

【哈希集合-英文输入法

题目描述与示例

题目

主管期望你来实现英文输入法单词联想功能，需求如下：

1. 依据用户输入的单词前缀，从已输入的英文语句中联想出用户想输入的单词。
2. 按字典序输出联想到的单词序列，如果联想不到，请输出用户输入的单词前缀。

注意：

1. 英文单词联想时区分大小写
2. 缩略形式如 "don't" 判定为两个单词 "don" 和 "t"
3. 输出的单词序列不能有重复单词，且只能是英文单词，不能有标点符号

输入

输入两行。

首行输入一段由英文单词 `word` 和 标点 构成的语句 `str`，接下来一行为一个英文单词前缀 `pre`。

`0 < word.length() <= 20`，`0 < str.length() <= 10000`，`0 < pre.length() <= 20`

输出

输出符合要求的单词序列或单词前缀。存在多个时，单词之间以单个空格分割

示例一

输入

```
1 I love you
2 He
```

输出

```
1 He
```

说明

用户已输入单词语句 "I love you"，可以提炼出 "I"，"love"，"you" 三个单词。接下来用户输入 "He"，

从已经输入信息中无法联想到符合要求的单词，所以输出用户输入的单词前缀。

示例二

输入

```
1 The furthest distance in the world,Is not between life and death,But when I
  stand in front of you,Yet you don't know that I love you.
2 f
```

输出

```
1 front furthest
```

解题思路

首先我们需要处理输入，将输入的字符串 `s` 根据标点符号和空格隔开，得到一个由若干单词 `word` 组成的单词列表 `lst`。这里稍微有点麻烦，不能再用我们熟悉的 `split()` 方法完成，而是改为较为麻烦的遍历写法。

首先我们初始化 `lst = [""]`，即单词列表中存放了一个空字符串。然后我们遍历字符串 `s` 中的字符 `ch`，当

- `ch` 是字母，则将其加入到 `lst` 最后一个元素的末尾，即延长当前单词。如果此时 `lst[-1]` 为一个空字符串 `""`，则 `ch` 充当了某个单词首位字母的角色。
- `ch` 不是字母，说明遇到一个标点符号，当前单词的获取已经结束，`lst` 的末尾插入一个新的空字符串 `""`。

上述思路整理为代码后即为：

```
1 lst = [""]
```

```
2
3 for ch in s:
4     if ch.isalpha():
5         lst[-1] += ch
6     else:
7         lst.append("")
```

当然这个过程也可用**正则表达式**以更加简短的代码来完成，但这部分知识已经超纲，大部分题目完全用不上，学有余力的同学可以自行研究一下。

得到 `lst` 之后，剩下的工作就相当简单了。由于 `lst` 中可能出现重复单词，我们**使用哈希集合进行去重操作**。又因为最后的输出要求按照字典序排序，因此去重之后再对哈希集合进行调用 `sorted()` 内置函数，再转化为列表。

```
1 lst_sorted = list(sorted(set(lst)))
```

对于 `lst_sorted` 中的每一个单词 `word`，我们可以使用切片来获得其前 `pre_length` 个字符所构成的字符串，并与 `pre` 进行比较，就能够得知 `word` 是否包含前缀 `pre` 了。

```
1 pre_length = len(pre)
2 for word in lst_sorted:
3     if word[:pre_length] == pre:
4         ans.append(word)
```

总体来说本题难度不大，甚至很难归类为哪一种具体的算法。

难点其实主要在于对输入的字符串处理，初始化 `lst = []` 实际上是一个颇有技巧的做法。

当然本题还存在着前缀树的最优解法，但也严重超纲，不要求掌握。

代码

解法一

Python

```

1 # 题目：2024E-英文输入法
2 # 分值：100
3 # 作者：许老师-闭着眼睛学数理化
4 # 算法：哈希集合
5 # 代码看不懂的地方，请直接在群上提问
6
7 s = input()
8 pre = input()
9
10 # 初始化列表lst用于存放所有单词
11 lst = [""]
12
13 # 遍历s中的所有字符ch，如果
14 # 1. ch是字母，则加入到lst最后一个元素的末尾，即延长当前单词
15 # 2. ch不是字母，说明遇到一个标点符号，结束当前单词的获取，lst的末尾插入一个新的空字符串""
16 # 这个过程也可以使用正则表达式来完成，不要求掌握，学有余力的同学可以自学一下
17 for ch in s:
18     if ch.isalpha():
19         lst[-1] += ch
20     else:
21         lst.append("")
22
23 # 用哈希集合去重lst中可能出现的重复单词
24 # 去重后进行排序，排序后在转化为列表lst_sorted
25 lst_sorted = list(sorted(set(lst)))
26
27 # 初始化答案数组
28 ans = list()
29
30 # 获得pre的长度，用于切片
31 pre_length = len(pre)
32 # 遍历lst_sorted中的每一个单词
33 for word in lst_sorted:
34     # 如果word前pre_length个字符的切片等于pre
35     # 说明word的前缀是pre，将其加入答案数组ans中
36     if word[:pre_length] == pre:
37         ans.append(word)
38
39 # 如果ans长度大于0，说明至少存在一个单词的前缀是pre，输出由所有单词组成的字符串
40 # 如果ans长度等于0，说明不存在任何一个单词的前缀是pre，返回pre
41 print(" ".join(ans) if len(ans) > 0 else pre)

```

```
1 import java.util.HashSet;
2 import java.util.ArrayList;
3 import java.util.Collections;
4 import java.util.Scanner;
5
6 public class Main {
7     public static void main(String[] args) {
8         Scanner scanner = new Scanner(System.in);
9
10        // 从输入获取字符串s和前缀pre
11        String s = scanner.nextLine();
12        String pre = scanner.nextLine();
13
14        // 初始化列表lst用于存放所有单词
15        ArrayList<String> lst = new ArrayList<>();
16        lst.add("");
17
18        // 遍历s中的所有字符ch
19        for (char ch : s.toCharArray()) {
20            // 如果ch是字母，则加入到lst最后一个元素的末尾，即延长当前单词
21            if (Character.isLetter(ch)) {
22                int lastIndex = lst.size() - 1;
23                lst.set(lastIndex, lst.get(lastIndex) + ch);
24            } else {
25                // 如果ch不是字母，说明遇到一个标点符号，结束当前单词的获取，lst的末尾
                // 插入一个新的空字符串""
26                lst.add("");
27            }
28        }
29
30        // 用哈希集合去重lst中可能出现的重复单词
31        HashSet<String> set = new HashSet<>(lst);
32        // 去重后进行排序，排序后在转化为列表lstSorted
33        ArrayList<String> lstSorted = new ArrayList<>(set);
34        Collections.sort(lstSorted);
35
36        // 初始化答案数组
37        ArrayList<String> ans = new ArrayList<>();
38        // 获得pre的长度，用于切片
39        int preLength = pre.length();
40
41        // 遍历lstSorted中的每一个单词
42        for (String word : lstSorted) {
43            // 如果word前preLength个字符的切片等于pre
44            // 说明word的前缀是pre，将其加入答案数组ans中
45            if (word.length() >= preLength && word.substring(0,
preLength).equals(pre)) {
```

```

46         ans.add(word);
47     }
48 }
49
50 // 如果ans长度大于0, 说明至少存在一个单词的前缀是pre, 输出由所有单词组成的字符串
51 // 如果ans长度等于0, 说明不存在任何一个单词的前缀是pre, 返回pre
52 if (ans.size() > 0) {
53     System.out.println(String.join(" ", ans));
54 } else {
55     System.out.println(pre);
56 }
57 }
58 }
59

```

C++

```

1  #include <iostream>
2  #include <unordered_set>
3  #include <vector>
4  #include <algorithm>
5  using namespace std;
6
7  int main() {
8      string s;
9      getline(cin, s);
10
11     string pre;
12     getline(cin, pre);
13
14     // 初始化列表lst用于存放所有单词
15     vector<string> lst;
16     lst.push_back("");
17
18     // 遍历s中的所有字符ch
19     for (char ch : s) {
20         // 如果ch是字母, 则加入到lst最后一个元素的末尾, 即延长当前单词
21         if (isalpha(ch)) {
22             int lastIndex = lst.size() - 1;
23             lst[lastIndex] += ch;
24         } else {
25             // 如果ch不是字母, 说明遇到一个标点符号, 结束当前单词的获取, lst的末尾插入一个
             // 新的空字符串""
26             lst.push_back("");
27         }
28     }
29 }

```

```

28     }
29
30     // 用哈希集合去重lst中可能出现的重复单词
31     unordered_set<string> set(lst.begin(), lst.end());
32     // 去重后进行排序, 排序后在转化为列表lstSorted
33     vector<string> lstSorted(set.begin(), set.end());
34     sort(lstSorted.begin(), lstSorted.end());
35
36     // 初始化答案数组
37     vector<string> ans;
38     // 获得pre的长度, 用于切片
39     int preLength = pre.length();
40
41     // 遍历lstSorted中的每一个单词
42     for (string word : lstSorted) {
43         // 如果word前preLength个字符的切片等于pre
44         // 说明word的前缀是pre, 将其加入答案数组ans中
45         if (word.length() >= preLength && word.substr(0, preLength) == pre) {
46             ans.push_back(word);
47         }
48     }
49
50     // 如果ans长度大于0, 说明至少存在一个单词的前缀是pre, 输出由所有单词组成的字符串
51     // 如果ans长度等于0, 说明不存在任何一个单词的前缀是pre, 返回pre
52     if (!ans.empty()) {
53         for (int i = 0; i < ans.size(); i++) {
54             cout << ans[i];
55             if (i != ans.size() - 1) {
56                 cout << " ";
57             }
58         }
59         cout << endl;
60     } else {
61         cout << pre << endl;
62     }
63
64     return 0;
65 }
66

```

时空复杂度

时间复杂度: $O(N \log N + NK)$ 。排序需要的时间复杂度为 $O(N \log N)$ 。遍历 `lst_sorted` 需要 $O(N)$ 的复杂度, 每次对 `word` 进行切片操作需要 $O(K)$ 的复杂度, 故遍历过程共需要 $O(NK)$ 的时

间复杂度。总的时间复杂度为两者相加，即 $O(N\log N + NK)$ ，如果 N 远大于 K ，也会退化成 $O(N\log N)$ 。

空间复杂度： $O(NM)$ 。主要为 `lst_sorted` 的所占空间。

N 为单词数目， M 为单词平均长度， K 为前缀单词 `pre` 的长度。

解法二*

（前缀树解法，不要求掌握，感兴趣的同学可以研究一下）

Python

```
1 # 题目：2024E-英文输入法
2 # 分值：100
3 # 作者：许老师-闭着眼睛学数理化
4 # 算法：前缀树
5 # 代码看不懂的地方，请直接在群上提问
6
7 # 构建前缀树节点类
8 class Trie():
9     def __init__(self) -> None:
10         self.children = [None] * 52      # 大小写均存在，需要构建长度为52的children
        列表
11         self.isEnd = False              # 结束标识符，True表示当前节点是一个单词的结
        尾
12
13     # 将单词word加入前缀树的函数
14     def addword(self, word):
15         node = self
16         # 遍历该单词中的所有字符
17         for ch in word:
18             # 获得ch在children列表中对应的索引
19             ch_idx = self.getIdx(ch)
20             # 如果对应位置为None
21             if node.children[ch_idx] is None:
22                 # 则为这个ch字符创建一个新的前缀树节点
23                 node.children[ch_idx] = Trie()
24             # 令前缀树节点前进到ch所在的节点
25             node = node.children[ch_idx]
26         # 完成该单词的添加，设置最后一个字符的节点的结束标识符为True
27         node.isEnd = True
28
29     # 根据字符ch获得在children列表中的对应索引的函数
```



```

30     def getIdx(self, ch):
31         # 如果ch是小写, 得到26-51的索引
32         if ch.islower():
33             ch_idx = ord(ch) - ord("a") + 26
34         # 如果ch是大写, 得到0-25的索引
35         else:
36             ch_idx = ord(ch) - ord("A")
37         return ch_idx
38
39
40     # 根据在children列表中的索引idx获得对应字符ch的函数
41     def getCh(self, idx):
42         # 如果idx大于等于26, 是一个小写字母
43         if idx >= 26:
44             ch = chr(idx + ord("a") - 26)
45         # 如果idx小于26, 是一个大写字母
46         else:
47             ch = chr(idx + ord("A"))
48         return ch
49
50
51     # 获得前缀prefix最后一个字符所在的节点
52     def getLastNode(self, prefix):
53         node = self
54         for ch in prefix:
55             ch_idx = self.getIdx(ch)
56             if node.children[ch_idx] is None:
57                 return None
58             node = node.children[ch_idx]
59         return node
60
61
62     # 对前缀树进行dfs前序遍历, 搜索得到所有后缀
63     def dfs(self, pre, ans, path):
64         node = self
65         # 遇到一个单词结束标识符, 将当前path合并为字符串后加入ans
66         if node.isEnd:
67             # 要注意path此时仅仅是后缀, 要得到完整的单词字符串还要在前面加上pre
68             ans.append(pre + "".join(path))
69         # 如果node.children存在任意一个非None节点, 需要对非空节点继续进行DFS搜索
70         if any(node.children):
71             # 遍历node.children中的所有下一个节点nxt_node
72             for nxt_idx, nxt_node in enumerate(node.children):
73                 # 如果nxt_node不为空, 则继续递归地进行DFS搜索
74                 if nxt_node is not None:
75                     # 根据nxt_idx获得对应的字符nxt_ch
76                     nxt_ch = self.getCh(nxt_idx)

```

```

77             # 将字符nxt_ch加在path末尾的结果，作为参数传入nxt_node的dfs递归
78             nxt_node.dfs(pre, ans, path + [nxt_ch])
79
80 s = input()
81 pre = input()
82
83 # 初始化列表lst用于存放所有单词
84 lst = [""]
85
86 # 遍历s中的所有字符ch，如果
87 # 1. ch是字母，则加入到lst最后一个元素的末尾，即延长当前单词
88 # 2. ch不是字母，说明遇到一个标点符号，结束当前单词的获取，lst的末尾插入一个新的空字符串""
89 # 这个过程也可以使用正则表达式来完成，不要求掌握，学有余力的同学可以自学一下
90 for ch in s:
91     if ch.isalpha():
92         lst[-1] += ch
93     else:
94         lst.append("")
95
96 # 对lst进行去重，因为使用前缀树，所以无需排序
97 lst = list(set(lst))
98
99 # 初始化前缀树根节点
100 root = Trie()
101
102 # 遍历lst中的每一个单词word，构建前缀树
103 for word in lst:
104     root.addword(word)
105
106 # 调用前缀树中的getLastNode()方法，得到前缀pre在树中的最后一个节点
107 lastNode = root.getLastNode(pre)
108
109 # 如果lastNode为空，说明在root前缀树中，不存在任何前缀为pre的单词，输出pre
110 if lastNode is None:
111     print(pre)
112 # 如果lastNode非空，说明在root前缀树中，存在前缀为pre的单词，要找到所有单词
113 else:
114     # 初始化答案数组
115     ans = list()
116     # 从lastNode开始，调用dfs，找到所有单词，按顺序储存在ans中
117     lastNode.dfs(pre, ans, [])
118     # 最后将ans用空格隔开合并为字符串后输出
119     print(" ".join(ans))

```

```

1 import java.util.*;
2
3 class TrieNode {
4     TrieNode[] children;
5     boolean isEnd;
6
7     // 构建前缀树节点类
8     public TrieNode() {
9         this.children = new TrieNode[52]; // 大小写均存在, 需要构建长度为52的
        children数组
10         this.isEnd = false; // 结束标识符, True表示当前节点是一个单词
        的结尾
11     }
12
13     // 将单词word加入前缀树的函数
14     public void addWord(String word) {
15         TrieNode node = this;
16         // 遍历该单词中的所有字符
17         for (char ch : word.toCharArray()) {
18             // 获得ch在children数组中对应的索引
19             int chIdx = getIdx(ch);
20             // 如果对应位置为null
21             if (node.children[chIdx] == null) {
22                 // 则为这个ch字符创建一个新的前缀树节点
23                 node.children[chIdx] = new TrieNode();
24             }
25             // 令前缀树节点前进到ch所在的节点
26             node = node.children[chIdx];
27         }
28         // 完成该单词的添加, 设置最后一个字符的节点的结束标识符为true
29         node.isEnd = true;
30     }
31
32     // 根据字符ch获得在children数组中的对应索引的函数
33     public int getIdx(char ch) {
34         // 如果ch是小写, 得到26-51的索引
35         if (Character.isLowerCase(ch)) {
36             return ch - 'a' + 26;
37         } else {
38             // 如果ch是大写, 得到0-25的索引
39             return ch - 'A';
40         }
41     }
42
43     // 根据在children数组中的索引idx获得对应字符ch的函数
44     public char getCh(int idx) {

```

```

45     // 如果idx大于等于26, 是一个小写字母
46     if (idx >= 26) {
47         return (char) (idx + 'a' - 26);
48     } else {
49         // 如果idx小于26, 是一个大写字母
50         return (char) (idx + 'A');
51     }
52 }
53
54 // 获得前缀prefix最后一个字符所在的节点
55 public TrieNode getLastNode(String prefix) {
56     TrieNode node = this;
57     // 遍历prefix中的每个字符
58     for (char ch : prefix.toCharArray()) {
59         int chIdx = getIdx(ch);
60         if (node.children[chIdx] == null) {
61             return null; // 如果某个字符不存在, 返回null
62         }
63         node = node.children[chIdx];
64     }
65     return node; // 返回前缀最后一个字符所在的节点
66 }
67
68 // 对前缀树进行dfs前序遍历, 搜索得到所有后缀
69 public void dfs(String pre, List<String> ans, List<Character> path) {
70     // 遇到一个单词结束标识符, 将当前path合并为字符串后加入ans
71     if (isEnd) {
72         StringBuilder sb = new StringBuilder(pre);
73         for (char ch : path) {
74             sb.append(ch);
75         }
76         ans.add(sb.toString());
77     }
78     // 如果children存在任意一个非null节点, 需要对非空节点继续进行DFS搜索
79     for (int i = 0; i < children.length; i++) {
80         if (children[i] != null) {
81             // 根据索引获得对应的字符
82             char nxtCh = getCh(i);
83             // 将字符加在path末尾的结果, 作为参数传入递归
84             List<Character> newPath = new ArrayList<>(path);
85             newPath.add(nxtCh);
86             children[i].dfs(pre, ans, newPath);
87         }
88     }
89 }
90 }
91

```

```

92 public class Main {
93     public static void main(String[] args) {
94         Scanner scanner = new Scanner(System.in);
95         String s = scanner.nextLine();
96         String pre = scanner.nextLine();
97
98         // 初始化列表lst用于存放所有单词
99         String[] words = s.split("[^a-zA-Z]+");
100         Set<String> set = new HashSet<>(Arrays.asList(words));
101
102         List<String> lst = new ArrayList<>(set);
103         // 初始化前缀树根节点
104         TrieNode root = new TrieNode();
105
106         // 遍历lst中的每一个单词word，构建前缀树
107         for (String word : lst) {
108             root.addWord(word);
109         }
110
111         // 调用前缀树中的getLastNode()方法，得到前缀pre在树中的最后一个节点
112         TrieNode lastNode = root.getLastNode(pre);
113
114         // 如果lastNode为空，说明在root前缀树中，不存在任何前缀为pre的单词，输出pre
115         if (lastNode == null) {
116             System.out.println(pre);
117         } else {
118             // 如果lastNode非空，说明在root前缀树中，存在前缀为pre的单词，要找到所有单
词
119             List<String> ans = new ArrayList<>();
120             // 从lastNode开始，调用dfs，找到所有单词，按顺序储存在ans中
121             lastNode.dfs(pre, ans, new ArrayList<>());
122             // 最后将ans用空格隔开合并为字符串后输出
123             System.out.println(String.join(" ", ans));
124         }
125     }
126 }
127

```

C++

```

1 #include <iostream>
2 #include <vector>
3 #include <unordered_set>
4
5 // 构建前缀树节点类

```

```

6 class TrieNode {
7 public:
8     TrieNode* children[52];
9     bool isEnd;
10
11     TrieNode() {
12         for (int i = 0; i < 52; ++i) {
13             children[i] = nullptr;
14         }
15         isEnd = false; // 结束标识符, true表示当前节点是一个单词的结尾
16     }
17
18     // 将单词word加入前缀树的函数
19     void addWord(const std::string& word) {
20         TrieNode* node = this;
21         for (char ch : word) {
22             int chIdx = getIdx(ch);
23             // 如果对应位置为nullptr
24             if (node->children[chIdx] == nullptr) {
25                 // 则为这个ch字符创建一个新的前缀树节点
26                 node->children[chIdx] = new TrieNode();
27             }
28             // 令前缀树节点前进到ch所在的节点
29             node = node->children[chIdx];
30         }
31         // 完成该单词的添加, 设置最后一个字符的节点的结束标识符为true
32         node->isEnd = true;
33     }
34
35     // 根据字符ch获得在children数组中的对应索引的函数
36     int getIdx(char ch) {
37         if (islower(ch)) {
38             return ch - 'a' + 26; // 如果ch是小写, 得到26-51的索引
39         } else {
40             return ch - 'A'; // 如果ch是大写, 得到0-25的索引
41         }
42     }
43
44     // 根据在children数组中的索引idx获得对应字符ch的函数
45     char getCh(int idx) {
46         if (idx >= 26) {
47             return idx + 'a' - 26; // 如果idx大于等于26, 是一个小写字母
48         } else {
49             return idx + 'A'; // 如果idx小于26, 是一个大写字母
50         }
51     }
52 }

```

```

53 // 获得前缀prefix最后一个字符所在的节点
54 TrieNode* getLastNode(const std::string& prefix) {
55     TrieNode* node = this;
56     for (char ch : prefix) {
57         int chIdx = getIdx(ch);
58         if (node->children[chIdx] == nullptr) {
59             return nullptr;
60         }
61         node = node->children[chIdx];
62     }
63     return node;
64 }
65
66 // 对前缀树进行dfs前序遍历, 搜索得到所有后缀
67 void dfs(const std::string& pre, std::vector<std::string>& ans, const
std::vector<char>& path) {
68     if (isEnd) {
69         std::string word = pre;
70         for (char ch : path) {
71             word += ch;
72         }
73         ans.push_back(word);
74     }
75     for (int i = 0; i < 52; ++i) {
76         if (children[i] != nullptr) {
77             char nxtCh = getCh(i);
78             std::vector<char> newPath(path);
79             newPath.push_back(nxtCh);
80             children[i]->dfs(pre, ans, newPath);
81         }
82     }
83 }
84 };
85
86 int main() {
87     std::string s, pre;
88     std::getline(std::cin, s);
89     std::getline(std::cin, pre);
90
91     // 初始化列表lst用于存放所有单词
92     std::unordered_set<std::string> wordSet;
93     size_t start = 0;
94     while (start < s.size()) {
95         while (start < s.size() && !isalpha(s[start])) {
96             ++start;
97         }
98         size_t end = start;

```

```

99         while (end < s.size() && isalpha(s[end])) {
100             ++end;
101         }
102         if (start < s.size()) {
103             wordSet.insert(s.substr(start, end - start));
104         }
105         start = end + 1;
106     }
107
108     // 对lst进行去重, 因为使用前缀树, 所以无需排序
109     std::vector<std::string> words(wordSet.begin(), wordSet.end());
110     // 初始化前缀树根节点
111     TrieNode* root = new TrieNode();
112
113     // 遍历lst中的每一个单词word, 构建前缀树
114     for (const std::string& word : words) {
115         root->addWord(word);
116     }
117
118     // 调用前缀树中的getLastNode()方法, 得到前缀pre在树中的最后一个节点
119     TrieNode* lastNode = root->getLastNode(pre);
120
121     // 如果lastNode为空, 说明在root前缀树中, 不存在任何前缀为pre的单词, 输出pre
122     if (lastNode == nullptr) {
123         std::cout << pre << std::endl;
124     } else { // 如果lastNode非空, 说明在root前缀树中, 存在前缀为pre的单词, 要找到所有
单词
125         // 初始化答案数组
126         std::vector<std::string> ans;
127         // 从lastNode开始, 调用dfs, 找到所有单词, 按顺序储存在ans中
128         lastNode->dfs(pre, ans, std::vector<char>());
129         for (const std::string& word : ans) {
130             std::cout << word << " ";
131         }
132         std::cout << std::endl;
133     }
134
135     delete root; // Don't forget to release memory
136     return 0;
137 }
138

```

时空复杂度

时间复杂度： $O(NM)$ 。建树、检查前缀的时间复杂度。

空间复杂度： $O(D)$ 。

N 为单词数目， M 为单词平均长度， D 为前缀树的节点数，远小于 NM 。