

## 1. 建立Tokenizer类

```
1 import jieba
2 import pandas as pd
3 import re
4
5
6 class Tokenizer:
7     def __init__(self, chars, coding='c', PAD=0):
8         self.coding = coding
9         self.chars = chars
10        self.PAD = PAD
11        d = {'[PAD]':PAD}
12        code = 1
13        if coding == 'c':
14            for string in chars:
15                for char in string:
16                    if char not in d:
17                        d[char] = code
18                        code += 1
19        elif coding == 'w':
20            for string in chars:
21                string = [word for word in jieba.lcut(string)]
22                for word in string:
23                    if word not in d:
24                        d[word] = code
25                        code += 1
26        self.vocabulary = d
27
28    def tokenize(self, sentence):
29        if self.coding == 'c':
30            list_of_chars = [char for char in sentence]
31        elif self.coding == 'w':
32            list_of_chars = [word for word in jieba.lcut(sentence)]
33        return list_of_chars
34
35    def encode(self, list_of_chars):
36        return [self.vocabulary[i] for i in list_of_chars]
37
38    def trim(self, tokens, seq_len):
39        n = seq_len - len(tokens)
40        if n < 0: # 超出部分截断
41            tokens = tokens[:seq_len]
42        else: # 不足部分补足
43            for i in range(n):
44                tokens.append(self.PAD)
45        return tokens
46
47    def decode(self, tokens):
48        voc_reverse = {v:k for k,v in self.vocabulary.items()}
49        return ''.join(voc_reverse[i] for i in tokens)
50
51    def encode_all(self, seq_len):
52        return [self.trim(self.encode(self.tokenize(i)), seq_len) for i in self.chars]
```



### 3. BERT 与 Word2Vec

Word2Vec 虽然相比于 one-hot 解决了高维度和语义表示问题，但没能解决同名词不同义的问题。

BERT 使用双向 Transformer，利用单词的上下文信息来做特征提取，解决了同名不同义的问题。