

# 词典分组

---

## 词典分组

1. 什么是词
2. 词典
  - 2.1 词典的加载
3. 切分算法
  - 3.1 完全切分
  - 3.2 正向最长匹配
  - 3.3 逆向最长匹配
  - 3.4 双向最长匹配
4. 字典树
  - 4.1 字典树定义
  - 4.2 字典树实现
  - 4.3 基于字典树的改进算法
5. HanLP的词典分词实现
  - 5.1 **DoubleArrayTrieSegment**
  - 5.2 去掉停用词

中文分词指的是将一段文本拆分成一系列单词的过程，这些单词顺序拼接后等于原文本。中文分词算法大致分为给予词典规则和基于机器学习这两大流派。

词典分组是最简单最常见的基于规则的分词算法，仅需一部词典和一套查词的规则即可。词典分词的重点不在于分词本身，而在于支撑词典的数据结构。

## 1. 什么是词

---

在语言学上，词语的定义是具备独立意义的最小单位。在基于词典的中文分词中，词典中的字符串就是词。

**齐夫定律：**一个单词的词频与它的词频排名成反比。MSR语料库（微软亚洲研究院语料库）上的统计结果验证了这一定律。



```

1 def load_dictionary():
2     """
3     加载HanLP中的mini词库
4     :return: 一个set形式的词库
5     """
6     IOUtil = JClass('com.hankcs.hanlp.corpus.io.IOUtil')
7     path = HanLP.Config.CoreDictionaryPath.replace('.txt', '.mini.txt')
8     dic = IOUtil.loadDictionary([path])
9     return set(dic.keySet())

```

## 3. 切分算法

词典分词常用的规则有正向最长匹配、逆向最长匹配和双向最长匹配，它们都属于完全切分过程。

### 3.1 完全切分

完全切分指的是，找出一段文本中的所有单词。

```

1 def fully_segment(text, dic):
2     word_list = []
3     for i in range(len(text)):           # i 从 0 到text的最后一个字的
        下标遍历
4         for j in range(i + 1, len(text) + 1): # j 遍历[i + 1, len(text)]
            区间
5             word = text[i:j]             # 取出连续区间[i, j]对应的字符
            串
6             if word in dic:               # 如果在词典中，则认为是一个词
7                 word_list.append(word)
8     return word_list
9
10 dic = load_dictionary()
11 print(fully_segment('商品和服务', dic))

```

输出：

```

1 ['商', '商品', '品', '和', '和服', '服', '服务', '务']

```

### 3.2 正向最长匹配

上面的输出并不是中文分词，中文分词想要的是那种有意义的词语序列，而不是所有出现在词典中的单词所构成的链表。在以某个下标为起点递增查词的过程中，优先输出更长的单词，这种规则被称为**最长匹配算法**。扫描顺序从前往后匹配则称为**正向最长匹配**，反之则称为**逆向最长匹配**。

```

1 def forward_segment(text, dic):
2     word_list = []
3     i = 0
4     while i < len(text):

```

```

5         longest_word = text[i] # 当前扫描位置的单字
6         for j in range(i + 1, len(text) + 1): # 所有可能的结尾
7             word = text[i:j] # 从当前位置到结尾的连续字
字符串
8             if word in dic: # 在词典中
9                 if len(word) > len(longest_word): # 并且更长
10                     longest_word = word # 则更优先输出
11             word_list.append(longest_word) # 输出最长词
12             i += len(longest_word) # 正向扫描
13         return word_list
14
15 dic = load_dictionary()
16 print(forward_segment('就读北京大学', dic))
17 print(forward_segment('研究生命起源', dic))

```

输出：

```

1 ['就读', '北京大学']
2 ['研究生', '命', '起源']

```

可以看出，第二句话产生了误差，希望提取的是“研究”，但是正向最长匹配算法就提取出了“研究生”，所以人们就想出了逆向最长匹配。

### 3.3 逆向最长匹配

逆向最长匹配与正向最长匹配唯一的区别就是扫描的方向不同。

```

1 def backward_segment(text, dic):
2     word_list = []
3     i = len(text) - 1
4     while i >= 0: # 扫描位置作为终点
5         longest_word = text[i] # 扫描位置的单字
6         for j in range(0, i): # 遍历[0, i]区间作为待查
询词语的起点
7             word = text[j: i + 1] # 取出[j, i]区间作为待查
询单词
8             if word in dic:
9                 if len(word) > len(longest_word): # 越长优先级越高
10                     longest_word = word
11                 break
12             word_list.insert(0, longest_word) # 逆向扫描，所以越先查出的
单词在位置上越靠后
13             i -= len(longest_word)
14         return word_list
15
16
17 if __name__ == '__main__':
18     dic = load_dictionary()

```

```
19     print(backward_segment('研究生命起源', dic))
20     print(backward_segment('项目的研究', dic))
```

输出：

```
1  ['研究', '生命', '起源']
2  ['项', '目的', '研究']
```

这次得到了正确的结果，['研究', '生命', '起源']，但是对于“项目的研究”的分词结果又产生了误差，因此，另一些人提出综合两种规则，期待它们取长补短，称为双向最长匹配。

## 3.4 双向最长匹配

这是一种融合两种匹配方法的复杂规则集，流程如下：

- 同时执行正向和逆向最长匹配，若两者的词数不同，则返回词数更少的那一个。
- 否则，返回两者中单字更少的那一个。当单字数也相同时，优先返回逆向最长匹配的结果。

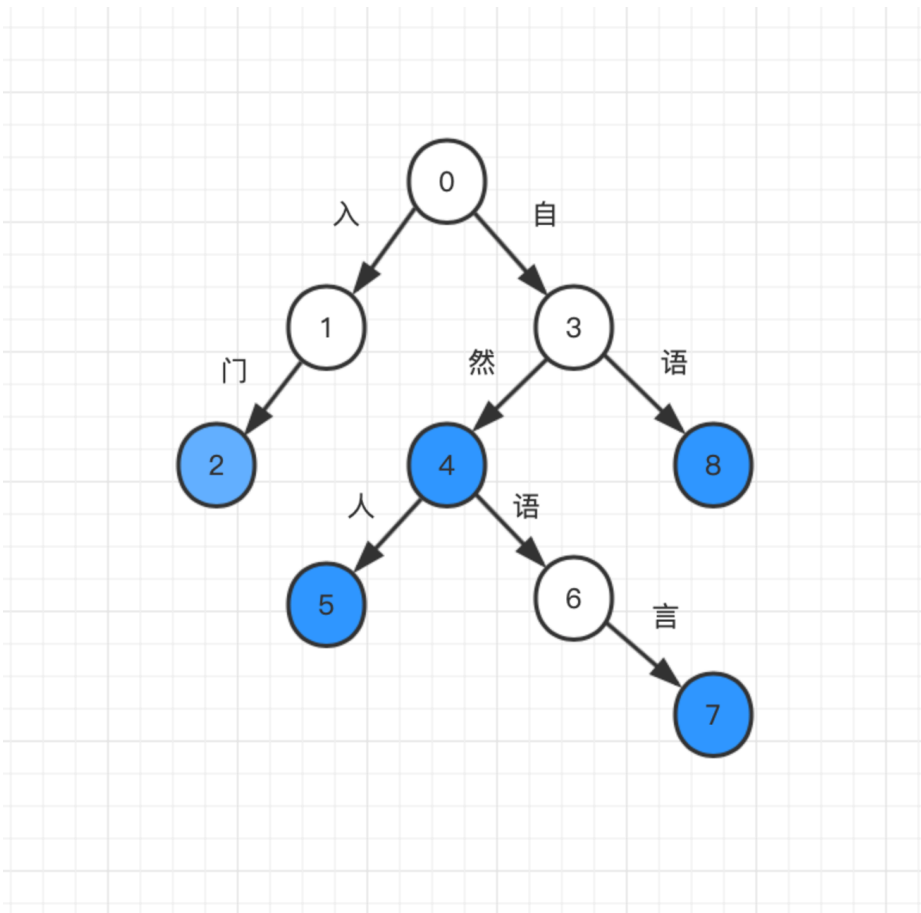
```
1  def count_single_char(word_list: list): # 统计单字成词的个数
2      return sum(1 for word in word_list if len(word) == 1)
3
4
5  def bidirectional_segment(text, dic):
6      f = forward_segment(text, dic)
7      b = backward_segment(text, dic)
8      if len(f) < len(b): # 词数更少优先级更高
9          return f
10     elif len(f) > len(b):
11         return b
12     else:
13         if count_single_char(f) < count_single_char(b): # 单字更少优先级更高
14             return f
15         else:
16             return b # 都相等时逆向匹配优先级更高
17
18     print(bidirectional_segment('研究生命起源', dic))
19     print(bidirectional_segment('项目的研究', dic))
```

## 4. 字典树

### 4.1 字典树定义

匹配算法的瓶颈之一在于如何判断集合(词典)中是否含有字符串。如果用有序集合(TreeMap)的话，复杂度是 $O(\log n)$  ( $n$ 是词典大小)；如果用散列表(Java的HashMap, Python的dict)的话，时间复杂度虽然下降了，但内存复杂度却上去了。有没有速度又快、内存又省的数据结构呢？这就是字典树。

字符串集合常用字典树(又称trie树、前缀树)存储，这是一种字符串上的树形数据结构。字典树中每条边都对应一个字，从根节点往下的路径构成一个个字符串。字典树并不直接在节点上存储字符串，而是将词语视作根节点到某节点之间的一条路径，并在终点节点(蓝色)上做个标记“该节点对应词语的结尾”。字符串就是一条路径，要查询一个单词，只需顺着这条路径从根节点往下走。如果能走到特殊标记的节点，则说明该字符串在集合中，否则说明不存在。一个典型的字典树如下图所示。



其中，蓝色标记着该节点是一个词的结尾，数字是人为的编号。按照路径我们可以得到如下表所示：

词语	路径
入门	0-1-2
自然	0-3-4
自然人	0-3-4-5
自然语言	0-3-4-6-7
自语	0-3-8

## 4.2 字典树实现

由上图可知，每个节点都至少应该知道自己的子节点与对应的边，以及自己是否对应一个词。如果要实现映射而不是集合的话，还需要知道自己对应的值。在HanLP中，约定用None值表示节点不对应词语，这样就不能插入值为None的键了，但是实现起来更加简洁。

节点实现用python描述如下：

```

1  class Node(object):
2      def __init__(self, value) -> None:
3          self._children = {} # 定义子节点
4          self._value = value # 定义当前节点的值
5
6      def _add_child(self, char, value, overwrite=False):
7          child = self._children.get(char)
8          if child is None:
9              child = Node(value) # 创建子节点
10             self._children[char] = child # 子节点赋值, 字->节点的映射
11         elif overwrite:
12             child._value = value # 节点上对应的词
13         return child
14
15
16 class Trie(Node):
17     def __init__(self) -> None:
18         super().__init__(None)
19
20     def __contains__(self, key):
21         return self[key] is not None
22
23     def __getitem__(self, key):
24         state = self
25         for char in key:
26             state = state._children.get(char)
27             if state is None:
28                 return None
29         return state._value
30
31     def __setitem__(self, key, value):
32         state = self
33         for i, char in enumerate(key):
34             if i < len(key) - 1:
35                 state = state._add_child(char, None, False)
36             else:
37                 state = state._add_child(char, value, True) # 最后一个表

```

示词

## 4.3 基于字典树的改进算法

在字典树的基础上，自然语言处理入门的作者何晗老师提出了多种改进优化算法，把分词速度推向了千万字每秒的级别，改进算法包括：

- 首字散列其余二分的字典树
- 双数组字典树
- AC自动机
- 基于双数组字典树的AC自动机

## 5. HanLP的词典分词实现

HanLP中所有的分词器都继承自Segment这个基类。

### 5.1 DoubleArrayTrieSegment

**DoubleArrayTrieSegment**分词器是对DAT最长匹配的封装，默认加载hanlp.properties中CoreDictionaryPath指定的词典。python调用如下：

```
1  from pyhanlp import *
2
3  # 不显示词性，显示词性为
4  HanLP.Config.ShowTermNature = False
5  # 可传入自定义字典，自定义字典放入list中作为参数，如[dir1, dir2]
6  segment = DoubleArrayTrieSegment()
7  # 激活数字和英文识别
8  segment.enablePartOfSpeechTagging(True)
9
10 print(segment.seg("江西鄱阳湖干枯，中国最大淡水湖变成大草原"))
11 print(segment.seg("上海市虹口区大连西路550号SISU"))
```

输出：

```
1  [江西, 鄱阳湖, 干枯, , , 中国, 最大, 淡水湖, 变成, 大草原]
2  [上海市, 虹口区, 大连, 西路, 550, 号, SISU]
```

### 5.2 去掉停用词

```
1  def load_from_file(path):
2      """
3      从词典文件加载DoubleArrayTrie
4      :param path: 词典路径
5      :return: 双数组trie树
6      """
7      map = JClass('java.util.TreeMap')() # 创建TreeMap实例
8      with open(path) as src:
9          for word in src:
10             word = word.strip() # 去掉Python读入的\n
11             map[word] = word
12     return JClass('com.hankcs.hanlp.collection.trie.DoubleArrayTrie')(map)
13
14
15 ## 去掉停用词
16 def remove_stopwords_termlist(termlist, trie):
17     return [term.word for term in termlist if not
18             trie.containsKey(term.word)]
18
```



```
19 trie = load_from_file('stopwords.txt')
20 termlist = segment.seg("江西鄱阳湖干枯了，中国最大的淡水湖变成了大草原")
21 print('去掉停用词前：', termlist)
22
23 print('去掉停用词后：', remove_stopwords_termlist(termlist, trie))
```

输出：

```
1  去掉停用词前： [江西，鄱阳湖，干枯，了，，，中国，最大，的，淡水湖，变成，了，大
   草原]
2  去掉停用词后： ['江西', '鄱阳湖', '干枯', '中国', '最大', '淡水湖', '变成', '大
   草原']
```