

4.中文分词&新词发现

1. 正向最大匹配

a. 步骤

i. 实现方式一

ii. 实现方式2

2. 反向最大匹配

3. 双向最大匹配

4. 分词的评价标准

5. 基于词表的分词缺点

6. 基于机器学习的分词

7. 其他

8. 新词发现

9. 重要词评价

a. TF 词频

b. IDF 逆文档频率

c. 特点

d. 应用场景

10. TF-IDF优势

11. TF-IDF劣势

12. 编程练习：

13. 代码分析

a. forward_segmentation_method1

b. forward_segmentation_method2

1. 正向最大匹配

a. 步骤

1. 收集一个词表

2. 对于一个待分词的字符串，从前向后寻找最长的，在此表中出现的词，在词边界做切分

3. 从切分处重复步骤2，直到字符串末尾

• 词表：

• 北京	生前
• 北京大学	前来
• 大学	报到
• 大学生	

句子：

北京大学学生前来报到

北京大学 / 生前 / 来 / 报到

i. 实现方式一

- 1. 找出词表中最大词长度
- 2. 从字符串开头开始选取最大词长度的窗口，检查窗口内的词是否在词表中
- 3. 如果在词表中，在词边界处进行切分，之后移动到词边界处，重复步骤2
- 4. 如果不在词表中，窗口右边界回退一个字符，之后检查窗口词是否在词表中

ii. 实现方式2

- 1. 从前向后进行查找
- 2. 如果窗口内的词是一个词前缀则继续扩大窗口
- 3. 如果窗口内的词不是一个词前缀，则记录已发现的词，并将窗口移动到词边界

2. 反向最大匹配

与正向最大匹配类似，但从后向前进行匹配

3. 双向最大匹配

同时进行正向和反向最大匹配，并比较二者结果

4. 分词的评价标准

- 如何比较?
- 1.单字词
 - 词表中可以有单字，从分词的角度，我们也会把它称为一个词
- 2.非字典词
 - 未在词表中出现过的词，一般都会被分成单字
- 3.词总量
 - 不同切分方法得到的词数可能不同
- 计算哪种切分方式总词频最高
- 词频事先根据分词后语料统计出来

5. 基于词表的分词缺点

1. 对词表极为依赖，如果没有词表，则无法进行；如果词表中缺少需要的词，结果也不会正确
2. 切分过程中不会关注整个句子表达的意思，只会将句子看成一个个片段
3. 如果文本中出现一定的错别字，会造成一连串影响
4. 对于人名等的无法枚举实体词无法有效的处理

6. 基于机器学习的分词

- 对于每一个字，我们想知道它是不是一个词的边界

- 上 | 海 | 自 | 来 | 水 | 来 | 自 | 海 | 上
- 0 1 0 0 1 0 1 1 1
- 蓝色表示不是词边界，红色表示是词边界

我们只需要得到一个与字符串等长的0, 1序列

7. 其他

- 自定义Dataset最关键的是重写__len__和__getitem__两个函数
- RNN中已经带有tanh激活函数，并不需要再定义激活函数

8. 新词发现

如YYDS、南方小土豆、泼天的富贵

- 假设没有词表，如何从文本中发现新词？
- 随着时间推移，新词会不断出现，固有词表会过时
- 补充词表有利于下游任务
- 词相当于一种固定搭配

- 词的内部应该是稳固的
- 内部凝固度

$$\frac{1}{n} \log \frac{p(W)}{p(c_1) \cdots p(c_n)}$$

- 词的外部应该是多变的
- 左右熵

$$H(U) = E[-\log p_i] = - \sum_{i=1}^n p_i \log p_i$$

1. $p(W)$ 表示该词的出现概率， $p(c_1)$ 表示词中的字出现的概率，如“魑魅魍魉”，作为单个字是十分少见的（概率低），但作为整个词概率则很高，因此内部凝固度高
2. $H(U)$ 表示左右熵，描述的是词左右的混乱程度， p_i 表示该词左边出现的字的概率（因为我们假设不进行分词，如“但是”，放在广泛的场景都可以使用，左右搭配可以很混乱，因此也是外部是多变的。

作为新词发现时，将上面两个指标做加权，这是为数不多的无监督任务

```
1  import math
2  from collections import defaultdict
3
4  class NewWordDetect:
5      def __init__(self, corpus_path):
6          self.max_word_length = 5
7          # 记录词频
8          self.word_count = defaultdict(int)
9          # 记录词左邻词频, 嵌套的dict
10         # key          value
11         # word          dict
12         #              key          value
13         #              leftWord      count
14         self.left_neighbor = defaultdict(dict)
15         # 记录词右邻词频, 同上
16         self.right_neighbor = defaultdict(dict)
17         self.load_corpus(corpus_path)
18         self.calc_pmi()
19         self.calc_entropy()
20         self.calc_word_values()
21
22
23         #加载语料数据, 并进行统计
24     def load_corpus(self, path):
25         with open(path, encoding="utf8") as f:
26             for line in f:
27                 sentence = line.strip()
28                 for word_length in range(1, self.max_word_length):
29                     self.ngram_count(sentence, word_length)
30
31         return
32
33     #按照窗口长度取词, 并记录左邻右邻
34     def ngram_count(self, sentence, word_length):
35         for i in range(len(sentence) - word_length + 1):
36             word = sentence[i:i + word_length]
37             self.word_count[word] += 1
38             # 如果大于1, 说明有左邻
39             if i - 1 >= 0:
40                 char = sentence[i - 1]
41                 self.left_neighbor[word][char] = self.left_neighbor[word].get(char, 0) + 1
42             # 如果小于length, 说明有右邻
43             if i + word_length < len(sentence):
44                 char = sentence[i + word_length]
```

```

44         self.right_neighbor[word][char] = self.right_neighbor[word].get(char, 0) + 1
45         return
46
47     #计算熵
48     def calc_entropy_by_word_count_dict(self, word_count_dict):
49         total = sum(word_count_dict.values())
50         entropy = sum([-(c / total) * math.log((c / total), 10) for c in word_count_dict.values()])
51         return entropy
52
53     #计算左右熵
54     def calc_entropy(self):
55         self.word_left_entropy = {}
56         self.word_right_entropy = {}
57         for word, count_dict in self.left_neighbor.items():
58             self.word_left_entropy[word] = self.calc_entropy_by_word_count_dict(count_dict)
59         for word, count_dict in self.right_neighbor.items():
60             self.word_right_entropy[word] = self.calc_entropy_by_word_count_dict(count_dict)
61
62
63     #统计每种词长下的词总数
64     def calc_total_count_by_length(self):
65         self.word_count_by_length = defaultdict(int)
66         # word表示词, count表示该词出现的个数
67         for word, count in self.word_count.items():
68             self.word_count_by_length[len(word)] += count
69         return
70
71     #计算互信息(pointwise mutual information)
72     def calc_pmi(self):
73         self.calc_total_count_by_length()
74         self.pmi = {}
75         # 词频是在对应词长度下单独计算的
76         for word, count in self.word_count.items():
77             p_word = count / self.word_count_by_length[len(word)]
78             p_chars = 1
79             # 统计词中的字频乘积
80             for char in word:
81                 p_chars *= self.word_count[char] / self.word_count_by_length[1]
82             self.pmi[word] = math.log(p_word / p_chars, 10) / len(word)
83         return
84
85     def calc_word_values(self):

```

```

86     self.word_values = {}
87     for word in self.pmi:
88         if len(word) < 2 or ", " in word:
89             continue
90         pmi = self.pmi.get(word, 1e-3)
91         le = self.word_left_entropy.get(word, 1e-3)
92         re = self.word_right_entropy.get(word, 1e-3)
93         # 这里直接乘积并不好，因为左右混论程度会相互抵消，应该取最小值，才能消除
        左右两端的相互影响
94         self.word_values[word] = pmi * min(le, re)
95
96 if __name__ == "__main__":
97     nwd = NewWordDetect("sample_corpus.txt")
98     # print(nwd.word_count)
99     # print(nwd.left_neighbor)
100    # print(nwd.right_neighbor)
101    # print(nwd.pmi)
102    # print(nwd.word_left_entropy)
103    # print(nwd.word_right_entropy)
104    # 对每个词的评分进行排序，并分别输出长度为2, 3, 4的前10个词
105    value_sort = sorted([(word, count) for word, count in nwd.word_value
        s.items()], key=lambda x:x[1], reverse=True)
106    print([x for x, c in value_sort if len(x) == 2][:10])
107    print([x for x, c in value_sort if len(x) == 3][:10])
108    print([x for x, c in value_sort if len(x) == 4][:10])

```

corpus: 语料库

9. 重要词评价

- 假如一个词在某类文本（假设为A类）中出现次数很多，而在其他类别文本（非A类）出现很少，那么这个词是A类文本的重要词（高权重词）。

恒星、黑洞 → 天文

- 反之，如果一个词出现在很多领域，则其对于任意类别的重要性都很差。

	政治?	?
中国	地理?	你好?
	经济?	?
	足球?	?

NLP中经典统计值: $TF \cdot IDF$

a. TF 词频

TF = 某个词在某类别中出现的次数/该类别词总数

b. IDF 逆文档频率

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

+1是为了防止分母为0。

$TF \cdot IDF$ 高代表该领域重要程度高！


```

1  import jieba
2  import math
3  import os
4  import json
5  from collections import defaultdict
6
7  """
8  tfidf的计算和使用
9  """
10
11
12  # 统计tf和idf值
13  def build_tf_idf_dict(corpus):
14      tf_dict = defaultdict(dict)
15      # key          value
16      # 文档序号      dict
17      #              key          value
18      #              词           词出现的次数
19      idf_dict = defaultdict(set)
20      # key          value
21      # 词           set : 文档序号, 最终用于计算每个词在多少篇文档中出现过
22      for text_index, text_words in enumerate(corpus):
23          for word in text_words:
24              if word not in tf_dict[text_index]:
25                  tf_dict[text_index][word] = 0
26                  tf_dict[text_index][word] += 1
27                  idf_dict[word].add(text_index)
28      # 计算每个词在多少个文档中出现过
29      idf_dict = dict([(key, len(value)) for key, value in idf_dict.items()])
30      return tf_dict, idf_dict
31
32
33  # 根据tf值和idf值计算tfidf
34  def calculate_tf_idf(tf_dict, idf_dict):
35      tf_idf_dict = defaultdict(dict)
36      # key          value
37      # 文档序号      dict
38      #              key          value
39      #              词           TF·IDF值
40      for text_index, word_tf_count_dict in tf_dict.items():
41          for word, tf_count in word_tf_count_dict.items():
42              # 每个词所在文档的词频
43              tf = tf_count / sum(word_tf_count_dict.values())

```

```

44         # tf-idf = tf * log(D/(idf + 1))
45         tf_idf_dict[text_index][word] = tf * math.log(len(tf_dict) /
(idf_dict[word] + 1))
46     return tf_idf_dict
47
48
49 # 输入语料 list of string
50 def calculate_tfidf(corpus):
51     # 先进行分词，自动输出分词结果，并用","作为间隔
52     corpus = [jieba.lcut(text) for text in corpus]
53     tf_dict, idf_dict = build_tf_idf_dict(corpus)
54     tf_idf_dict = calculate_tf_idf(tf_dict, idf_dict)
55     return tf_idf_dict
56
57
58 # 根据tfidf字典，显示每个领域topK的关键词
59 def tf_idf_topk(tfidf_dict, paths=[], top=10, print_word=True):
60     topk_dict = {}
61     for text_index, text_tfidf_dict in tfidf_dict.items():
62         # dict.items
63         # 这里的key相当于java中Arrays.sort的比较器，这里是根据第二个维度的元素(词的
TF·IDF值)进行比较
64         word_list = sorted(text_tfidf_dict.items(), key=lambda x: x[1], re
verse=True)
65         topk_dict[text_index] = word_list[:top]
66         if print_word:
67             print(text_index, paths[text_index])
68             for i in range(top):
69                 print(word_list[i])
70             print("-----")
71     return topk_dict
72
73
74 def main():
75     dir_path = r"category_corpus/"
76     corpus = []
77     paths = []
78     for path in os.listdir(dir_path):
79         path = os.path.join(dir_path, path)
80         if path.endswith(".txt"):
81             corpus.append(open(path, encoding="utf8").read())
82             paths.append(os.path.basename(path))
83     tf_idf_dict = calculate_tfidf(corpus)
84     tf_idf_topk(tf_idf_dict, paths)
85
86
87 if __name__ == "__main__":

```

c. 特点

- 1. $tf-idf$ 的计算非常依赖分词结果，如果分词出错，统计值的意义会大打折扣
- 2. 每个词，对于每篇文档，有不同的 $tf-idf$ 值，所以不能脱离数据讨论 $tfidf$
- 3. 假如只有一篇文本，不能计算 $tf-idf$
- 4. 类别数据均衡很重要
- 5. 容易受各种特殊符号影响，最好做一些预处理

d. 应用场景

• 搜索引擎

- 1. 对于已有的所有网页（文本），计算每个网页中，词的 $TFIDF$ 值
- 2. 对于一个输入 $query$ 进行分词
- 3. 对于文档 D ，计算 $query$ 中的词在文档 D 中的 $TFIDF$ 值总和，作为 $query$ 和文档的相关性得分

```
1  import jieba
2  import math
3  import os
4  import json
5  from collections import defaultdict
6  from calculate_tfidf import calculate_tfidf, tf_idf_topk
7
8  """
9  基于tfidf实现简单搜索引擎
10 """
11
12 jieba.initialize()
13
14
15 # 加载文档数据（可以想象成网页数据），计算每个网页的tfidf字典
16 ▼ def load_data(file_path):
17     corpus = []
18     ▼ with open(file_path, encoding="utf8") as f:
19         documents = json.loads(f.read())
20     ▼ for document in documents:
21         # 将标题和内容进行拼接，这里相当于标题和内容看成一样的权重
22         corpus.append(document["title"] + "\n" + document["content"])
23     tf_idf_dict = calculate_tfidf(corpus)
24     return tf_idf_dict, corpus
25     # key          value
26     # 文档序号      dict
27     #              key    value
28     #              词      TF·IDF值
29
30
31 ▼ def search_engine(query, tf_idf_dict, corpus, top=3):
32     # 先对要搜索的内容进行分词
33     query_words = jieba.lcut(query)
34     res = []
35     # 遍历所有文档，统计得分
36     ▼ for doc_id, tf_idf in tf_idf_dict.items():
37         score = 0
38         # 对要搜索的字符串中的每个词在某个文档内TF·IDF值的加和
39     ▼ for word in query_words:
40         score += tf_idf.get(word, 0)
41         # 该字符串在每个文档的TF·IDF值
42         res.append([doc_id, score])
43     # 对每个文档的TF·IDF值进行排序
44     res = sorted(res, reverse=True, key=lambda x: x[1])
```

```

45     for i in range(top):
46         doc_id = res[i][0]
47         print(corpus[doc_id])
48         print("-----")
49     return res
50
51
52 if __name__ == "__main__":
53     # 包含标题, 内容两个key
54     path = "news.json"
55     tf_idf_dict, corpus = load_data(path)
56     while True:
57         query = input("请输入您要搜索的内容:")
58         search_engine(query, tf_idf_dict, corpus)
59

```

• 文本摘要

- 1.通过计算TFIDF值得到每个文本的关键词。
- 2.将包含关键词多的句子, 认为是关键句。
- 3.挑选若干关键句, 作为文本的摘要。

```
1  import jieba
2  import math
3  import os
4  import random
5  import re
6  import json
7  from collections import defaultdict
8  from calculate_tfidf import calculate_tfidf, tf_idf_topk
9
10 """
11 基于tfidf实现简单文本摘要
12 """
13
14 jieba.initialize()
15
16
17 # 加载文档数据（可以想象成网页数据），计算每个网页的tfidf字典
18 ▼ def load_data(file_path):
19     corpus = []
20     with open(file_path, encoding="utf8") as f:
21         documents = json.loads(f.read())
22     ▼ for document in documents:
23         assert "\n" not in document["title"]
24         assert "\n" not in document["content"]
25         corpus.append(document["title"] + "\n" + document["content"])
26     tf_idf_dict = calculate_tfidf(corpus)
27     return tf_idf_dict, corpus
28
29
30 # 计算每一篇文章的摘要
31 # 输入该文章的tf_idf词典，和文章内容
32 # top为人为定义的选取的句子数量
33 # 过滤掉一些正文太短的文章，因为正文太短在做摘要意义不大
34 ▼ def generate_document_abstract(document_tf_idf, document, top=3):
35     sentences = re.split("? |! |。", document)
36     # 过滤掉正文在五句以内的文章
37     ▼ if len(sentences) <= 5:
38         return None
39     result = []
40     ▼ for index, sentence in enumerate(sentences):
41         sentence_score = 0
42         words = jieba.lcut(sentence)
43     ▼ for word in words:
44         # 统计word在该文档内的TF·IDF值
```

```

45         sentence_score += document_tf_idf.get(word, 0)
46         # 做一个归一化，否则越长的句子的得分会越高
47         sentence_score /= (len(words) + 1)
48         result.append([sentence_score, index])
49     result = sorted(result, key=lambda x: x[0], reverse=True)
50     # 权重最高的可能依次是第10，第6，第3句，将他们调整为出现顺序比较合理，即3,6,10
51     # 选得分最高的top个句子，根据index再排一次序
52     important_sentence_indexes = sorted([x[1] for x in result[:top]])
53     return "。".join([sentences[index] for index in important_sentence_indexes])
54
55
56 # 生成所有文章的摘要
57 def generate_abstract(tf_idf_dict, corpus):
58     res = []
59     for index, document_tf_idf in tf_idf_dict.items():
60         title, content = corpus[index].split("\n")
61         abstract = generate_document_abstract(document_tf_idf, content)
62         if abstract is None:
63             continue
64         corpus[index] += "\n" + abstract
65         res.append({"标题": title, "正文": content, "摘要": abstract})
66     return res
67
68
69 if __name__ == "__main__":
70     path = "news.json"
71     tf_idf_dict, corpus = load_data(path)
72     res = generate_abstract(tf_idf_dict, corpus)
73     writer = open("abstract.json", "w", encoding="utf8")
74     # ensure_ascii 设置为 False，以便生成更易读的 JSON 字符串，不用强行转换为非ascii字符
75     writer.write(json.dumps(res, ensure_ascii=False, indent=2))
76     writer.close()
77

```

• 文本相似度

- 对所有文本计算tfidf后，从每个文本选取tfidf较高的前n个词，得到一个词的集合S。
- 对于每篇文本D，计算S中的每个词的词频，将其作为文本的向量。
- 通过计算向量夹角余弦值，得到向量相似度，作为文本的相似度

- 向量夹角余弦值计算：

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

- 代码演示*


```

1  #coding:utf8
2  import jieba
3  import math
4  import os
5  import json
6  from collections import defaultdict
7  from calculate_tfidf import calculate_tfidf, tf_idf_topk
8
9  """
10  基于tfidf实现文本相似度计算
11  """
12
13  jieba.initialize()
14
15  #加载文档数据（可以想象成网页数据），计算每个网页的tfidf字典
16  #之后统计每篇文档重要在前10的词，统计出重要词词表
17  #重要词词表用于后续文本向量化
18  def load_data(file_path):
19      corpus = []
20      with open(file_path, encoding="utf8") as f:
21          documents = json.loads(f.read())
22          for document in documents:
23              corpus.append(document["title"] + "\n" + document["content"])
24      tf_idf_dict = calculate_tfidf(corpus)
25      # 每个文取前5个词构成词表
26      topk_words = tf_idf_topk(tf_idf_dict, top=5, print_word=False)
27      vocab = set()
28      for words in topk_words.values():
29          for word, _ in words:
30              vocab.add(word)
31      print("词表大小: ", len(vocab))
32      return tf_idf_dict, list(vocab), corpus
33
34
35  #passage是文本字符串
36  #vocab是词列表
37  #向量化的方式：计算每个重要词在文档中的出现频率
38  def doc_to_vec(passage, vocab):
39      vector = [0] * len(vocab)
40      passage_words = jieba.lcut(passage)
41      for index, word in enumerate(vocab):
42          vector[index] = passage_words.count(word) / len(passage_words)
43      return vector
44

```

```

45 #先计算所有文档的向量
46 def calculate_corpus_vectors(corpus, vocab):
47     corpus_vectors = [doc_to_vec(c, vocab) for c in corpus]
48     return corpus_vectors
49
50 #计算向量余弦相似度
51 def cosine_similarity(vector1, vector2):
52     x_dot_y = sum([x*y for x, y in zip(vector1, vector2)])
53     sqrt_x = math.sqrt(sum([x ** 2 for x in vector1]))
54     sqrt_y = math.sqrt(sum([x ** 2 for x in vector2]))
55     if sqrt_y == 0 or sqrt_y == 0:
56         return 0
57     # 分母加上很小的数, 防止出现除0
58     return x_dot_y / (sqrt_x * sqrt_y + 1e-7)
59
60
61 #输入一篇文本, 寻找最相似文本
62 def search_most_similar_document(passage, corpus_vectors, vocab):
63     input_vec = doc_to_vec(passage, vocab)
64     result = []
65     for index, vector in enumerate(corpus_vectors):
66         score = cosine_similarity(input_vec, vector)
67         result.append([index, score])
68     result = sorted(result, reverse=True, key=lambda x:x[1])
69     return result[:4]
70
71
72 if __name__ == "__main__":
73     path = "news.json"
74     tf_idf_dict, vocab, corpus = load_data(path)
75     corpus_vectors = calculate_corpus_vectors(corpus, vocab)
76     passage = "魔兽争霸"
77     for corpus_index, score in search_most_similar_document(passage, corpus_vectors, vocab):
78         print("相似文章:\n", corpus[corpus_index].strip())
79         print("得分: ", score)
80         print("-----")
81

```

10. TF·IDF优势

- 1.可解释性好
 - 可以清晰地看到关键词
 - 即使预测结果出错，也很容易找到原因
- 2.计算速度快
 - 分词本身占耗时最多，其余为简单统计计算
- 3.对标注数据依赖小
 - 可以使用无标注语料完成一部分工作
- 4.可以与很多算法组合使用
 - 可以看做是词权重

11. TF·IDF劣势

- 1.受分词效果影响大
 - 2.词与词之间没有语义相似度
 - 3.没有语序信息（词袋模型）
 - 4.能力范围有限，无法完成复杂任务，如机器翻译和实体挖掘等
 - 5.样本不均衡会对结果有很大影响
 - 6.类内样本间分布不被考虑
2. 如“您好”和“你好”在TF·IDF中是两个完全不同的词
3. 没有序列关系，像Pooling一样
4. 在标注数据少的情况下，TF·IDF比较常用

12. 编程练习：

给定词表，给出所有可能的分词结果

- 感觉能用回溯

13. 代码分析

a. forward_segmentation_method1

前向分割第一种方式，通过最大词长度，回退搜索

b. forward_segmentation_method2

前向分割第二种方式，记录词是否是其他词的前缀或真词