

中文分词作业

魏少杭

学号：20373594

学院：人工智能研究院

2022 年 10 月 13 日

目录

1	思路与过程	2
1.1	最大匹配算法实现	2
1.1.1	正向最大匹配 (FMM) 实现细节	2
1.1.2	逆向最大匹配 (BMM) 实现细节	3
1.1.3	双向最大匹配 (BM) 算法	3
1.1.4	模型评估	4
1.2	最少分词法/最短路径法	6
1.2.1	生成有向图	6
1.2.2	最短路径-dijkstra 算法	7
1.2.3	最少分词法整体算法	8
1.2.4	计算评估模型	8
1.3	词云图	9
1.3.1	最大匹配算法-词云图	9
1.3.2	最少分词法-词云图	9
1.4	切分词的条件判断方法	9
2	总结和说明	13
2.1	两类算法的比较	13
2.2	算法上的局限性	13
2.3	基于验证集和结果的改进	13
2.4	避免时间浪费的一些技巧	14

1 思路与过程

1.1 最大匹配算法实现

首先分别通过正向最大匹配和逆向最大匹配，获得各自算法的结果，结果用**分词后的词语表**（list 类型，且存在词语的重复情况）和**分词后的区间表示**。

1.1.1 正向最大匹配（FMM）实现细节

FMM 算法

```
1 #BM.py
2 def FMM(entry_list):
3     """
4     Through FMM, get the partition array(type:list) for words in the range
        given.
5     param: entry_list:待分词的句子列表
6     return: 分词结果词语表FMM_res(list) 、分词区间FMM_intervals(list)
7     """
8     # 初始化FMM中存词典的词
9     FMM_res = []
10    # 为了后续count区间，所以需要把所有的FMM预测区间进行存储
11    FMM_intervals = []
12    sentence_begin = 0
13    line_num = 0
14    for entry in entry_list:
15        # 每一行进行处理，但为了和标准答案能够有对比，在计算区间intervals的
            时候下标仍然是整个文档中的下标
16        line_num += 1
17        print('FMM: ',line_num)      # 打印输出当前处理到第几行了
18        sentence = entry.strip()      # 去除尾部的\n
19        N = len(sentence)
20        i = 0
21        while i < N:      #注意：当用for i in range(N)的时候，每次是从一个迭代
            器里面查找的
22            for j in range(min(N-1, i+21), i-1, -1):
23                if j-i+1 > 22: # 由于词典最长的词长度只有22，这个continue可
                    以减少不必要的时间浪费
24                    continue
25                if judge(sentence, i, j): # 判断是否切分
26                    FMM_res.append(sentence[i:j+1])
27                    FMM_intervals.append((sentence_begin + i, sentence_begin
                        + j))
28                    i = j
29                    break
30            i += 1
```

```

31     sentence_begin += N
32     return (FMM_res, FMM_intervals)

```

传入待分词的句子列表,然后对每一个句子从前往后匹配词语。输出匹配结果 FMM_res 和分词区间 FMM_intervals。

判断是否需要切分一个词的时候,进行了两个 if 条件判断,第一个 if 语句是由于词典最长词长度为 22,故 continue 可以减少不必要的时间浪费。

第二个 if 语句调用了 judge() 函数,i 和 j 分别为当前待判断的区间的左右端点,sentence 为当前的一个句子。为了能够与标准答案有对比,在计算区间 intervals 的时候下标仍然使用整个文档 txt 的下标,所以每次 intervals 列表中插入的时候都是用绝对下标 sentence_begin+i 和 sentence_begin+j 来表示。judge() 函数将在后面进行具体说明。

1.1.2 逆向最大匹配 (BMM) 实现细节

逆向最大匹配算法实现基本与 FMM 一致,其区别主要在于每一个句子的遍历切分的方向不同。每一个句子分词结果的产生顺序是由右向左。

以下是 BMM 算法中与 FMM 算法不同的地方:

BMM 算法

```

1  # BM.py
2      i = N - 1
3      while i >= 0:
4          for j in range(max(0, i-21), i+1):
5              if i-j+1 > 22:
6                  continue
7              # if isWord(sentence[j:i+1]) or i == j:
8              if judge(sentence, j, i): # 判断是否拆分
9                  BMM_res.append(sentence[j:i+1])
10                 BMM_intervals.append((sentence_begin + j, sentence_begin
11                                     + i))
12                 i = j
13                 break
14             i -= 1
15     sentence_begin += N

```

1.1.3 双向最大匹配 (BM) 算法

双向最大匹配就是将分别调用一次 FMM 和 BMM,对返回的分词结果进行比较并返回更优的结果。

首先比较 FMM 和 BMM 算法分词结果数目,选择较小者作为最大匹配的结果。若两个算法结果数目相同,则对分词的词结果去重,即将相同的结果取一次返回。

其次，若去重后的结果不同，则取去重后分词数较少者作为最大匹配的结果。若去重之后的结果相同，则返回 BMM 算法结果。

因为研究表明 BMM 算法的结果从经验上讲准确率更高。

下面是实现：

BM 算法实现

```
1 # BM.py
2 def BM(txt):
3     """
4     输入为一个文章段落，首先对文章进行预处理操作TODO：
5     传入txt文件变量即可
6     res表示BM的结果
7     下面选取两个算法得到的结果中最好的结果
8     flag： 如果为0，返回FMM，如果为1，返回BMM
9     """
10    entry_list = []
11    for entry in txt:
12        entry_list.append(entry.strip())
13    FMM_res, FMM_intervals = FMM(entry_list)
14    BMM_res, BMM_intervals = BMM(entry_list)
15    if len(FMM_res) < len(BMM_res):
16        flag = 0
17    elif len(FMM_res) > len(BMM_res):
18        flag = 1
19    else:
20        # 去重，取去重之后结果最小的那个作为结果
21        if len(list(dict.fromkeys(FMM_res))) < len(list(dict.fromkeys(
22            BMM_res))):
23            flag = 0
24        else:
25            flag = 1
26    if flag == 0:
27        return (FMM_res, FMM_intervals)
28    else:
29        return (BMM_res, BMM_intervals)
```

1.1.4 模型评估

第一部分：获取分词的标准答案

主要在 get_std_intervals() 函数中实现，其输出为分词的标准答案中的区间结果。

实现如下：

get_std_intervals

```
1 # BM.py
```

```

2 def get_std_intervals():
3     with open('./分词对比文件/gold.txt', 'r', encoding='utf-8') as fr:
4         std_string = ''
5         for entry in fr:
6             std_string = std_string + entry.strip('\n')
7         len_passage = len(std_string)
8         std_intervals = []
9         left = 0
10        j = -1
11        i = 0
12        std_words = []
13        while i < len_passage:
14            if i+2 < len_passage and std_string[i+1] == ' ' and std_string[i
15                +2] == ' ':
16                len_string = i - j
17                std_intervals.append((left, left + len_string - 1))
18                std_words.append(std_string[j+1: i+1])
19                left = left + len_string
20                i = i + 2
21                j = i
22            i += 1
23        std_words = std_words
24        return std_intervals

```

在标准答案的区间结果（list 类型）已经获取后，对模型进行计算评估。

在评估之前，必须要注意如何定义 Positive 类和 Negative 类。

Positive 即预测分词区间在标准分词结果区间中的个数；*Negative* 即预测分词区间不在标准分词结果区间中的个数；标准分词结果中的所有区间的总数即为 $TP+FN$ ；预测分词结果中的所有区间的总数即为 $TP+FP$ ；预测分词结果中的所有区间与标准分词结果中的所有区间的交集大小即为 TP ；

第二部分：计算 TP

calc_TP

```

1 # BM.py
2 def calc_TP(std_intervals, prd_intervals):
3     TP = 0
4     for interval in prd_intervals:
5         if interval in std_intervals:
6             TP += 1
7     return TP

```

有了 TP 之后，就可以计算 P,R 和 F1 了。

第三部分，计算 P,R,F1

eval_PRF1

```
1 # BM.py
2     TPaddFN = len(std_intervals)
3     TPaddFP = len(prd_intervals)
4     TP = calc_TP(std_intervals, prd_intervals)
5     precision = TP / TPaddFP
6     recall = TP / TPaddFN
7     F1 = 2 * precision * recall / (precision + recall)
8     return (precision, recall, F1)
```

计算结果在 20373594-魏少杭/评估结果/eval_res_BM.txt 中。即：

Precision:90.1214%

Recall:92.7852%

F1:91.4339%

1.2 最少分词法/最短路径法

1.2.1 生成有向图

由于整个方法是将词模型转换为图模型，所以首先将一个句子转换为有向图，数据结构是一个矩阵 Graph， v_{ij} 为矩阵的元素，即从 i 到 j 结点是否存在路径，换句话说，即从 c_i 到 c_{j-1} 是否存在。初始时，设定每两个直接相邻结点之间存在路径，即 $v_{i,i+1} = 1$ ，矩阵其他元素均为 0。

其中判断是否有路径使用与最大匹配法一样的方法，在 *judge* 函数中实现。

get_graph

```
1 # leastSegment.py
2 def get_graph(sentence):
3     '''
4     input one sentence
5     output a graph:
6     '''
7     # 初始化一个图，是一个点-点矩阵，每一个元素为1或0，表示是否有从一个点到
        另一个点的有向边
8     M = len(sentence) # M表示节点最大下标，恰好为sentence的长度值（因为从0
        开始算）
9     graph = []
10    for i in range(M + 1):
11        graph.append([])
12        for j in range(M + 1):
13            graph[i].append(0)
14        if i < M:
15            graph[i][i + 1] = 1
```

```

16     # 扫描整个句子，判断结点之间是否为有向边
17     for i in range(M):
18         for j in range(i+1, min(i + 23, M + 1)):      # 词最长为22
19             if judge(sentence, i, j-1):
20                 graph[i][j] = 1                        # 如果判断需要分词的话，就分出来
21
22     return graph

```

1.2.2 最短路径-dijkstra 算法

分别用 succession (list) 存最短路径上的前继结点，found (list) 存已经找到最短路径的结点，not_found (list) 存尚未找到最短路径的结点，shortest (ndarray) 存整个图中已经探测到的最短路径长度。

初始化和主要算法部分如下：

dijkstra 算法主体

```

1  # leastSegment.py->dijkstra()
2  # 初始化
3  for i in range(M + 1):
4      succession.append(i - 1)  # 初始化，用i-1来作为前继结点 succession
5      [0] = -1
6      if graph[0][i] == 1:
7          shortest[i] = 1
8  # 开始搜索
9  while M not in found:
10     # 如果从0结点到最后一个结点M的最短路径没有找到的话，不断找最短路径
11     # 从没有找到最短路径的点里面找到当前探测的最短距离对应的第一个点
12     # v_new_found = list(shortest).index(np.min(shortest[not_found]))
13     temp_v_list = np.argwhere(shortest == np.min(shortest[not_found])).
14     reshape(-1)
15     for ver in temp_v_list:
16         if ver in not_found:
17             v_new_found = ver
18             break
19     # 需要新加入的点从not_found里删除，并加入到found中
20     not_found.remove(v_new_found)
21     found.append(v_new_found)
22     # 找这个新加入的点的邻接节点并判断是否更新最短距离
23     # 如果需要更新最短距离，则将这些结点的前继结点进行更新为v_new_found
24     for v in not_found:
25         if v_new_found == M:      # 找到了M结点的最短路径后就结束了
26             break
27         if (graph[v_new_found][v] == 1
28             and 1 + shortest[v_new_found] < shortest[v]):

```

```

27         shortest[v] = shortest[v_new_found] + 1
28         succession[v] = v_new_found
29
30     return get_path(succession, M)

```

其中还定义并调用了 `get_path`，通过传入前继结点集合，返回一条最短路径（结点序列）。

1.2.3 最少分词法整体算法

如下面代码所示，传入文本文件，输出预测的分词结果的列表、预测的分词区间：

`least_segment`

```

1 # leastSegment.py
2 def least_segment(txt):
3     '''
4     传入一个文本文件
5     按行进行get_graph 并 dijkstra 来get_path
6     然后用path来get_intervals作为这一行的intervals
7     '''
8     line_counter = 0
9     prd_intervals = []
10    prd_words = []
11    for entry in txt:
12        sentence = entry.strip()
13        graph = get_graph(sentence)
14        path = dijkstra(graph, len(sentence))
15        line_intervals = get_intervals(path)
16        prd_intervals.append(line_intervals)
17        for (left, right) in line_intervals:
18            prd_words.append(sentence[left: right+1])
19        line_counter += 1
20
21    return (prd_words, prd_intervals)

```

1.2.4 计算评估模型

与最大匹配算法相同的计算思路。

结果如下：

Presicion:87.8124%

Recall:91.8922%

F1:89.8060%

1.3 词云图

选取了鲁迅先生的《铸剑》全文，大约一万字。

1.3.1 最大匹配算法-词云图

word_cloud

```
1 # BM.py
2 def word_cloud():
3     stop_words = [line.strip() for line in open('./停用词/stop_word.txt', 'r',
4         encoding='utf-8')]
5     # 获取想要得到的
6     with open('./喜欢的文章段落/《铸剑》by鲁迅.txt', 'r', encoding='utf-8')
7         as f:
8         words, __ = BM(f) # 将str格式用BM算法来获取分词结果
9         # 去除停用词，得到sentences
10        sentences = ""
11        for word in words:
12            if word in stop_words:
13                continue
14            sentences += str(word)+' '
15        # 生成词云图
16        cloud_img = WordCloud(background_color='white',
17            font_path='./字体/
18            plun04wkyso54jx5893e17t1zkz94ix1.otf',
19            width=2000,
20            height=2000,).generate(sentences)
21        plt.imshow(cloud_img)
22        plt.axis('off')
23        plt.savefig('./词云/《铸剑》by鲁迅')
24        plt.show()
```

该算法实现的词云图如图 1:

1.3.2 最少分词法-词云图

结果如图 2:

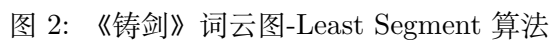
1.4 切分词的条件判断方法

对于一个给定的句子片段，判断为需要切分的词，需要满足以下几种条件之一：

1. 如果这个片段在词典中。
2. 如果这个片段长度仅为 1。



图 1: 《铸剑》词云图-BM 算法



3. 如果这个片段为纯数字。在标准答案中，所有纯数字都分为一个词，且这些纯数字均为半角数字。而在词典中的所有纯数字均为全角数字，故仅查找词典不可能匹配成功。所以需要自己实现一个判断是否全为包含中文、阿拉伯数字的函数 *isAllNum()*。

4. 如果这个片段是日期。具体而言，末尾是‘年’，‘月’，‘日’，‘时’，‘分’等字的时候，就是一个日期。

按照以上思路，构建函数 *judge()*，传入需要判断的句子、待判断的左右端 *i,j*，输出布尔值 *True|False*。

判断是否分词 *judge()*

```
1 # leastSegment.py
2 def judge(sentence, i, j):
3     # 此处的i<=j
4     if (isWord(sentence[i:j+1])
5         or i == j # 如果只有单个字了
6         or (isAllNum(sentence[i : j]) # 如果是 十九万八千 或
7           10238 整个串全是数字
8           and (((sentence[j] in time_tags) or (sentence[j] == '年' and j -
9             i == 4))
10              or isAllNum(sentence[j]))) # 如果是如 10月/23日/23时 或
11              如果是如 2001年或二〇〇一年
12     ):
13         return True
14     return False
```

judge() 调用了判断片段是否全为数字（包括中文数字大写、中文进制数、百分号、全角 0、全角负号、中文小数标志‘点’字符、分数标志‘分’‘之’）的函数 *isAllNum()*。

isAllNum()

```
1 def isAllNum(word):
2     for char in word:
3         if char in Chinese_num:
4             continue
5         elif isNum(char):
6             continue
7         else:
8             return False
9     return True
```

isNum()

```
1 def isNum(char):
2     '''判断是否是数字'''
3     try:
4         float(char)
5         return True
```

```
6     except ValueError:
7         pass
8     try:
9         import unicodedata
10        unicodedata.numeric(char)
11        return True
12    except (TypeError, ValueError):
13        pass
14
15    return False
```

2 总结和说明

2.1 两类算法的比较

双向最大匹配算法的三个评价标准优于最少分词算法。

但是两种算法作为基于词典的算法，在实现时并未经过消歧等操作，所以在实际预测结果中很多地方会出现不符合句子原意的分词方法。这需要额外的消歧手段，或者一些基于统计的算法也应该会有更好的效果。

而最少分词算法还存在多条路径情况下选择最优路径的原则不明确等情况。

2.2 算法上的局限性

在最少分词法中，我当前使用的是传统的单源最短路径的贪心算法-dijkstra 最短路径算法。由于这种贪婪算法实际上对于一个目标点计算出的最短路径只有一条，无法获取多条并依据一定的原则取舍一条最短路径作为输出结果，所以结果不一定可信。

由于时间紧迫，我无法实现广度优先搜索方法 BFS 实现最短路径搜索。实际上，如果一个句子较长时，可能出现许多条从第一个结点到最后一个结点的最短路径。如果使用 BFS 可以避免 dijkstra 算法仅计算一条最短路径的情况。BFS 可以实现一次性多条路径生成，但同样需要处理的问题是，如何制定多条最优路径的选择。经过查阅资料，可能选择标准有：去掉重复词后分词结果个数最少标准；分词结果集合中的分词长度的方差最小标准；分词结果中单字符结果最少标准等。

这两种基于词典的分词方法，效率十分低，时间复杂度极高。每个算法运行时间（cpu 跑）均在 20 分钟以上。

2.3 基于验证集和结果的改进

在验证集的对比文件中，发现了词典中不存在的词，如“多云转晴”在词典中不存在，这极大影响了模型评估的准确率等指标，因此，可以将这些在标准答案中反复出现的词加入到词典中。

但是对于较为大量的训练集或在验证集中出现次数较少的词语来说这种方式不切实际，这需要大量的人工去判断，操作较为复杂。

2.4 避免时间浪费的一些技巧

首先从词典中获取最长的词的长度，这样可以对后面每一次 i,j 所切分出的句子片段进行判断时，指定 i,j 的距离在词典最长词的长度之内，有效避免了长句子中的无效分词带来的时间浪费。