

上海对外经贸大学



机器学习与数据挖掘实战课程项目

基于词频统计法和 TextRank 算法的抽取式单文本自动摘要

学院：统计与信息学院

班级：应用统计 1701

姓名：朱强强

学号：17064001

2019 年 12 月

一、项目简介

文本摘要任务旨在从一篇或多篇相同主题的文本中抽取能够迅速反映主题的精简压缩版本，可以帮助用户快速形成对特定主题文本内容的全面了解，提高浏览信息和获取知识的效率。文本摘要根据抽取方式的不同可分为抽取式摘要和生成式摘要。抽取式摘要直接从原文本中不加修改地抽取文本片段（通常是句子）组成摘要；生成式摘要是重新组织句子形成比抽取式摘要更加精简的形式。

目前中文生成式自动摘要的公开数据集只有 LCSTS 和 NLPCC 数据集，由于数据集申请并未收到回复，没有数据集训练中文生成式自动摘要模型，所以本项目为基于词频统计算法和 TextRank 算法的中文抽取式自动摘要。

二、算法介绍

1. TF-IDF

判断句子重要性一个很常见的方法就是统计关键词在句子中出现的频率，而查找关键词最常用的方法之一就是 TF-IDF 算法。

$$\text{词频}(TF) = \frac{\text{某个词在文章中的出现次数}}{\text{文章的总词数}}$$

但是，仅仅根据词频抽出的关键词很可能是一些像“的”，“是”，“在”这一类最常用的词，这对结果毫无帮助，需要被过滤掉。

所以我们还需要一个重要性系数来进一步衡量一个词是不是常见词。如果某个词比较少见，但是它在这篇文章中多次出现，那么它很可能就反映了这篇文章的特性，这才是我们需要的关键词。这个重要性系数被称为逆文档概率（IDF）。

$$\text{逆文档概率}(IDF) = \log \left(\frac{\text{语料库中的文档总数}}{\text{包含该词的文档数} + 1} \right)$$

$$TF - IDF = TF \times IDF$$

TF-IDF 与一个词在文档中的出现次数成正比，与该词在整个语言中的出现次数成反比。计算出文档中每个词的 TF-IDF 值，按降序排列取排在最前面的几个词就是关键词。

2. 余弦相似性

余弦相似性通过测量两个向量的夹角的余弦值来度量它们之间的相似性。我们将句子通过词向量的方法转化为向量形式，然后计算不同向量之间的余弦值。余弦值越接近 1，就表明两个向量所代表的句子相似性越高。

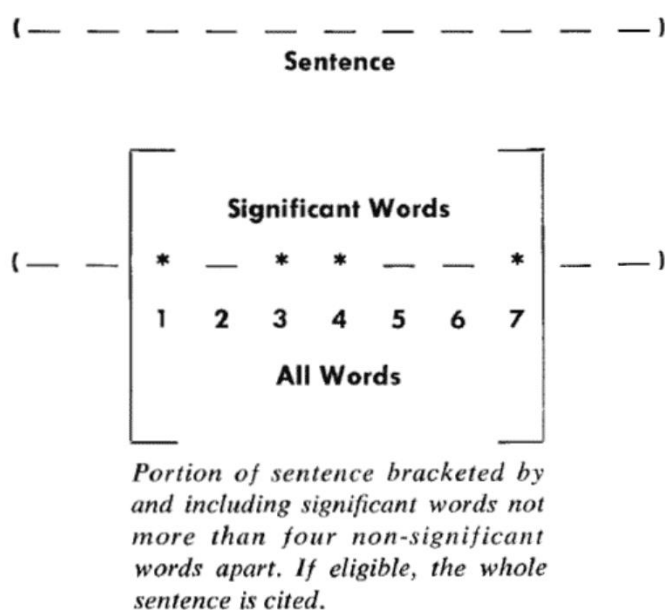
假设 A 和 B 是两个 n 维向量 A 是 $[A_1, A_2, \dots, A_n]$ ，B 是 $[B_1, B_2, \dots, B_n]$ ，则 A 与 B 的夹角 θ 的余弦等于：

$$\cos\theta = \frac{\sum_{i=1}^n (A_i \times B_i)}{\sqrt{\sum_{i=1}^n A_i^2} \times \sqrt{\sum_{i=1}^n B_i^2}} = \frac{A \cdot B}{|A| \times |B|}$$

3. 词频统计方法

词频统计方法最早出自 1958 年的 IBM 公司科学家 H.P. Luhn 的论文《The Automatic Creation of Literature Abstracts》。

作者认为，文章的信息都包含在句子中，有些句子包含的信息越多，就说明这个句子越重要。Luhn 提出用“簇”来表示关键词的聚集。



上图中被框起来的部分就是一个“簇”，只有关键词之间的距离小于一个阈值，它们就被认为处于同一个簇之中。Luhn 建议的阈值是 4 或 5，即如果两个关键词之间有 4 个或 5 个以上的关键词，就可以把这两个关键词分在两个簇中。

之后对于每一个簇，计算它们的重要性得分。

$$\text{簇的重要性} = \frac{(\text{包含的关键词数量})^2}{\text{簇的长度}}$$

找出包含分值最高的簇的句子，把它们结合在一起就构成了该文章的自动摘要。

4. TextRank 算法

TextRank 算法是一种基于图的用于文本自动摘要的排序算法，通过把文本分割成若干个组成单元（句子），构建节点连接图，用句子之间的相似度作为边的权重，通过循环迭代计算句子的 TextRank 值，最后抽取排名高的句子合成文本摘要。

$$TR(v_i) = \frac{1-d}{n} + d \left(\sum_{v_j \in In(v_i)} \frac{w_{ji}}{\sum_{v_k \in Out(v_j)} w_{jk}} TR(v_j) \right), i = 1, 2, \dots, n$$

d 是阻尼系数，代表从某一句子指向其他任意句子的概率，一般设置为 0.85。

w_{ji} 为边的权重，也就是句子的之间的相似度。

$In(v_j)$ 是指向句子 v_j 的句子集合。

$Out(v_j)$ 是句子 v_j 指出去的句子集合。

三、 算法思路

词频统计法

1. 进行文本预处理，将文本切割成单个句子，并进行分词。
2. 去除停用词，使用的是哈工大公开的停用词表。
3. 利用 TF-IDF 算法找出关键词。本文由于缺少语料库，无法计算逆文档概率，故仅使用 TF 寻找关键词。
4. 计算簇的重要性得分，并为句子的重要性打分并降序排序，由此得出的前 n 个句子组成该文档的摘要。

TextRank 算法

1. 文本预处理步骤同词频统计法。
2. 加载经过训练的 word2vec 词向量，并由此得到句子的向量表示。
3. 计算句子之间的余弦相似度，构成相似度矩阵。
4. 利用句子相似度矩阵构建图结构，句子为节点，句子相似度为转移概率。
5. 得到所有句子的 TextRank 值，并由此对句子进行排序，由此得出的前 n 个句子组成该文档的摘要。

四、 算法结果

-----top_n_summary-----

中共中央政治局常委、国务院总理李克强日前就在自由贸易试验区开展“证照分离”改革全覆盖试点作出重要批示。

批示指出：推进“证照分离”改革，是深化“放管服”改革、优化营商环境、更大激发市场活力的重大举措。

必须持续推进市场化、法治化、国际化营商环境建设，更大力度推进“放管服”改革，扎实做好“证照分离”改革试点扩面工作，着力推进照后减证和简化审批，

坚决克服“准入不准营”现象，促进更多新企业开办和发展壮大。

韩正强调，要高质量开展自贸试验区“证照分离”改革全覆盖试点，提升改革工作成效。

要把握好“进四扇门”的办法，按照直接取消审批、审批改为备案、实行告知承诺、优化审批服务等方式，从实际出发，分类实施、稳步推进改革。

-----mean_scored_summary-----

批示指出：推进“证照分离”改革，是深化“放管服”改革、优化营商环境、更大激发市场活力的重大举措。

韩正强调，要高质量开展自贸试验区“证照分离”改革全覆盖试点，提升改革工作成效。

-----TextRank-----

韩正表示，推进“证照分离”改革，是深入贯彻落实党的十九届四中全会精

神，推进国家治理体系和治理能力现代化的一项重要基础性工作。

必须持续推进市场化、法治化、国际化营商环境建设，更大力度推进“放管服”改革，扎实做好“证照分离”改革试点扩面工作，着力推进照后减证和简化审批，

坚决克服“准入不准营”现象，促进更多新企业开办和发展壮大。

韩正强调，要高质量开展自贸试验区“证照分离”改革全覆盖试点，提升改革工作成效。

批示指出：推进“证照分离”改革，是深化“放管服”改革、优化营商环境、更大激发市场活力的重大举措。

要把握好全覆盖这个基本要求，建立清单管理制度，将涉企经营许可事项全部纳入清单管理，明确改革目标，增加改革透明度，倒逼改革深化。

两个算法得出的摘要结果存在一定差异，但目前对抽取式摘要精确度还并没有一个统一、精确的度量方法，只能靠人工判别。

Code

```
1  # -*- coding:utf-8 -*-
2
3  import nltk
4  import jieba
5  import codecs
6  import numpy as np
7  import networkx as nx
8  from sklearn.metrics.pairwise import cosine_similarity
9
10
11 # 文本预处理
12 def sent_tokenizer(texts):
13     start = 0
14     i = 0 # 每个字符的位置
15     sentences = []
16     punt_list = ".!?. ! ? ; "
17     for text in texts:
18         if i < len(texts) - 1:
19             if text in punt_list:
20                 sentences.append(texts[start:i + 1]) # 当前标点符
号位置
21                 start = i + 1 # start标记到下一句的开头
22                 i += 1
23             else:
24                 i += 1 # 若不是标点符号，则字符位置继续前移
25         else:
26             sentences.append(texts[start:]) # 处理文本末尾没有标点
符号的情况
27         break
28     return sentences
29
30
31 # 停用词
32 def load_stop_words_list(path):
33     stop_list = [
34         line.strip()
35         for line in codecs.open(path, 'r',
encoding='utf8').readlines()
36     ]
37     stop_words = {}.fromkeys(stop_list)
38     return stop_words
39
40
41 # 摘要
42 def summarize(text):
```

```

43 stopwords = load_stop_words_list("stopwords.txt")
44 sentences = sent_tokenizer(text)
45 words = [
46     w for sentence in sentences for w in jieba.cut(sentence)
47     if w not in stopwords if len(w) > 1 and w != '\t'
48 ]
49 word_frequency = nltk.FreqDist(words)
50 top_words = [
51     w[0] for w in sorted(
52         word_frequency.items(), key=lambda d: d[1],
reverse=True)
53    ][:n]
54 scored_sentences = _score_sentences(sentences, top_words)
55 # 利用均值和标准差过滤非重要句子
56 avg = np.mean([s[1] for s in scored_sentences])
57 std = np.std([s[1] for s in scored_sentences])
58 mean_scored = [(sent_idx, score) for (sent_idx, score) in
scored_sentences
59                 if score > (avg + 1.5 * std)]
60 top_n_scored = sorted(scored_sentences,
61                       key=lambda s: s[1])[-
num_top_sentence:]
62 top_n_scored = sorted(top_n_scored, key=lambda s: s[0])
63 return dict(
64     top_n_summary=[sentences[idx] for (idx, score) in
top_n_scored],
65     mean_scored_summary=[sentences[idx] for (idx, score) in
mean_scored])
66
67
68 # 句子得分
69 def _score_sentences(sentences, top_words):
70     scores = []
71     sentence_idx = -1
72     for s in [list(jieba.cut(s)) for s in sentences]:
73         sentence_idx += 1
74         word_idx = []
75         for index in range(len(s)):
76             if s[index] in top_words:
77                 word_idx.append(index)
78             else:
79                 pass
80         word_idx.sort()
81         if len(word_idx) == 0:
82             continue
83         # 对于两个连续的单词，利用单词位置索引，通过距离阈值计算簇
84         clusters = []
85         cluster = [word_idx[0]]
86         i = 1

```



```

87         while i < len(word_idx):
88             if word_idx[i] - word_idx[i - 1] < cluster_distance:
89                 cluster.append(word_idx[i])
90             else:
91                 clusters.append(cluster[:])
92                 cluster = [word_idx[i]]
93             i += 1
94         clusters.append(cluster)
95         # 对每个簇打分，每个簇类的最大分数是对句子的打分
96         max_cluster_score = 0
97         for c in clusters:
98             key_words = len(c)
99             total_words = c[-1] - c[0] + 1
100             score = key_words * key_words / total_words
101             if score > max_cluster_score:
102                 max_cluster_score = score
103             scores.append((sentence_idx, max_cluster_score))
104         return scores
105
106
107 # TextRank算法
108 def text_rank_summarize(text):
109     stopwords = load_stop_words_list('E:/Python Files/stopwords/
中文停用词表.txt')
110     sentences = sent_tokenizer(text)
111     words_list = []
112     for sentence in sentences:
113         word_list = [
114             w for w in jieba.cut(sentence) if w not in stopwords
115             if len(w) > 1 and w != '\t'
116         ]
117         words_list.append(word_list)
118
119     # 加载word2vec词向量
120     word_embeddings = {}
121     with codecs.open("sgns.renmin.char", encoding='utf-8') as f:
122         for line in f:
123             # 把第一行的内容去掉
124             if '355996 300\n' not in line:
125                 values = line.split()
126                 word = values[0]
127                 embedding = np.asarray(values[1:],
dtype='float32')
128                 word_embeddings[word] = embedding
129         f.close()
130
131     # 得到句子的向量表示
132     sentence_vectors = []
133     for word_list in words_list:

```

```

134         if len(word_list) != 0:
135             # 如果句子中的词语不在字典中，那就把embedding设为300维元素
            为0的向量。
136             # 得到句子中全部词的词向量后，求平均值，得到句子的向量表示
137             v = sum(
138                 [word_embeddings.get(w, np.zeros([
139                     300,
140                     ])) for w in word_list]) / (len(word_list))
141         else:
142             # 如果句子为[],那么就向量表示为300维元素为0个向量。
143             v = np.zeros([
144                 300,
145             ])
146             sentence_vectors.append(v)
147
148         # 计算句子之间的余弦相似度，构成相似度矩阵
149         sim_mat = np.ones([len(words_list), len(words_list)])
150         for i in range(len(words_list)):
151             for j in range(len(words_list)):
152                 if i != j:
153                     sim_mat[i][j] = cosine_similarity(
154                         sentence_vectors[i].reshape(1, 300),
155                         sentence_vectors[j].reshape(1, 300))[0, 0]
156
157         # 利用句子相似度矩阵构建图结构，句子为节点，句子相似度为转移概率
158         nx_graph = nx.from_numpy_array(sim_mat)
159
160         # 得到所有句子的TextRank值
161         scores = nx.pagerank(nx_graph)
162
163         # 根据TextRank值对句子进行排序
164         ranked_sentences = sorted(
165             ((scores[i], s) for i, s in enumerate(sentences)),
            reverse=True)
166
167         result = []
168         for i in range(num_top_sentence):
169             result.append(ranked_sentences[i][1])
170
171         return result
172
173
174     # 读取文本文件
175     def read_text(path):
176         with codecs.open(path, "r", encoding="utf-8") as f:
177             text = f.read()
178         return text
179
180

```

```
181 if __name__ == '__main__':
182     n = 30 # 关键词数量
183     cluster_distance = 5 # 单词间的距离
184     num_top_sentence = 5 # 返回句子的数量
185     text = read_text("text.txt")
186     algorithm1 = summarize(text)
187     algorithm2 = text_rank_summarize(text)
188     print('-----top_n_summary-----')
189     for sent in algorithm1['top_n_summary']:
190         print(sent)
191     print('-----mean_scored_summary-----')
192     for sent in algorithm1['mean_scored_summary']:
193         print(sent)
194     print('-----TextRank-----')
195     for sent in algorithm2:
196         print(sent)
197
```