1. 进行初始化

传递进来字符串列表后，此时根据模式（w 为词模式，c 为字模式）判断进行词典的初始化。同时要判断字符是否为汉字，如果不是汉字不记录进入词典

```python
class Tokenizer:

    def __init__(self,chars,coding='c',PAD=0) :
        self.chars=chars
        self.coding=coding
        self.PAD=PAD
        self.dict={}
        self.dict["PAD"]=PAD
        self.count=1
        if self.coding == 'c':
            for i in self.chars:
                for j in list(i):
                    if j not in self.dict.keys() and '\u4e00' <= j <= '\u9fff':
                        self.dict[j]=self.count
                        self.count += 1
        else:
            for i in self.chars:
                for j in jieba.lcut(i):
                    if j not in self.dict.keys() and '\u4e00' <= j <= '\u9fff':
                        self.dict[j]=self.count
                        self.count += 1
```

运行结果：
字模式：



词模式：

2. 2. tokenize(self, sentence) 输入一句话，返回分词(字）后的字符列表(list_of_chars)。

**根据不同的模式进行分词（分字），同时筛除不是汉字的字符**

```python
def tokenize(self,sentence):
    res=[]
    if self.coding == 'c':
        for i in list(sentence):
            if '\u4e00' <= i <= '\u9fff':
                res.append(i)
    else:
        for i in jieba.lcut(sentence):
            if '\u4e00' <= i <= '\u9fff':
                res.append(i)
    return res
```

运行结果：

字模式：



词模式：



3. encode(self, list_of_chars) 输入字符(字或者词）的字符列表，返回转换后的数字列表
(tokens)

**根据词典生成相应的数字列表**

```python
def encode(self,list_of_chars):
    res=[]
    for i in list_of_chars:
        res.append(self.dict[i])
    return res
```

运行结果：

```
chars=[]
with open(filepath,'r',encoding="utf-8") as f:
    for line in f.readlines():
        chars.append(line.strip())

t=Tokenizer(chars,'w')

list_of_chars=t.tokenize(chars[55])
tokens=t.encode(list_of_chars)
print(tokens)
```

```
Building prefix dict from the default dictionary ...
Loading model from cache C:\Users\86186\AppData\Local\Temp\jieba.cache
Loading model cost 0.866 seconds.
Prefix dict has been built successfully.
[9, 128, 5, 1087, 96, 1119, 35, 168, 77, 164, 1201, 1202, 891, 38, 433, 1203, 46, 134, 106, 9, 975, 5, 711, 404, 464, 38, 923, 22, 1204, 134, 106, 9, 51, 237, 5, 482, 9, 14,
 212, 44, 62, 488, 24, 792]
Backend TkAgg is interactive backend. Turning interactive mode on.
```

4. trim(self, tokens, seq_len) 输入数字列表 tokens，整理数字列表的长度。不足 seq_len 的部分用 PAD 补足，超过的部分截断。

```
def trim(self,tokens,seq_len):
    while(len(tokens) < seq_len):
        tokens.append(self.PAD)
    while(len(tokens) > seq_len):
        tokens.pop()
    return tokens
```

运行结果：

```
Building prefix dict from the default dictionary ...
Loading model from cache C:\Users\86186\AppData\Local\Temp\jieba.cache
Loading model cost 0.794 seconds.
Prefix dict has been built successfully.
[9, 128, 5, 1087, 96, 1119, 35, 168, 77, 164, 1201, 1202, 891, 38, 433, 1203, 46, 134, 106, 9, 975, 5, 711, 404, 464, 38, 923, 22, 1204, 134, 106, 9, 51, 237, 5, 482, 9, 14,
 212, 44, 62, 488, 24, 792, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0]
Backend TkAgg is interactive backend. Turning interactive mode on.
```

5. decode(self, tokens) 将模型输出的数字列表翻译回句子。如果有 PAD，输出'[PAD]'。
**注意到词典中的值刚好对应键的索引值，故生成键的列表，利用键值当作索引值**

```
def decode(self,tokens):
    res=[]
    ls=list(self.dict.keys())
    for i in tokens:
        res.append(ls[i])

    for j in res:
        print(j,end='')
```

运行结果：

Building prefix dict from the default dictionary ...
Loading model from cache C:\Users\86186\AppData\Local\Temp\jieba.cache
Loading model cost 0.768 seconds.
Prefix dict has been built successfully.
很好的一次购物下单后秒发货由于下班时间晚晚上了才开始装客服都很配合的指导安装问了很多问题售后客服都很有耐心的系统很不错反应快是正版可以更新PADPADPADPADPADPADPADPADPADPADPA
DPADPADPADPADPADPADPADPADPADPADPADPADPADPADPADPADPADPADPADPADPADPADPADPADPADPADPADPADPADPADBackend TkAgg is interactive backend. Turning interactive mode on.

6. encode_all(self，seq_len) 返回所有文本（chars)的长度为 seq_len 的 tokens。

```python
def encode_all(self,seq_len):
    res=[]
    for i in self.chars:
        if len(i) == seq_len:
            list_of_chars=self.tokenize(i)
            tokens=self.encode(list_of_chars)
            res.append(self.trim(tokens,len(i)))
    return res
```

运行结果：

```python
chars=[]
with open(filepath,'r',encoding="utf-8") as f:
    for line in f.readlines():
        chars.append(line.strip())

t=Tokenizer(chars,'w')

list_of_chars=t.tokenize(chars[55])
tokens=t.encode(list_of_chars)
tokens=t.trim(tokens,len(chars[55]))
print(t.encode_all(55))
```

[[9, 1196, 1197, 1153, 114, 856, 25, 38, 1198, 296, 568, 569, 5, 1199, 106, 75, 46, 38, 225, 226, 5, 1200, 15, 1133, 5, 367, 32, 125, 368, 388, 0, 0, 0, 0, 0, 0, 0, 0, 0,
0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0], [9, 128, 5, 1496, 62, 406, 926, 14, 1021, 852, 433, 1497, 60, 9, 98, 1498, 1285, 272, 1499, 5, 1496, 128, 0, 0, 0, 0, 0, 0,
0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0], [1007, 821, 652, 38, 209, 1529, 211, 1530, 1531, 5, 173, 1532, 1158, 126, 38, 1533, 390, 346,
83, 24, 1534, 38, 192, 5, 1014, 1535, 1483, 211, 1536, 5, 346, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0], [405, 5, 346, 14, 1038, 32, 54, 1631
, 1692, 5, 446, 38, 376, 1693, 405, 5, 106, 1325, 854, 534, 132, 534, 1694, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0],
[9, 128, 9, 1515, 272, 192, 183, 5, 1695, 382, 128, 1696, 1697, 1698, 1699, 9, 1294, 177, 1365, 412, 68, 69, 1700, 947, 1701, 62, 1702, 5, 0, 0, 0, 0, 0, 0, 0, 0, 0
, 0, 0, 0, 0, 0, 0, 0, 0], [646, 573, 62, 60, 9, 98, 5, 189, 3312, 3313, 888, 1604, 2474, 225, 111, 3314, 9, 1099, 225, 3082, 1298, 2073, 190, 62, 3299,
3315, 763, 696, 763, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0], [3316, 78, 12, 496, 44, 1203, 68, 1625, 2764, 2599, 692, 76, 3317, 3318,
3319, 54, 1249, 38, 139, 205, 795, 696, 639, 316, 3320, 5, 181, 9, 2013, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0], [9, 128, 5, 36, 11,
3794, 1323, 367, 244, 1878, 1879, 68, 104, 458, 3795, 5, 887, 412, 81, 871, 382, 1884, 458, 422, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0
, 0, 0, 0, 0], [72, 43, 814, 75, 54, 652, 38, 139, 137, 4104, 77, 3863, 703, 88, 25, 121, 1146, 419, 9, 1196, 834, 1001, 9, 128, 25, 24, 1066, 4105, 591, 0, 0, 0, 0, 0, 0
, 0, 0, 0, 0, 0, 0], [62, 192, 9, 10, 5, 1087, 96, 78, 44, 76, 82, 128, 941, 652, 596, 91, 657, 14, 212, 4417, 5, 1460, 13, 44, 104
, 5, 4418, 1760, 84, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0], [25, 38, 1198, 433, 568, 569, 8, 14, 11, 75, 293, 289, 168, 173, 4922, 16
74, 380, 1505, 71, 1505, 421, 13, 802, 482, 83, 124, 667, 1055, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0], [83, 14, 60, 915, 5, 32, 12
5, 1297, 83, 944, 24, 190, 2866, 2867, 370, 170, 83, 14, 521, 2622, 1173, 38, 3822, 367, 388, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0
, 0, 0, 0], [611, 258, 2755, 2911, 1171, 29, 167, 905, 5, 1122, 5172, 5, 338, 675, 57, 54, 124, 178, 126, 38, 23, 650, 57, 167, 4883, 107, 5, 1122, 0, 0, 0, 0, 0, 0, 0
, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0], [612, 1289, 1429, 1271, 1427, 54, 1285, 9, 10, 38, 192, 484, 745, 126, 32, 612, 473, 1242, 83, 1378, 1378, 121, 211,
3304, 1115, 89, 1679, 921, 5406, 242, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0], [25, 38, 857, 170, 338, 128, 71, 2, 583, 81, 6246, 5, 795,
4280, 845, 6247, 6248, 6249, 1412, 1615, 188, 1976, 3050, 68, 6250, 136, 6251, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0], [78, 221
7, 38, 1144, 1119, 471, 1072, 54, 127, 38, 44, 4475, 3101, 128, 192, 471, 652, 54, 3151, 38, 6442, 5, 3463, 1099, 1220, 103, 124, 650, 126, 128, 0, 0, 0, 0, 0, 0, 0, 0
, 0, 0, 0], [170, 62, 6478, 5, 3500, 5, 103, 128, 740, 231, 9, 6479, 170, 211, 192, 1352, 5, 6480, 81, 5306, 5, 4152, 619, 23, 205, 1397,
24, 5029, 421, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0], [83, 14, 1338, 169, 192, 1929, 13, 9, 1294, 6530, 799, 6531, 38, 1307, 6532, 1
929, 5, 474, 4084, 757, 845, 6533, 2576, 650, 5, 78, 103, 1117, 10, 5, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0]]
[102 535  28 543 544  51  52  53  54  55  56  57  58  59  60  61  62  63
  64  65  66  67  68  69  70  71  72  73  74  75  76  77  78  79  80  81
  82  83  84  85  86  87  88  89  90  91  92  93  94  95  96  97  98  99
 609 101 100 103 104 105 106 107 108 109 110 111 112 113 114 115 116 117
 118 119 120 121 122 123 124 125 126 127 128 129 130 131 132 133 134 135
 136 138 139 140 141 143 144 145 146 147 148 150 151 153 154 155 159 160
 161 162 163 165 166 167 168 169 170 171 172 173 174 175 176 177 178 181
 182 183 184 185 186 187 191 192 193 195 196 197 198 199 200 201 203 204
 205 207 208 209 210 211 216 219 229 231 233 234 237 241 242 248 254 256
 257 258 263 264 271 274 281 284 287 289 290 294 297 298 304 313 317 332
 338 344 351 355 371 375 377 379 398 432 435 469 472 493 498 500 502 503
 508] [ 2  1  1 19 29 18 15 23 18 15 23 18 25 17 20 18 21 18 13 15 17 18 15

7. seq_len 是句子的长度，实际任务中一般怎么来确定一个合适的长度，请以前次作业中的评论文本或微博文本为例，通过文本的长度分布来进行观察和讨论。

**在上次作业的基础上，在 plotnodes 模块中加入对文本长度分布的函数，同时在本次作业的环境变量中加入包的路径，导入对应的模块**

```python
def plotlen(len_list):
    x=list(set(len_list))
    y=[]
    for i in x:
        caly=len_list.count(i)
        y.append(caly)
    x=np.array(x)
    y=np.array(y)
    print(x,y)
    plt.bar(x,y,color='purple')
    plt.xlabel("文本长度")
    plt.ylabel("频数")
    plt.title("文本长度的分布")
    plt.grid(alpha=0.5,linestyle='-.')
    plt.show()
```
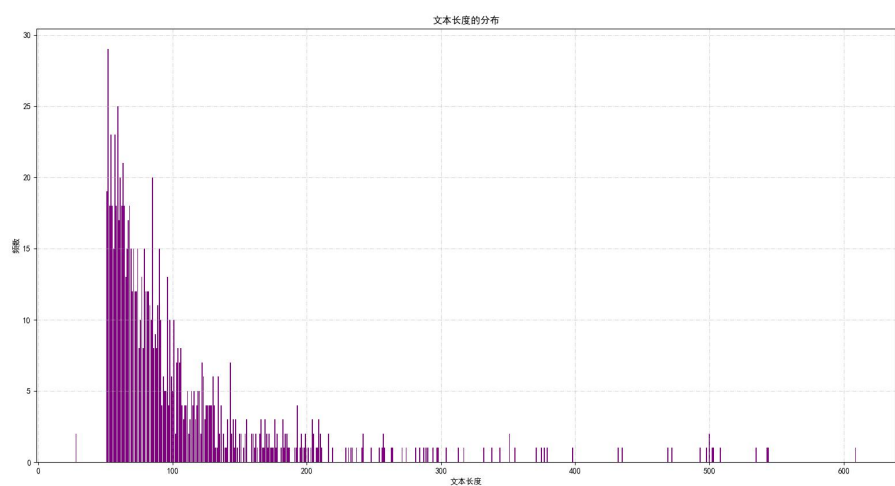
```python
import sys
sys.path.append(r"C:\Users\86186\Desktop\现代程序设计\4")

from GraphStat.Visualization import plotnodes
import jieba
```

运行结果:

```python
chars=[]
with open(filepath,'r',encoding="utf-8") as f:
    for line in f.readlines():
        chars.append(line.strip())
```

```python
len_list=[]
for i in chars:
    len_list.append(len(i))
plotnodes.plotlen(len_list)
```

文本长度的分布

将图放大可以发现，京东的评论中，长度为 52 的频数最多，为 29 次