

作业1: 中文的平均信息熵计算

谭天一

ZY2203511

billytan@buaa.edu.cn

1. 实验介绍

阅读论文《An Estimate of an Upper Bound for the Entropy of English》，分别以字和词为单位计算中文的平均信息熵。语料库为16部小说：《白马啸西风》《碧血剑》《飞狐外传》《连城诀》《鹿鼎记》《三十三剑客图》《射雕英雄传》《神雕侠侣》《书剑恩仇录》《天龙八部》《侠客行》《笑傲江湖》《雪山飞狐》《倚天屠龙记》《鸳鸯刀》《越女剑》。

2. 实验过程

实验分为文本加载、文本预处理、字/词频统计、信息熵计算等步骤。

2.1 文本加载

文本加载函数负责读取作为语料库的16部小说，并将其合并为一个字符串 `content_mix`。由于Mac电脑的文件夹下自带难以删除的.Ds_Store文件，因此特意设置只读取txt文件。

```
1 def load_data(path):
2     # 这个程序项目由Mac构建，文件路径采用正斜杠 '/'，在Windows系统运行请改成 '\\'
3     content_mix = ''
4     novels = os.listdir(path)
5     for novel in novels:
6         if 'txt' in novel: # Mac文件夹中含有.Ds_store，应该滤除
7             novel_path = path + '/' + novel
8             with open(novel_path, 'r', encoding='gb18030') as file:
9                 content = file.read()
10                content_mix += content
11                content_mix += '\n'
12            file.close()
13    return content_mix
```

2.2 文本预处理

首先从给定路径中加载停词列表cn_stopwords.txt，将其读取为列表 `stopwords`。基于停词列表对文本加载得到的 `content_mix` 进行处理，删掉停词列表中包含的关键词。由于停词列表中不包含回车符 `'\n'`，因此手动扩展列表 `stopwords`，加入回车符。

```

1 def data_process(content, stopwords_path):
2     with open(stopwords_path, 'r', encoding='utf-8') as file:
3         stopwords = file.read().split('\n')
4     file.close()
5     stopwords.append('\n')
6     for stopword in stopwords:
7         content = content.replace(stopword, '')
8     return content

```

2.3 字/词频统计

一元模型不考虑其他词的影响。

通过get函数对字/词进行计数，然后构建词典 `freq_one`，格式为 `{word: frequency}`。借助lambda表达式依据词典的 `frequency` 对词典的内容降序排列。

字/词频统计函数同时适用于以字为单位和以词为单位，当以字为单位时，输入应为 `<str>` 类型的文本，当以词为单位时，输入应为 `<list>` 类型的已完成分词的文本

```

1 def freq_count_one(content):
2     freq_one = {}
3     for word in content:
4         freq_one[word] = freq_one.get(word, 0) + 1
5     freq_one = dict(sorted(freq_one.items(), key=lambda x: x[1], reverse=True))
6     return freq_one

```

二元模型考虑前一个词和本词之间的关系。

```

1 def freq_count_two(content):
2     freq_two = {}
3     for i in range(len(content) - 1):
4         freq_two[content[i], content[i + 1]] = freq_two.get((content[i], content[i + 1]), 0) + 1
5     freq_two = dict(sorted(freq_two.items(), key=lambda x: x[1], reverse=True))
6     return freq_two

```

三元模型同理。

```

1 def freq_count_three(content):
2     freq_three = {}
3     for i in range(len(content) - 2):
4         freq_three[content[i], content[i + 1], content[i + 2]] = freq_three.get(
5             (content[i], content[i + 1], content[i + 2]), 0) + 1
6     freq_three = dict(sorted(freq_three.items(), key=lambda x: x[1], reverse=True))
7     return freq_three

```

2.4 信息熵计算

一元模型下，信息熵的计算公式为 $H(x) = - \sum_{x \in X} P(x) \log P(x)$ 。

```
1 def entropy_calculate_one(frequency1):
2     entropy = 0
3     content_length = sum(frequency1.values())
4     for word, freq in frequency1.items():
5         prob = freq / content_length
6         entropy -= prob * math.log(prob, 2)
7     return entropy
```

二元模型下，信息熵的计算公式为 $H(X | Y) = - \sum_{x \in X, y \in Y} P(x, y) \log P(x | y)$ ，其中 $P(x | y)$ 为二元词组的频数与以该二元词组的第一个词为首的三元词组的频数的比。而分母，也就是以该二元词组的第一个词为首的三元词组的频数，即可看作该二元词组的第一个词在一元模型下的频数，程序中可以直接以此代替。三元模型同理。

```
1 def entropy_calculate_two(frequency2, frequency1):
2     entropy = 0
3     content_length = sum(frequency2.values())
4     for word1, word2 in frequency2:
5         prob = frequency2[(word1, word2)] / content_length
6         prob_conditional = frequency2[(word1, word2)] / frequency1[word1]
7         entropy -= prob * math.log(prob_conditional, 2)
8     return entropy
```

三元模型下，信息熵的计算公式为 $H(X | Y, Z) = - \sum_{x \in X, y \in Y, z \in Z} P(x, y, z) \log P(x | y, z)$ 。

```
1 def entropy_calculate_three(frequency3, frequency2):
2     entropy = 0
3     content_length = sum(frequency3.values())
4     for word1, word2, word3 in frequency3:
5         prob = frequency3[(word1, word2, word3)] / content_length
6         prob_conditional = frequency3[(word1, word2, word3)] / frequency2[(word1,
7 word2)]
8         entropy -= prob * math.log(prob_conditional, 2)
9     return entropy
```

3. 实验结果

以字为单位	信息熵
一元模型	9.86566397937682
二元模型	6.957636263565023
三元模型	3.511847456690315

以词为单位	信息熵
一元模型	13.488345992329663
二元模型	6.15746702884708
三元模型	1.2965720270876406

4. 总结

随着关联的字词增加，信息熵逐渐变小。当用单字进行统计的时候，由于事实上中文存在很多词语，所以二元模型和三元模型的信息熵仍然较大。而以词为单位进行统计的时候，由于中文的词与词之间的关联不是很紧密，所以二元模型和三元模型的信息熵迅速变小。

本次实验仍然存在一些不足。可能由于停词列表中的关键词不够丰富，导致文本预处理程序对非中文符号的剔除不够彻底，在字/词频统计时统计出了大量奇奇怪怪的符号，可能影响计算结果的精度。

5. 参考文献

1. Brown P F, Della Pietra S A, Della Pietra V J, et al. An estimate of an upper bound for the entropy of English[J]. Computational Linguistics, 1992, 18(1): 31-40.
2. 在完成本次作业时曾阅读过网友“红衣青蛙”的博客https://blog.csdn.net/shzx_55733/article/details/115744123

附录：全部代码

```
1  # @Author: 谭天一
2  # @Description: 根据参考文献Entropy_of_English_PeterBrown和给定语料库计算中文的一元、二元、三元信息熵
3  # @Date: 2023-3-27
4
5  import jieba
6  import math
7  import os
8
9
10 def load_data(path):
11     # 这个程序项目由Mac构建，文件路径采用正斜杠 '/', 在Windows系统运行请改成 '\\'
12     content_mix = ''
13     novels = os.listdir(path)
14     for novel in novels:
15         if 'txt' in novel: # Mac文件夹中含有.Ds_store, 应该滤除
16             novel_path = path + '/' + novel
17             with open(novel_path, 'r', encoding='gb18030') as file:
18                 content = file.read()
19                 content_mix += content
20                 content_mix += '\n'
21             file.close()
22     return content_mix
```

```

23
24
25 def data_process(content, stopwords_path):
26     with open(stopwords_path, 'r', encoding='utf-8') as file:
27         stopwords = file.read().split('\n')
28     file.close()
29     stopwords.append('\n')
30     for stopword in stopwords:
31         content = content.replace(stopword, '')
32     return content
33
34
35 def freq_count_one(content):
36     """
37     同时适用于以字为单位和以词为单位,
38     当以字为单位时, 输入应为<str>类型的文本
39     当以词为单位时, 输入应为<list>类型的已完成分词的文本
40     本注释也适用于下方的freq_count_two(content)、freq_count_three(content)
41     """
42     freq_one = {}
43     for word in content:
44         freq_one[word] = freq_one.get(word, 0) + 1
45     freq_one = dict(sorted(freq_one.items(), key=lambda x: x[1], reverse=True))
46     return freq_one
47
48
49 def freq_count_two(content):
50     freq_two = {}
51     for i in range(len(content) - 1):
52         freq_two[content[i], content[i + 1]] = freq_two.get((content[i],
53 content[i + 1]), 0) + 1
54     freq_two = dict(sorted(freq_two.items(), key=lambda x: x[1], reverse=True))
55     return freq_two
56
57
58 def freq_count_three(content):
59     freq_three = {}
60     for i in range(len(content) - 2):
61         freq_three[content[i], content[i + 1], content[i + 2]] = freq_three.get(
62             (content[i], content[i + 1], content[i + 2]), 0) + 1
63     freq_three = dict(sorted(freq_three.items(), key=lambda x: x[1],
64 reverse=True))
65     return freq_three
66
67
68 def entropy_calculate_one(frequency1):
69     entropy = 0
70     content_length = sum(frequency1.values())
71     for word, freq in frequency1.items():

```

```

70     prob = freq / content_length
71     entropy -= prob * math.log(prob, 2)
72     return entropy
73
74
75 def entropy_calculate_two(frequency2, frequency1):
76     entropy = 0
77     content_length = sum(frequency2.values())
78     for word1, word2 in frequency2:
79         prob = frequency2[(word1, word2)] / content_length
80         prob_conditional = frequency2[(word1, word2)] / frequency1[word1]
81         #  $P(x|y)$  可近似看作每个二元词组在语料库中出现的频数与以该二元词组的第一个词为词首的二元
词组的频数的比值
82         # 分母等价于该二元词组的第一个词在一元词组统计中的频数
83         entropy -= prob * math.log(prob_conditional, 2)
84     return entropy
85
86
87 def entropy_calculate_three(frequency3, frequency2):
88     entropy = 0
89     content_length = sum(frequency3.values())
90     for word1, word2, word3 in frequency3:
91         prob = frequency3[(word1, word2, word3)] / content_length
92         prob_conditional = frequency3[(word1, word2, word3)] / frequency2[(word1,
word2)]
93         #  $P(x|y,z)$  可近似看作每个三元词组在语料库中出现的频数与以该三元词组的前两个词为词首的三
元词组的频数的比值
94         # 分母等价于该三元词组的前两个词在二元词组统计中的频数
95         entropy -= prob * math.log(prob_conditional, 2)
96     return entropy
97
98
99 if __name__ == '__main__':
100     content = load_data('./jyxstxtqj_downcc')
101     content = data_process(content, './cn_stopwords.txt')
102
103     # 以字为单位求一元模型的信息熵
104     frequency1 = freq_count_one(content)
105     entropy1 = entropy_calculate_one(frequency1)
106     print('以字为单位, 一元模型: ', entropy1)
107
108     # 以字为单位求二元模型的信息熵
109     frequency2 = freq_count_two(content)
110     entropy2 = entropy_calculate_two(frequency2, frequency1)
111     print('以字为单位, 二元模型: ', entropy2)
112
113     # 以字为单位求三元模型的信息熵
114     frequency3 = freq_count_three(content)
115     entropy3 = entropy_calculate_three(frequency3, frequency2)

```

```
116     print('以字为单位, 三元模型: ', entropy3)
117
118     # jieba分词
119     content_cut = jieba.lcut(content)
120
121     # 以词为单位求一元模型的信息熵
122     frequency1 = freq_count_one(content_cut)
123     entropy1 = entropy_calculate_one(frequency1)
124     print('以词为单位, 一元模型: ', entropy1)
125
126     # 以词为单位求二元模型的信息熵
127     frequency2 = freq_count_two(content_cut)
128     entropy2 = entropy_calculate_two(frequency2, frequency1)
129     print('以词为单位, 二元模型: ', entropy2)
130
131     # 以词为单位求三元模型的信息熵
132     frequency3 = freq_count_three(content_cut)
133     entropy3 = entropy_calculate_three(frequency3, frequency2)
134     print('以词为单位, 三元模型: ', entropy3)
135
```