪ str(p.val)[1:-1])) # 在文件中直接写入

```

```

6     if p.child[1]:
7         self.dfs(p.child[1], delete_node)
8     if delete_node: # 用于清空整棵树
9         del p
10
11 def clear(self):
12     self.dfs(self.root, delete_node=True)
13     self.__init__() # 调用初始化函数, 清空 BST
14
15 def printInorder(self, file):
16     if file is None:
17         return None
18     self.file = file
19     self.dfs(self.root)
20     return # 返回中序遍历

```

3.1.2 加入节点

直接通过递归找到对应键值位置, 然后修改 value 值即可.

```

1 def add(self, p, key, val):
2     if p is None: # 当前节点为空节点, 开始创建
3         p = self.Node()
4         p.key = key
5     if key == p.key: # 找到 key 值对应节点, 更新节点 val
6         if isinstance(p.val, list): # 如果 val 当前是列表, 则加入值
7             p.val.append(val)
8         else: # 否则直接修改当前值
9             p.val = val
10    elif key < p.key: # key 节点在左子树中
11        p.child[0] = self.add(p.child[0], key, val)
12    else: # key 节点在右子树中
13        p.child[1] = self.add(p.child[1], key, val)
14    return p
15
16 def insert(self, key, val):
17     if key is None or val is None:
18         return
19     self.root = self.add(self.root, key, val) # 每次从根节点开始查找插入位置

```

3.1.3 查询节点

直接通过递归搜索即可.

```

1 def find(self, p, key):
2     if p == None: # 空节点
3         return None
4     elif key == p.key: # 找到了 key 节点
5         return p.val
6     elif key < p.key: # key 节点在左子树中

```

```

7         return self.find(p.child[0], key)
8     return self.find(p.child[1], key) # 否则只能在右子树中
9
10 def search(self, key):
11     if key is None:
12         return None
13     return self.find(self.root, key)

```

3.1.4 删除节点

删除节点稍微有点麻烦，若删除的节点有两个子节点，为了保持二叉搜索树的性质，需要用左子树中的最大值，或者右子树的最小值替代删除掉的节点；若删除的节点仅有一个儿子，则直接用它代替即可；若删除的节点为叶子结点，直接删去即可。

```

1 def delete_min(self, p): # 查找最小值
2     if p.child[0] is None: # 找到最小值
3         return p, p.child[1] # 返回的第一个参数为找到的最小值点，第二个参数为
        ↪ 更新点编号
4     min_p, update_p = self.delete_min(p.child[0])
5     p.child[0] = update_p # 用第二个参数更新该点
6     return min_p, p # 返回的第一个参数为找到的最小值点，第二个参数为更新点编号
7
8 def delete(self, p, key):
9     if p.key == key:
10         if p.child[0] and p.child[1]: # 如果有两个儿子节点，就需要找到左子树
            ↪ 中的最大值或右子树的最小值替代
            # 这里找右子树最小值替代
11             min_p, update_p = self.delete_min(p.child[1])
12             min_p.child[0] = p.child[0]
13             min_p.child[1] = update_p
14             del p # 删除掉当前节点 p
15             return min_p
16         else:
17             ret = p.child[0] if p.child[0] else p.child[1]
18             del p # 删除掉当前节点 p
19             return ret
20         elif key < p.key: # key 节点在左子树中
21             p.child[0] = self.delete(p.child[0], key)
22         else: # key 节点在右子树中
23             p.child[1] = self.delete(p.child[1], key)
24     return p
25
26
27 def remove(self, key):
28     val = self.search(key)
29     if key is None or val is None: # 如果找不到该 key 值也返回 False
30         return None
31     self.root = self.delete(self.root, key)
32     return val

```

3.1.5 查询 BST 结构

直接通过递归即可完成。

```

1 def struct(self, p, h): # 查找 bst 的节点数和高度
2     self.size += 1 # 总节点数 +1
3     self.height = max(self.height, h) # 更新树的深度
4     if p.child[0]:
5         self.struct(p.child[0], h + 1)
6     if p.child[1]:
7         self.struct(p.child[1], h + 1)
8
9 def showStructure(self, file): # 返回树的总节点数, 返回树的高度
10    if file is None:
11        return None
12    self.size, self.height = 0, 0 # 先初始化为 0
13    self.struct(self.root, 1)
14    file.write('-----\n')
15    file.write(f'There are {self.size} nodes in this BST.\nThe height of
16    ↪ this BST is {self.height}.\n')
17    file.write('-----\n')

```

3.2 实验 1

主要是对读入字符串的划分和条件判断, 执行对应的 BST 函数即可。

```

1 import my_bst
2
3 with open('BST_testcases.txt', 'r', encoding='utf-8') as file,\
4     open('my_result1.txt', 'w', encoding='utf-8') as outfile:
5     bst = my_bst.BST()
6     while True:
7         line = file.readline()
8         if not line:
9             break
10        opt = line[0]
11        if opt == '#':
12            bst.showStructure(outfile)
13            continue
14        key = line.split(' ')[1] # 提取 key 值
15        if opt == '+':
16            value = line.split('\n')[1] # 提取 value 值
17            bst.insert(key, value)
18        elif opt == '-':
19            value = bst.remove(key)
20            if value is not None:
21                outfile.write(f'remove success ---{key} {value}\n')
22            else:
23                outfile.write(f'remove unsuccess ---{key}\n')
24        elif opt == '?':
25            value = bst.search(key)

```

```

26         if value is not None:
27             outfile.write(f'search success ---{key} {value}\n')
28         else:
29             outfile.write(f'search success ---{key}\n')
30     elif opt == '=':
31         value = line.split('\n')[1] # 提取 value 值
32         flag = bst.update(key, value)
33         if flag:
34             outfile.write(f'update success ---{key} {value}\n')
35         else:
36             outfile.write(f'update success ---{key}\n')

```

3.3 实验 2

主要先对文本进行清洗，清洗方法非常简单，按空格划分出单词，然后去除掉首尾的非字母符号即可。

```

1  import my_bst
2
3  def wash(word): # 将其他字符都去掉，只剩下拉丁字母
4      while len(word):
5          if not word[-1].isalpha():
6              word = word[:-1]
7          elif not word[0].isalpha():
8              word = word[1:]
9          else:
10             break
11     return word
12
13 with open('article.txt', 'r', encoding='utf-8', errors='ignore') as file:
14     bst = my_bst.BST(val0=[])
15     cnt = 0 # 用于记录行数
16     while True:
17         # print(cnt)
18         line = file.readline()
19         if not line:
20             break
21         cnt += 1
22         line = wash(line) # 将每一行也洗一下，把多余的空格去掉
23         words = line.split(' ') # 将一行中的单词分离出来
24         for key in words: # 逐个遍历
25             key = wash(key) # 将 key 中的其他标点去掉，只剩下单词
26             if not len(key): # 洗没了
27                 continue
28             bst.insert(key, cnt)
29
30 with open("my_result2.txt", 'w') as file:
31     bst.printInorder(file)

```

4 结论与讨论

最后通过比对, `my_reslut1.txt` 与 `BST_result.txt` 完全一致, 说明 BST 算法应该基本正确, 实验 2 的输出文件为 `my_result2.txt`, 去重后总计有 16580 个单词.

本次实验, 我学会了如何使用 Python 实现 BST 的指针操作, 相比 C 语言虽然不直观, 但是实现效果非常简洁. Python 在处理文本上相较于 C 也更为方便.