型 $p(w_t|w_{t-1}) \rightarrow -\sum_{t=1}^{|w_t|} \log p(w_t|w_{t-1})$ 防止浮点数连乘 相乘之后向下溢出, 取对数

OOV : (out of vocabulary) 在字典外

$\begin{cases} pagerank & 有向图 \\ textrank & 无向图 \end{cases}$

文本生成也可以生成摘要 不过要进行语句检测。

textrank 也可以进行关键词提取。

文本聚类 (demo_text_clustering.py): 没有标签的聚类。K-means (K 均值聚类) ← 无监督

距离、相似性、迭代求解,随机总 K 个对象水机始的聚类中心。※ microsoft edge 可以画图.
↳ 腾讯会议 可以开批注

把点划分最匹配的那类。 从 K 中哪几点 根据簇间距离判断, 停止迭代。 目标函数

K-means : jianshu.com/p/fc81fed8c77b, 簇内尽量小, 簇间尽量大。

可计算文本相似度 文本聚类 ⇒ 一般用主题模型,

切词, 去停用词, 文档向量

代码. cnblogs.com/biaoding/p/8878080.html     doc2vec

sigmoid 分类, .demo_word2vec 词向量, word2vec 词向量, 一个词用一个向量表示 (13年谷歌)

跳字 (skip-gram) 和 词袋模型 (CBOW), 其中 skip-gram 用中心词预测周围词.

用周围词预测中心词.     softmax 多分类.   上文

丰度: 是指一种化学元素在其个自然体中重量占这个自然体总重量 的相对份额。

命名实体识别, 复用到 NES 标注集, 不构成命名实体的单词, 统一标注为 O (outside)
↓                                                    不是命名词汇外
现实存在的实体 ,

```
class ClusterAnalyzer():          def read_txt(self):      def cal_tfidf(self):
    def __init__(self):               pass                     pass
        pass    ⟨ self.document={}
                ⟨ self.document_tfidf={}
                                 def text_process(self):    def cal_dist(self
    def add_document(self, path):     pass                      arr1, arr2):
        self.document=self.read_text(path)                      pass
                                                           def cal_zone(self):
                                                               pass
```

```python
def save_text(self):
    pass

if __name__ == '__main__':
    data_path = os.path.join(data)
    kmeans = ClusterAanalyzer()
    kmeans.add_document(data_path)
    print(kmeans.document)
```

```python
def add_document(self, path):
    for i in range(100):
        file_path = os.path.join(path, '{}.txt'.format(i))
        with open(file_path, encoding='utf-8') as f:
            self.document[str(i)] = f.read()
        self.text_process()
```

文件名字要和模块名一样 否则会出错

os.path.dirname(__file__)

___file___ 上一层
有 debug 进程有 debug probe 可以用

init中再用 self.process_document={}, self.word_table=[]

segment = DoubleArrayTrieSegment()    stopwords = Stopwords()

```python
def text_process(self):
    for file_id, text in self.document.items():
        self.process_document[file_id] = [term.word for term in self.segment.seg(text)
                                          if not self.stopwords.containsKey(term.word)]
        word_list.extend
        [self.process_document[file_id]

    self.word_table = list(set(word_list))

    def fit(self, k):  #训练 kmeans
        self.cal_tfidf()    #要表示为固定的维度
```

```python
def cal_tfidf(self):
    for file_id, text in self.process_document.items():
        for word in text:
            tfidf_list = [0] * len(self.word_table)
```
同生成 这长的都0的张    词    成本
word, text    word, text
```python
            tfidf_list.insert(self.word_table.index(word), self.tf() * self.idf())
```
不是插入,是更新
```python
    self.document_itidf[file_id] = tfidf_list
```

```python
def tf( self ):                          # word
c   return text.count(word) /len(text)

def idf (self, word):
                                         # file_id, text
                                         #   ∨ process
    return np.log ( len(self.document)/sum([ 1 for self.document.items() if word in
                                                                              text]))
def fit (self, num=10):   # k=5     # 迭代次数
    self.cal_tfidf()                                 # 文件id
    self.cluster_dict=dict.fromkeys        {'1':[1,2,3], '2':[]}
              (range(1, k+1), {})                                    改为 {'center':[   ],'file_id:[   ]}
                  'center': None, 'file_id': None
                                           dict. fromkeys ( seq [, value])
                                                                 → 对应值
                                                          键
    # 随机初始化k个质中心:
    self. cluster_dict = dict.fromkeys (range (1, k+1),{'center': None, 'file_name': [None]})
                                                      self.process_document
    print (random. sample(range(1,len()), K ))
                                         # 遍历所有生成
    for i in (random. sample (range (1, len(self.process_document)),k):

        self. cluster_dict[str(index+1)]['center'] = self. document_tfidf['i']  → 改为str(i)
    count=1
    while count<num:    # 计算所有样本到质中心的距离

        for file_id, text in  self.document_tfidf.items():
          # 遍历所有质中心
          for __ in  self. cluster_dict.items():
              cluster_id, cluster
               text_dist.(self.dist (text, cluster['center']))
    min_dist_idx                 :append
          text_dist. index (min (text_dist))

          self. cluster_dict[min_dist_idx]['file_id']. append (file_id)
    count+=1
def cal_dist (self, arr1, arr2):                   用numpy更快, 底层用cpath,
                              np.array(arr1)
    return np. sqrt (np.sum((arr1 - arr2)**2))
                              np.array (arr2)
```

```python
for i in range (1, k+1):      # fromkeys: 对应创建了的值都是一排的
    self.cluster_dict[i] ={ "center": [], "file_id":[]} #前面的初始化,
                                                          有问题
更新聚类中心: self.update_cluster_center()
def update_cluster_center (self):
    for cluster_id, cluster in self.cluster_dict.items():
        for file_id in cluster['file_id']:    text_ifidf=[]
text_ifidf.append
text=>> self.document_itidf [file_id]] = tfidf_list
                         str
    arr_text_ifidf=np.array (text_ifidf)
    np.mean (arr_text_ifidf, axis =0).
    avg_cluster_center
print(avg_cluster_center.shape)  验证是2333维

self.cluster_dict[cluster_id] ['center']= avg_cluster_center,
更新完毕后, 在 while count<num 之后清空 self.cluster_dict
for i in range (1, k+1):
    self.cluster_dict [i] ['file_id']= []
def save_text (self):
    for cluster_id, cluster in self.cluster_dict.items ():
        if os. path. exists (str(cluster_id)):
            os. mkdir (str (cluster_id)):  #创建文件夹
        for file_id in cluster['cluster'] :                    (file_id))
            file_path = os. path. join ( str(cluster_id), 'f}. text'format
            with open (file_path,'w', encoding = 'utf-8') as f:
```

f.write (self.document [file_id])

主题模型 分布问题 (用的较多)

现在文本分类一般用深度学习,不用机器学习。

word2vec有上下文关系的,(在 pyhanlp下 static下 data中 word2vec 文件)

nlu : 自然语言理解 { 新词发现 (互信息,信息熵)

词法分析 ← { 分词分析 → 拼写纠正
          { 词性标注

句法分析 { 成分句法分析
         { 依存句法分析 (从从属词指向支配词)

语义分析 { 词义消歧 (elmo 解决一词多义)

          语义角色标注

          词汇/句子/段落的向量化表示

文本分类/情感分析                           应用 { 机器翻译
  文本匹配                                        { 搜索推荐 (工资高)
                                crf 实体识别       { 知识图谱
(难) { 信息抽取 { 实体抽取                         { 文本生成
              { 关系抽取                          { 问答系统
              { 事件抽取

阅读理解 (bert)

NLG : { 文本摘要
      { 机器翻译