

PYTHON自然语言处理（中文版）

一、语言处理与Python

- 1、NLTK入门
- 2、自然语言处理

二、获得文本语料和词汇资料

- 1、单语料库使用
 - 书籍
 - 网络文本
 - 即时聊天会话语料库
 - 布朗语料库
 - 路透社语料库
 - 就职演说语料库
 - 标注文本语料库: <https://www.nltk.org/howto/>
 - 词汇列表语料库
 - 其他语言语料库
- 2、使用自己的语料库
- 3、生成随机文本
- 4、条件概率
- 5、词典语料
- 6、WordNet

三、加工原料文本

- 1、从网络和硬盘访问文本
- 2、编码
- 3、正则表达式
- 4、词干提取器（中文不需要）
- 5、词形归并/词形还原（中文不需要）
- 6、分词

四、编写结构化程序

五、分类和标注词汇

- 默认标注器
- 正则表达式标注器
- 查询标注器
- N-gram标注器
- 组合标注器
- 存储标注器
- 基于转换的标注 - Brill标注
- 确定词性

六、学习分类文本

- 性别鉴定
- 词性分析
- 序列分类
- 句子分割

七、文法

- 使用文法
- 交互式文法编辑器
- 依存文法
- 特征结构

十一、语言数据管理

- TIMIT
- XML
- Toolbox
- OLAC元数据

PYTHON自然语言处理（中文版）

一、语言处理与Python

1、NLTK入门

下载数据集：输入命令：`nltk.download()`，之后选择要下载的数据集（如book）

若出现拒接连接则将https://github.com/nltk/nltk_data/tree/gh-pages/packages文件全部复制替换到C:\Users[用户名]\AppData\Roaming\nltk_data下即可解决

查看模块内容：`from nltk.book import *`

检索指定词：`text1.concordance("monstrous")`

检索指定词的上下文：`text2.similar("life")`

检索多个词共同的上下文：`text2.common_contexts(["love", "dear"])`

计算指定词从开头算起其前面有多少词：`text4.dispersion_plot(["citizens", "democracy", "freedom", "duties", "America"])`

随机生成指定文本风格的文字：`text3.generate()`

计算特定词出现的次数：`text3.count("smote")`

获取指定文本的词频信息（字典）：`fdist1=FreqDist(text1)`

检索只出现一次的词：`fdist1.hapaxes()`

前50高频词累计频率图（cumulative=False是频率分布图）：`fdist1.plot(50, cumulative=True)`

获取相邻词对：`bigrams(['more', 'is', 'than', 'done'])`

检索高频双连词：`text4.collocations()`

检索指定长度的词的数量：`fdist.items()`

计算词频：`fdist.freq('I')`

创建频率分布表：`fdist.tabulate()`

2、自然语言处理

文本对齐：给出一个双语文档，可以自动配对组成句子的过程

二、获得文本语料和词汇资料

1、单语料库使用

书籍

```

1 from nltk.corpus import gutenberg # 书籍文本库
2 gutenberg.fileids() # 获取语料库标识符
3
4 emma = gutenberg.words('austen-emma.txt') # 获取《爱玛》语料
5 emma = nltk.Text(emma) # 获取单个文本检索信息对象
6 emma.concordance('surprise')
7
8 gutenberg.raw('austen-emma.txt') # 原文
9 gutenberg.words('austen-emma.txt') # 词汇集
10 gutenberg.sents('austen-emma.txt') # 句子（词链表）

```

网络文本

```

1 from nltk.corpus import webtext # 网络文本集
2 for fileid in webtext.fileids():
3     print(fileid, webtext.raw(fileid)[:50])

```

即时聊天会话语料库

10-19-20s_706posts.xml表示包括2006年10月19日从20多岁聊天室收集的706个帖子。

```

1 from nltk.corpus import nps_chat # 即时聊天回话预料库
2 chatroom = nps_chat.posts('10-19-20s_706posts.xml')
3 chatroom[123]

```

布朗语料库

ID	文件	文体	描述
A16	ca16	新闻 news	Chicago Tribune: Society Reportage
B02	cb02	社论 editorial	Christian Science Monitor: Editorials
C17	cc17	评论 reviews	Time Magazine: Reviews
D12	cd12	宗教 religion	Underwood: Probing the Ethics of Realtors
E36	ce36	爱好 hobbies	Norling: Renting a Car in Europe
F25	cf25	传说 lore	Boroff: Jewish Teenage Culture
G22	cg22	纯文学 belles_lettres	Reiner: Coping with Runaway Technology
H15	ch15	政府 government	US Office of Civil and Defence Mobilization: The Family Fallout Shelter
J17	cj19	博览 learned	Mosteller: Probability with Statistical Applications
K04	ck04	小说 fiction	W.E.B. Du Bois: Worlds of Color
L13	cl13	推理小说 mystery	Hitchens: Footsteps in the Night
M01	cm01	科幻 science_fiction	Heinlein: Stranger in a Strange Land
N14	cn15	探险 adventure	Field: Rattlesnake Ridge
P12	cp12	言情 romance	Callaghan: A Passion in Rome
R06	cr06	幽默 humor	Thurber: The Future, If Any, of Comedy

```

1 from nltk.corpus import brown # 英语电子预料库
2 brown.categories()
3
4 brown.words(categories='news')
5
6 brown.words(fileids=['cg22'])
7
8 brown.sents(categories=['news', 'editorial', 'reviews'])

```

比较不同文体中情态动词的用法:

```

1 from nltk.corpus import brown
2
3 news_text = brown.words(categories='news')
4 fdist = nltk.FreqDist([w.lower() for w in news_text]) # 统计词频
5 modals = ['can', 'could', 'may', 'might', 'must', 'will']
6 # 显示指定文体中的情态动词词频信息
7 for m in modals:
8     print(m + ':' + str(fdist[m]))
9
10 # 获取不同文体以及对应的词
11 cfd = nltk.ConditionalFreqDist(
12     (genre, word) for genre in brown.categories()
13     for word in brown.words(categories=genre))
14
15 genres = ['news', 'religion', 'hobbies', 'science_fiction', 'romance',
16 'humor']
17 modals = ['can', 'could', 'may', 'might', 'must', 'will']
18 # 对指定文体的指定词计算词频
19 cfd.tabulate(conditions=genres, samples=modals)

```

路透社语料库

```

1 from nltk.corpus import reuters # 路透社预料库
2 reuters.fileids()[:10]
3
4 reuters.categories()[:20]
5
6 reuters.categories(['training/9865', 'test/14826']) # 获取指定集类别信息
7
8 reuters.fileids(['barley', 'corn'])[:10] # 指定类别查询语料集
9
10 reuters.words('training/9865')[:15] # 开头大写的是题目
11
12 reuters.words(['training/9865', 'training/9880'])[:15]
13
14 reuters.words(categories=['barley'])[:15]

```

就职演说语料库

```

1 from nltk.corpus import inaugural # 就职演说预料库
2 inaugural.fileids()[10]
3
4 [fileid[:4] for fileid in inaugural.fileids()][10] # 获取时间
5
6 # 绘制不同词在随时间演讲时的变换
7 cfd = nltk.ConditionalFreqDist(
8     (target, fileid[:4]) for fileid in inaugural.fileids()
9                             for w in inaugural.words(fileid)
10                                for target in ['america', 'citizen']
11                                    if w.lower().startswith(target))
12 cfd.plot()

```

标注文本语料库: <https://www.nltk.org/howto/>

词汇列表语料库

```

1 # 过滤高频词汇
2 def unusual_words(text):
3     text_vocab = set(w.lower() for w in text if w.isalpha()) # 不重复英文单词
4     english_vocab = set(w.lower() for w in nltk.corpus.words.words()) # 所有
    不重复单词
5     # 返回text_vocab不同于english_vocab的词
6     unusual = text_vocab.difference(english_vocab)
7     return sorted(unusual)
8
9 unusual_words(nltk.corpus.gutenberg.words('austen-sense.txt'))[10]

```

```

1 from nltk.corpus import stopwords # 停用词
2 stopwords.words('english')[10]

```

其他语言语料库

```

1 from nltk.corpus import udhr # 引入世界人权宣言语料
2
3 udhr.fileids()[20] # 检索语言
4
5 # 绘制不同语言在《世界人权宣言》的字长差异
6 languages = ['Chickasaw', 'English', 'German_Deutsch']
7 cfd = nltk.ConditionalFreqDist(
8     (lang, len(word)) for lang in languages
9                             for word in udhr.words(lang + '-Latin1'))
10 cfd.conditions() # 查看条件,对于每个cfd['xxx']都是一个频率分布
11 cfd.plot(cumulative=True)

```

2、使用自己的预料库

```

1 from nltk.corpus import PlaintextCorpusReader
2 corpus_root = '.' # 设置路径
3 wordlists = PlaintextCorpusReader(corpus_root, ".*") # 查询
4 wordlists.fileids() # 查看文件夹内容
5
6 wordlists.words('my_text.txt') # 使用

```

3、生成随机文本

```
1 def generate_model(cfdist, word, num=15): # 生成随机文本
2     for i in range(num):
3         print(word)
4         word = cfdist[word].max()
5
6 text = nltk.corpus.genesis.words('english-kjv.txt')
7 bigrams = nltk.bigrams(text) # 相邻词对
8 cfd = nltk.ConditionalFreqDist(bigrams)
9
10 cfd['living'] # 查看'living'后面可能跟哪些词
11
12 generate_model(cfd, 'living')
```

4、条件概率

```
1 pairs = [('a', 'b'), ('b', 'c'), ('c', 'a'), ('c', 'd')]
2 cfd = nltk.ConditionalFreqDist(pairs)
3 cfd # 有3个条件
4
5 cfd.conditions() # 将条件按字母排序
6
7 cfd['c'] # 查看此条件下的频率分布
8
9 cfd['c']['b'] # 查看指定条件和样本的频率
10
11 cfd.tabulate() # 为条件频率分布制表
12
13 cfd.plot() # 为条件频率分布绘图
14
15 # 获取不大于指定词频的单词，即得到词的词频都<=target的词频
16 target = 'abcde'
17 wordlist = nltk.corpus.words.words()
18 for w in wordlist:
19     if nltk.FreqDist(w) <= nltk.FreqDist(target):
20         print(w)
```

5、词典语料

```
1 names = nltk.corpus.names # 名字预料库
2 names.fileids()
3
4 male_names = names.words('male.txt')
5 male_names[:10]
6
7 # =====
8 entries = nltk.corpus.cmudict.entries() # 发音语料
9
10 for entry in entries[:10]:
11     print(entry) # 每个词后面是声音的标签，类似音节（同一个词可能多种发音）
12
13 # 找到所有发音开头与aalseth相似的词汇
14 syllable = 'AA1'
15 [word for word, pron in entries if syllable in pron][:10]
```

```

16
17 def stress(pron):
18     return [char for phone in pron for char in phone if char.isdigit()]
19
20 # 提取重音词，其中主重音（1）、次重音（2）和无重音（0）
21 [w for w, pron in entries if stress(pron) == ['0', '1', '2', '0']]
22
23 # 找到词汇的最小受限集合
24 p3 = [(pron[0] + '-' + pron[2], word) for (word, pron) in entries
25       if pron[0] == 'P' and len(pron) == 3]
26
27 cfd = nltk.ConditionalFreqDist(p3)
28 for template in cfd.conditions(): # 获取所有条件，如'P-TH'、'P-K'等
29     if len(cfd[template]) > 0: # 对应条件有元素
30         words = cfd[template]
31         wordlist = ' '.join(words)
32         print(template, wordlist[:70] + ' ...')
33
34 # 获取发音词典（若无则可以手动添加，但是对NLTK预料库不影响）
35 prondict = nltk.corpus.cmudict.dict()
36 prondict['fire']
37
38 # =====
39 from nltk.corpus import swadesh # 比较词表，多种语言都含200左右常用词
40 swadesh.fileids()
41
42 swadesh.words('en')[:10] # 英语
43
44 fr2en = swadesh.entries(['fr', 'en']) # 不同语言的同源词
45 fr2en[:10]
46
47 translate = dict(fr2en) # 翻译 fr -> en
48 translate['chien']
49
50 # 更新翻译词典
51 de2en = swadesh.entries(['de', 'en'])
52 translate.update(dict(de2en)) # 加入德语
53 translate['Hund']

```

6、WordNet

面向语义的英语词典，具有丰富的结构。下述的上位词和下位词都是同义词之间的空间关系。

同义词集：形容对某一事物的特征

同义词：相同特征的不同描述语

词条：同义词集和词的配对

```

1 from nltk.corpus import wordnet as wn
2
3 wn.synsets('motorcar') # 查找'motorcar'的同义词集
4 wn.synset('car.n.01').lemma_names() # 查看同义词
5 wn.synset('car.n.01').definition() # 定义
6 wn.synset('car.n.01').examples() # 例句
7 wn.synset('car.n.01').lemmas() # 同义词集的所有词条
8 wn.lemmas('car') # 查看包含'car'的所有词条
9 wn.lemma('car.n.01.automobile') # 查找特定的词条

```

```

10
11 wn.lemma('car.n.01.automobile').synset() # 得到词条的对应同义词集
12 wn.lemma('car.n.01.automobile').name() # 得到词条的"名字"
13
14 wn.synsets('car') # 与'motorcar'不同,有5个词集
15 for synset in wn.synsets('car'): # 查看每个同义词集的同义词
16     print(synset.lemma_names())

```

```

1 motorcar = wn.synset('car.n.01')
2 type_of_motorcar = motorcar.hyponyms() # 下位词(包含的词集)
3
4 motorcar.hypernyms() # 上位词(父词集)
5
6 paths = motorcar.hypernym_paths() # 上位词路径
7
8 motorcar.root_hypernyms() # 根上位词

```

```

1 wn.synset('tree.n.01').part_meronyms() # 包含的部分
2 wn.synset('tree.n.01').substance_meronyms() # 其实质
3 wn.synset('tree.n.01').member_holonyms() # 组成的整体
4
5 for synset in wn.synsets('mint', wn.NOUN):
6     print(synset.name(), ': ', synset.definition())
7 wn.synset('mint.n.04').part_holonyms() # 叶子是薄荷的一部分
8 wn.synset('mint.n.04').substance_holonyms() # 用薄荷油制作的糖材质是薄荷
9
10 wn.synset('walk.v.01').entailments() # 走路的"需求"包括抬脚
11 wn.synset('eat.v.01').entailments()
12 wn.lemma('supply.n.02.supply').antonyms() # "供给"的反义词是"需求"
13 dir(wn.synset('harmony.n.02')) # 查看指定词集的所有方法

```

```

1 right = wn.synset('right_whale.n.01')
2 orca = wn.synset('orca.n.01')
3 minke = wn.synset('minke_whale.n.01')
4 tortoise = wn.synset('tortoise.n.01')
5 novel = wn.synset('novel.n.01')
6 right.lowest_common_hypernyms(minke) # 求语义最相近的
7
8 wn.synset('baleen_whale.n.01').min_depth() # 同义词集深度
9 wn.synset('whale.n.02').min_depth()
10 wn.synset('entity.n.01').min_depth()
11
12 right.path_similarity(minke) # 求取相似度
13
14 help(wn)

```

三、加工原料文本

1、从网络和硬盘访问文本

访问本地文本可以直接使用 `open('xxx').read()` 来获取raw

```

1 from urllib.request import urlopen
2

```



```

3 url = "http://www.gutenberg.org/files/2554/2554-0.txt"
4 raw = urlopen(url).read()
5 raw = str(raw)
6
7 tokens = nltk.word_tokenize(raw) # 生成token列表
8 nltk.tokenwrap(raw) # 生成token字符串
9
10 text = nltk.Text(tokens)
11 text.collocations() # 检测高频双连词
12
13 raw.find("PART I")
14 raw.rfind("the subject of a new story")
15
16 raw = raw[5866: 1338204]
17 raw.find("PART I")
18
19 raw = raw[5866: 1338204]
20 raw.find("the subject of a new story")
21
22 '''
23 需要清除html情况
24 '''
25 url = "http://www.gutenberg.org/files/2554/2554-h/2554-h.htm"
26 html = urlopen(url).read()
27
28 from bs4 import BeautifulSoup
29
30 raw = BeautifulSoup(html).get_text()
31 tokens = nltk.word_tokenize(raw)
32
33 text = nltk.Text(tokens)
34 text.concordance('forbidden') # 检索指定词
35
36 '''
37 读取NLTK语料库文件
38 '''
39 path = nltk.data.find('corpora/gutenberg/melville-moby_dick.txt')
40 raw = open(path, 'r').read()

```

2、编码

解码：翻译为Unicode

编码：将Unicode转化为其他编码的过程

文件开头添加：# -*- coding: utf-8 -*-

3、正则表达式

^：表示以后面的字符为开头，若放在方括号里则代表除了括号内的字符之外

\$：表示以前面的字符为结尾

.：匹配任意单个字符

+: 匹配至少一个字符，也被成为闭包

*: 匹配至少零个字符，也被成为闭包

\: 表示后面的一个字符不具备特殊匹配含义

{a, b}: 表示前面的项目指定重复次数

字符串前面加'r'表示原始字符串

符号	功能
\b	词边界（零宽度）
\d	任一十进制数字（相当于[0-9]）
\D	任何非数字字符（等价于[^ 0-9]）
\s	任何空白字符（相当于[\t\n\r\f\v]）
\S	任何非空白字符（相当于[^ \t\n\r\f\v]）
\w	任何字母数字字符（相当于[a-zA-Z0-9_]）
\W	任何非字母数字字符（相当于[^a-zA-Z0-9_]）
\t	制表符
\n	换行符

```
1 import re
2
3 wordlist = [w for w in nltk.corpus.words.words() if w.islower()]
4 [w for w in wordlist if re.search('ed$', w)][:10] # 查找以ed结尾的词
5
6 [w for w in wordlist if re.search('^.j..t..$', w)][:10] # 查找指定位置的词
7
8 [w for w in wordlist if re.search('^a(b|c)k', w)][:10]
9
10 word = 'kfbhfuicbhflchbirlblirb'
11 re.findall(r'[aeiou]', word)
12
13 [int(n) for n in re.findall(r'[0-9]{2,4}', '2019-12-31')]
```

```
1 # 从罗托卡特语词汇中提取所有辅音-元音序列
2 cv_word_pairs = [(cv, w) for w in nltk.corpus.toolbox.words('rotokas.dic')
3                   for cv in re.findall(r'[ptksvr][aeiou]', w)]
4 cv_index = nltk.Index(cv_word_pairs) # 转换为索引表
5
6 cv_index['po']
```

搜索已分词文本

```
1 from nltk.corpus import gutenberg, nps_chat
2 moby = nltk.Text(gutenberg.words('melville-moby_dick.txt'))
3 # 尖括号用于标记标识符的边界,
4 # 尖括号之间的所有空白都被忽略（这只对nltk中的findall()方法有效）
5 moby.findall(r'<a><man>')
6 moby.findall(r'<a><.*><man>')
7 moby.findall(r'<a><(<.*>)<man>')
8
9 w = nps_chat.words()
10 # 找出以'bro'结尾的三个词组成的短语
11 chat = nltk.Text(w)
12 chat.findall(r'<.*><.*><bro>')
13 # 找出以字母'l'开始的三个或更多词组成的序列
14 chat.findall(r'<l.*>{3,}')
```

```

1 # 标记字符串指定模式
2 nltk.re_show(r'^a|b', 'aaasfgsgdhabbbdsdg')
3
4 # 提供正则匹配的图形化程序
5 nltk.app.nemo()

```

4、词干提取器（中文不需要）

词干：词的原形，不包含后缀如ly,s,es等

```

1 raw = '''
2 a word or unit of text to be carried over to a new line automatically as
3 the margin is reached, or to fit around embedded features such as pictures.
4 '''
5 tokens = nltk.word_tokenize(raw)
6
7 # 两种不同的内置词干提取器
8 porter = nltk.PorterStemmer()
9 lancaster = nltk.LancasterStemmer()
10
11 [porter.stem(t) for t in tokens][:10]
12 [lancaster.stem(t) for t in tokens][:10]

```

```

1 # 使用词干提取器索引文本
2 class IndexedText(object):
3     def __init__(self, stemmer, text):
4         self._text = text
5         self._stemmer = stemmer
6         self._index = nltk.Index((self._stem(word), i)
7                                   for (i, word) in enumerate(text))
8
9     def concordance(self, word, width=40):
10        key = self._stem(word)
11        wc = width // 4 # 前后各多少词
12        for i in self._index[key]: # 目的是对齐输出显示
13            lcontext = ' '.join(self._text[i-wc:i])
14            rcontext = ' '.join(self._text[i:i+wc])
15            ldisplay = '%*s' % (width, lcontext[-width:])
16            rdisplay = '%-*s' % (width, rcontext[:width])
17            print(ldisplay, rdisplay)
18
19        def _stem(self, word):
20            return self._stemmer.stem(word).lower()
21
22
23 poter = nltk.PorterStemmer()
24 grail = nltk.corpus.webtext.words('grail.txt')
25 text = IndexedText(poter, grail)
26 text.concordance('lie')

```

5、词形归并/词形还原（中文不需要）

词形还原：将复杂的词形表示转换为最基本的状态，如women装换为woman，但是不转换普通的过去式等

```

1 # wordNet词形归并器删除词缀产生的词都是在它的字典中的词
2 raw = '''
3 a women or unit of text to be carried lying to a new line automatically as
4 the margin is reached, or to fit around embedded features such as pictures.
5 '''
6 tokens = nltk.word_tokenize(raw)
7 wn1 = nltk.WordNetLemmatizer()
8 [wn1.lemmatize(t) for t in tokens][:10]

```

6、分词

模拟退火算法：迭代求解策略的一种随机寻优算法。模拟退火算法从某一较高初温出发，伴随温度参数的不断下降，结合概率突跳特性在解空间中随机寻找目标函数的全局最优解，即在局部最优解能概率性地跳出并最终趋于全局最优。模拟退火算法是一种通用的优化算法，理论上算法具有概率的全局优化性能

```

1 text = 'doyouseethekittyseethedoggydoyoulikethekittylikethedoggy'
2 seg1 = '00000000000000010000000000100000000000000010000000000'
3 seg2 = '01001001001000010010010000101001000100100001000100100'
4
5 # 根据segs对text进行分段/词
6 def segment(text, segs):
7     words = []
8     last = 0
9     for i in range(len(segs)):
10         if segs[i] == '1':
11             words.append(text[last: i+1])
12             last = i + 1
13     words.append(text[last:])
14     return words
15
16 # 通过合计每个词项与推导表的字符数，作为分词质量的得分。值越小越好
17 def evaluate(text, segs):
18     words = segment(text, segs)
19     text_size = len(words)
20     lexicon_size = len(' '.join(list(set(words)))) # 查找不重复的所有词
21     return text_size + lexicon_size
22
23 # 使用模拟退火算法的非确定形搜索
24 # 一开始仅搜索短语分词：随机扰动1/0；
25 # 它们与“温度”成正比，每次迭代温度都会降低，扰动边界会减少
26 from random import randint
27
28 def flip(segs, pos): # 更改pos位置的数字：0到1或1到0
29     return segs[:pos] + str(1-int(segs[pos])) + segs[pos+1:]
30
31 def flip_n(segs, n): # 对字符串随机更改n位
32     for i in range(n):
33         segs = flip(segs, randint(0, len(segs) - 1))
34     return segs
35
36 def anneal(text, segs, iterations, cooling_rate):
37     temperature = float(len(segs))
38     while temperature > 0.5:
39         best_segs, best = segs, evaluate(text, segs)

```

```

40         for i in range(iterations):
41             guess = flip_n(segs, int(round(temperature))) # 产生一个随机更改n
位的串
42             score = evaluate(text, guess)
43             if score < best:
44                 best, best_segs = score, guess
45             score, segs = best, best_segs
46             temperature = temperature / cooling_rate # 更改位数递减
47             print(evaluate(text, segs), segment(text, segs))
48         print()
49         return segs
50
51     anneal(text, seg1, 5000, 1.2)

```

四、编写结构化程序

```

1  # 绘制条形图来统计词频
2  import nltk
3
4  colors = 'rbcmky' # 多种颜色的组合, 字符串形式
5  def bar_chart(categories, words, counts):
6      import pylab
7      ind = pylab.arange(len(words))
8      width = 1 / (len(categories))
9      bar_groups = []
10     for c in range(len(categories)):
11         bars = pylab.bar(ind + c * width, counts[categories[c]],
12                          width, color=colors[c % len(colors)])
13         bar_groups.append(bars)
14     pylab.xticks(ind + width, words)
15     pylab.legend([b[0] for b in bar_groups], categories, loc='upper left')
16     pylab.ylabel('Frequency')
17     pylab.title('Frequency of Six Modal Verbs by Genre')
18     pylab.show()
19
20     genres = ['news', 'religion', 'hobbies', 'government', 'adventure']
21     modals = ['can', 'could', 'may', 'might', 'must', 'will']
22     cfdist = nltk.ConditionalFreqDist((genre, word)
23                                       for genre in genres
24                                       for word in
nltk.corpus.brown.words(categories=genre)
25                                       if word in modals)
26     counts = {}
27     for genre in genres:
28         counts[genre] = [cfdist[genre][word] for word in modals]
29
30     bar_chart(genres, modals, counts)

```

```

1  # 绘制树形图构建上下文关系
2  import networkx as nx
3  import matplotlib
4  from nltk.corpus import wordnet as wn
5
6  def traverse(graph, start, node):
7      graph.depth[node.name] = node.shortest_path_distance(start)

```



```

18 tag_fd = nltk.FreqDist(tag for (word, tag) in brown_news_tagged)
19 tag_fd
20 tag_fd.plot(cumulative=False)

```

```

1 # 找出最频繁的名词
2 def findtags(tag_prefix, tagged_text):
3     cfd = nltk.ConditionalFreqDist((tag, word) for (word, tag) in tagged_text
4                                     if tag.startswith(tag_prefix))
5     return dict((tag, list(cfd[tag].keys())[:5]) for tag in
6                 cfd.conditions())
7 tagdict = findtags('NN', nltk.corpus.brown.tagged_words(categories='news'))
8 for tag in sorted(tagdict):
9     print(tag, tagdict[tag])

```

默认标注器

```

1 brown_tagged_sents = brown.tagged_sents(categories='news')
2 brown_sents = brown.sents(categories='news')
3
4 # 查看最多词性是什么，指定为默认
5 tags = [tag for (word, tag) in brown.tagged_words(categories='news')]
6 nltk.FreqDist(tags).max()
7
8 # 使用
9 raw = 'i do not like green eggs and ham, i do not like them sam i am!'
10 tokens = nltk.word_tokenize(raw)
11 default_tagger = nltk.DefaultTagger('NN')
12 default_tagger.tag(tokens)
13
14 # 评估
15 default_tagger.evaluate(brown_tagged_sents)

```

正则表达式标注器

```

1 patterns = [
2     (r'.*ing$', 'VBG'),
3     (r'.*ed$', 'VBD'),
4     (r'.*es$', 'VBZ'),
5     (r'.*ould$', 'MD'),
6     (r'.*\'s$', 'NN$'),
7     (r'.*s$', 'NNS'),
8     (r'^-?[0-9]+(.[0-9]+)?$', 'CD'),
9     (r'.*', 'NN')
10 ]
11
12 regexp_tagger = nltk.RegexpTagger(patterns)
13 regexp_tagger.tag(brown_sents[3])
14
15 regexp_tagger.evaluate(brown_tagged_sents)

```

查询标注器

```

1 fd = nltk.FreqDist(brown.words(categories='news')) # 统计词频

```

```

2  cfd = nltk.ConditionalFreqDist(brown.tagged_words(categories='news')) # 词
   - 词性 统计
3
4  most_freq_words = list(fd.keys())[:100] # 最高频的前100词
5  likely_tags = dict((word, cfd[word].max()) for word in most_freq_words) #
   高频词及高频词性
6
7  baseline_tagger = nltk.UnigramTagger(model=likely_tags)
8
9  baseline_tagger.evaluate(brown_tagged_sents)
10
11 sent = brown.sents(categories='news')[3]
12 baseline_tagger.tag(sent) # 不在查询表中时会被标记为空
13 # 从上可知有些词被标注为空，因为我们需要将未查到的词性标注为默认的
14 baseline_tagger = nltk.UnigramTagger(model=likely_tags,
   backoff=nltk.DefaultTagger('NN'))
15
16 def performance(cfd, wordlist):
17     lt = dict((word, cfd[word].max()) for word in wordlist)
18     baseline_tagger = nltk.UnigramTagger(model=lt,
19                                           backoff=nltk.DefaultTagger('NN'))
20     return baseline_tagger.evaluate(brown.tagged_sents(categories='news'))
21
22 def display():
23     import pylab
24     words_by_freq = list(nltk.FreqDist(brown.words(categories='news')))
25     cfd = nltk.ConditionalFreqDist(brown.tagged_words(categories='news'))
26     sizes = 2 ** pylab.arange(15)
27     perfs = [performance(cfd, words_by_freq[:size]) for size in sizes]
28     pylab.plot(sizes, perfs, '-bo')
29     pylab.title('Lookup Tagger Performance with Varying Model Size')
30     pylab.xlabel('Model Size')
31     pylab.ylabel('Performance')
32     pylab.show()
33
34 display()

```

N-gram标注器

挑选在给定上下文中最可能的标记

上下文：当前词和它前面的n-1个标识符的词性标记

问题：当标注的数据在训练数据中不存在且占比较高时，被称为数据稀疏问题，在NLP中是普遍的。因此研究结果的精度和覆盖范围之间需要有一个权衡

```

1  bigram_tagger = nltk.BigramTagger(train_sents) # 二元标注（未出现的不会统计）
2  bigram_tagger.tag(brown_sents[2007])
3
4  unseen_sent = brown_sents[4203]
5  bigram_tagger.tag(unseen_sent)
6
7  bigram_tagger.evaluate(test_sents)

```

组合标注器

如首先使用二元标注器，否则使用一元标注器，最后使用默认标注器


```

1 t0 = nltk.DefaultTagger('NN')
2 t1 = nltk.UnigramTagger(train_sents, backoff=t0)
3 t2 = nltk.BigramTagger(train_sents, backoff=t1)
4 t2.evaluate(test_sents)
5
6 # 将会丢弃那些只看到一次或两次的上下文
7 t2 = nltk.BigramTagger(train_sents, cutoff=2, backoff=t1)
8 t2.evaluate(test_sents)

```

存储标注器

```

1 from pickle import dump
2 from pickle import load
3
4 # 存储
5 output = open('t2.pkl', 'wb')
6 dump(t2, output, -1)
7 output.close()
8
9 # 取出
10 inputs = open('t2.pkl', 'rb')
11 tagger = load(inputs)
12 inputs.close()
13
14 # 使用
15 text = '''
16 One notable effort in increasing the interoperability of biomedical
17 ontologies has been the creation of logical definitions[71]. This is an
18 initiative
19 '''
20 tokens = text.split()
21 tagger.tag(tokens)

```

给定当前单词及其前两个标记，根据训练数据，在5%的情况下，有一个以上的标记可能合理的分配给当前词。假设我们总是挑选在这种含糊不清的上下文中最可能的标记，可以得出trigram标注器性能的下界

训练数据中的歧义导致标注器性能的上限

另一种查看标注错误的方法：查看混淆矩阵

基于转换的标注 - Brill标注

思想：猜每个词的标记，然后返回和修复错误。这种方式下，Brill标注器陆续将一个不良标注的文本转换成一个更好的。同样是有监督的学习方法，不过和n-gram不同的是，它不计数观察结果，只编制一个转换修正规则的链表。

以画画类比：以大笔画开始，然后修复细节，一点点的修改

我们将研究两个规则的运作：(a) 当前面的词是 TO 时，替换 NN 为 VB；(b) 当下一个标记是 NNS 时，替换 TO 为 IN。表 5-6 说明了这一过程，首先使用 unigram 标注器标注，然后运用规则修正错误。

表 5-6. Brill 标注的步骤

Phrase	to	increase	grants	to	states	for	vocational	rehabilitation
Unigram	TO	NN	NNS	TO	NNS	IN	JJ	NN
Rule1		VB						
Rule2				IN				
Output	TO	VB	NNS	IN	NNS	IN	JJ	NN
Gold	TO	VB	NNS	IN	NNS	IN	JJ	NN

确定词性

根据词的形式（如-ness、-ing等）、根据句法（词的顺序）以及语义等综合分析。没有正确的方式来分配标记，只有根据目标不同或多或少有用的方法

六、学习分类文本

性别鉴定

```

1 import nltk
2 import random
3 from nltk.corpus import names
4
5 # 提取特征
6 def gender_features(word):
7     return {'last_letter': word[-1]}
8
9 name_set = [(name, 'male') for name in names.words('male.txt')] +
10             [(name, 'female') for name in names.words('female.txt')]
11 random.shuffle(name_set)
12
13 featuresets = [(gender_features(n), g) for (n, g) in name_set]
14 train_set, test_set = featuresets[500:], featuresets[:500]
15
16 classifier = nltk.NaiveBayesClassifier.train(train_set)
17
18 # 查看模型效果
19 classifier.classify(gender_features('Neo'))
20 classifier.classify(gender_features('Trinity'))
21
22 nltk.classify.accuracy(classifier, test_set)
23
24 # 显示的比率为似然比，用于比较不同特征-结果关系
25 classifier.show_most_informative_features(5)
26
27 # 在处理大型预料库时，构建一个包含每一个实例的特征的单独的链表会使用大量的内存。
28 # 下述方式不会在内容中存储所有的特征集对象
29 from nltk.classify import apply_features
30 train_set = apply_features(gender_features, name_set[500:])
31 test_set = apply_features(gender_features, name_set[:500])

```

词性分析

```

1  from nltk.corpus import brown
2
3  # 找出最常见的后缀
4  suffix_fdist = nltk.FreqDist()
5  for word in brown.words():
6      word = word.lower()
7      suffix_fdist[word[-1:]] += 1
8      suffix_fdist[word[-2:]] += 1
9      suffix_fdist[word[-3:]] += 1
10
11 common_suffixes = list(suffix_fdist.keys())[:100]
12 common_suffixes
13
14 # 词性特征提取器
15 def pos_features(word):
16     features = {}
17     for suffix in common_suffixes:
18         features['endwith(%s)' % suffix] = word.lower().endwith(suffix)
19     return features
20
21 tagged_words = brown.tagged_words(categories='news')
22 featuresets = [(pos_features(n), g) for (n, g) in tagged_words]
23
24 size = int(len(featuresets) * 0.1)
25 train_set, test_set = featuresets[size:], featuresets[:size]
26
27 classifier = nltk.DecisionTreeClassifier.train(train_set)
28 nltk.classify.accuracy(classifier, test_set)
29
30 classifier.classify(pos_features('cats'))

```

上下文语境

```

1  # 基于句子的词特征提取
2  # 输入分别为: 句链表, 句中词的索引
3  # 输出为特征
4  def pos_features(sentence, i):
5      features = {'suffix(1)': sentence[i][-1:],
6                  'suffix(2)': sentence[i][-2:],
7                  'suffix(3)': sentence[i][-3:]}
8      if i == 0:
9          features['prev-word'] = '<START>'
10     else:
11         features['prev-word'] = sentence[i - 1]
12     return features
13
14 brown.sents()[0]
15 pos_features(brown.sents()[0], 8)
16
17 tagged_sents = brown.tagged_sents(categories='news')
18 featuresets = []
19 for tagged_sent in tagged_sents:
20     untagged_sent = nltk.tag.untag(tagged_sent)
21     for i, (word, tag) in enumerate(untagged_sent):
22         featuresets.append((pos_features(untagged_sent, i), tag))
23
24 size = int(len(featuresets) * 0.1)

```

```

25 train_set, test_set = featuresets[size:], featuresets[:size]
26
27 classifier = nltk.NaiveBayesClassifier.train(train_set)
28 nltk.classify.accuracy(classifier, test_set)

```

序列分类

```

1 def pos_features(sentence, i, history):
2     features = {'suffix(1)': sentence[i][-1:],
3                 'suffix(2)': sentence[i][-2:],
4                 'suffix(3)': sentence[i][-3:]}
5     if i == 0:
6         features['prev-word'] = '<START>'
7         features['prev_tag'] = '<START>'
8     else:
9         features['prev-word'] = sentence[i - 1]
10        features['prev_tag'] = history[i - 1]
11    return features
12
13 # 序列分类器
14 class ConsecutivePosTagger(nltk.TaggerI):
15     def __init__(self, train_sents):
16         train_set = []
17         for tagged_sent in train_sents:
18             untagged_sent = nltk.tag.untag(tagged_sent)
19             history = []
20             for i, (word, tag) in enumerate(untagged_sent):
21                 featureset = pos_features(untagged_sent, i, history)
22                 train_set.append((featureset, tag))
23                 history.append(tag)
24             self.classifier = nltk.NaiveBayesClassifier.train(train_set)
25
26     def tag(self, sentence):
27         history = []
28         for i, word in enumerate(sentence):
29             featureset = pos_features(sentence, i, history)
30             tag = self.classifier.classify(featureset)
31             history.append(tag)
32         return zip(sentence, history)
33
34 tagged_sents = brown.tagged_sents(categories='news')
35
36 size = int(len(tagged_sents) * 0.1)
37 train_sents, test_sents = tagged_sents[size:], tagged_sents[:size]
38
39 tagger = ConsecutivePosTagger(train_sents)
40 tagger.evaluate(test_sents)

```

隐马尔可夫模型：是统计模型，它用来描述一个含有隐含未知参数的马尔可夫过程。其难点是从可观察的参数中确定该过程的隐含参数。然后利用这些参数来作进一步的分析，例如模式识别。

在正常的马尔可夫模型中，状态对于观察者来说是直接可见的。这样状态的转换概率便是全部的参数。而在隐马尔可夫模型中，状态并不是直接可见的，但受状态影响的某些变量则是可见的。每一个状态在可能输出的符号上都有一概率分布。因此输出符号的序列能够透露出状态序列的一些信息。

实例：假设你有一个住得很远的朋友，他每天跟你打电话告诉你他那天做了什么。你的朋友仅仅对三种活动感兴趣：公园散步，购物以及清理房间。他选择做什么事情只凭天气。你对于他所住的地方的天气情况并不了解，但是你知道总的趋势。在他告诉你每天所做的事情基础上，你想要猜测他所在地的天气情况。

你认为天气的运行就像一个马尔可夫链。其有两个状态 "雨"和"晴",但是无法直接观察它们,也就是说,它们对于你是隐藏的。每天，你的朋友有一定的概率进行下列活动:"散步"、"购物"、"清理"。因为你朋友告诉你他的活动，所以这些活动就是你的观察数据。这整个系统就是一个隐马尔可夫模型HMM。

在这个例子中，如果今天下雨,那么明天天晴的概率只有30%，表示了你朋友每天做某件事的概率。如果下雨，有 50% 的概率他在清理房间；如果天晴，则有60%的概率他在外头散步。

应用：语音识别、中文分词、光学字符识别和机器翻译等

句子分割

```
1 sents = nltk.corpus.treebank_raw.sents()
2 tokens = [] # 存储句子
3 boundaries = set() # 对应句子的词索引
4 offset = 0 # 总词数
5 for sent in nltk.corpus.treebank_raw.sents():
6     tokens.extend(sent)
7     offset += len(sent)
8     boundaries.add(offset - 1)
9
10 def punct_features(tokens, i):
11     return {'next-word-capitalized': tokens[i + 1][0].isupper(),
12            'prevword': tokens[i - 1].lower(),
13            'punct': tokens[i],
14            'prev-word-is-one-char': len(tokens[i - 1]) == 1}
15
16 # 提取可能是句子结束符的特征
17 # 特征：句索引链表
18 featuresets = [(punct_features(tokens, i), (i in boundaries))
19                for i in range(1, len(tokens) - 1)
20                if tokens[i] in '?!')]
21
22 size = int(len(featuresets) * 0.1)
23 train_set, test_set = featuresets[size:], featuresets[:size]
24
25 classifier = nltk.NaiveBayesClassifier.train(train_set)
26 nltk.classify.accuracy(classifier, test_set)
```

七-十、文法

使用文法

```
1 import nltk
2 from nltk import CFG
3
4 groucho_grammar = CFG.fromstring("""
5 S -> NP VP
6 PP -> P NP
7 NP -> Det N | Det N PP | 'I'
```

```

8  VP -> V NP | VP PP
9  Det -> 'an' | 'my'
10 N -> 'elephant' | 'pajamas'
11 V -> 'shot'
12 P -> 'in'
13 """
14
15 sent = ['I', 'shot', 'an', 'elephant', 'in', 'my', 'pajamas']
16 # 使用递归下降分析器
17 parser = nltk.ChartParser(groucho_grammar)
18 trees = parser.parse(sent)
19 for tree in trees:
20     print(tree)

```

```

1  # cfg文件内容与上块5-12行代码类似
2  grammar = nltk.data.load('file:mygrammar.cfg')
3
4  sent = 'Mary saw Bob'.split()
5  rd_parser = nltk.RecursiveDescentParser(grammar)
6  for tree in rd_parser.parse(sent):
7      print(tree)

```

```

1  # 使用移进-规约分析器
2  sr_parse = nltk.ShiftReduceParser(grammar)
3  for tree in sr_parse.parse(sent):
4      print(tree)
5
6  sr_parse = nltk.ShiftReduceParser(grammar, trace=2)
7  for tree in sr_parse.parse(sent):
8      print(tree)

```

RecursiveDescentParser 递归下降分析器：

- 左递归产生式，如 NP -> NP PP，会进行死循环
- 浪费了很多时间处理不符合输入句子的词和结构
- 回溯过程中可能会丢弃分析过的成分，它们将需要在之后再次重建

ShiftReduceParser 移进-规约分析器：

- 不执行任何回溯，所以不能保证一定能找到一个文本的解析（即使真的存在）
- 即使有多个解析最多也只能找到一个
- 每个结构只建立一次

通过优先执行规约操作来解决移进-规约冲突

交互式文法编辑器

```
nltk.app.chartparser()
```

依存文法

```

1  groucho_dep_grammar = nltk.grammar.DependencyGrammar.fromstring("""
2  'shot' -> 'I' | 'elephant' | 'in'
3  'elephant' -> 'an' | 'in'
4  'in' -> 'pajamas'
5  'pajamas' -> 'my'

```

```

6  """')
7
8  print(groucho_dep_grammar)
9
10 pdp = nltk.ProjectiveDependencyParser(groucho_dep_grammar)
11 sent = 'I shot an elephant in my pajamas'.split()
12 trees = pdp.parse(sent)
13 for tree in trees:
14     print(tree)

```

特征结构

```

1  print(nltk.FeatStruct(''
2  [A = 'a',
3  B = (1)[C = 'c'],
4  D -> (1),
5  E -> (1)]
6  '''))
7
8  fs1 = nltk.FeatStruct(NUMBER = 74,
9                        STREET = 'run Pascal')
10 fs2 = nltk.FeatStruct(CITY = 'Paris')
11
12 print(fs1.unify(fs2))
13
14 fs0 = nltk.FeatStruct(A = 'a')
15 fs1 = nltk.FeatStruct(A = 'b')
16 fs2 = fs0.unify(fs1)
17 print(fs2)
18
19 fs1 = nltk.FeatStruct(''
20 [ADDRESS1 = [NUMBER = 74, STREET = 'run Pascal']]
21 ''')
22 fs2 = nltk.FeatStruct(''
23 [ADDRESS1 = ?x,
24 ADDRESS2 = ?x]
25 ''')
26 print(fs1)
27 print(fs2)
28 print(fs2.unify(fs1))

```

十一、语言数据管理

```

1  s1 = '01010101'
2  s2 = '00001111'
3  s3 = '00001110'
4
5  # 在指定窗口大小下计算两个字符串的差异大小
6  nltk.windowdiff(s1, s2, 3)
7  nltk.windowdiff(s2, s3, 3)

```

TIMIT

语音语料库，数据是文本形式

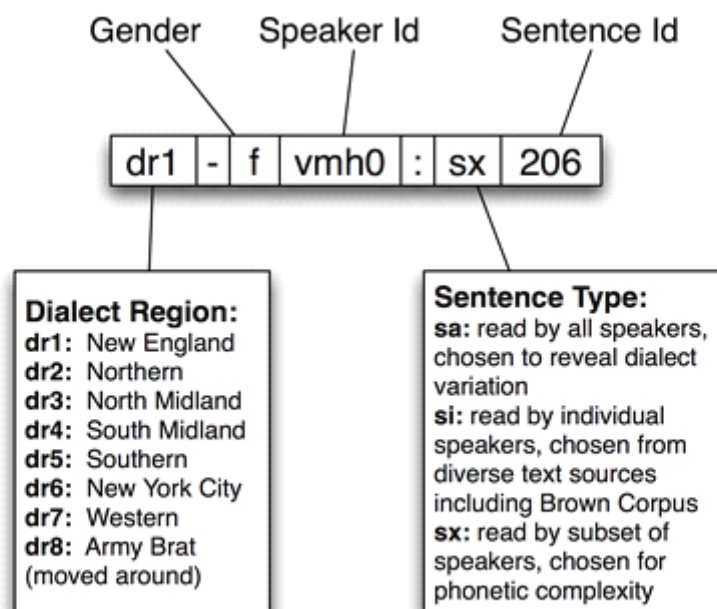


图 11-1. TIMIT 标识符的结构：每个记录使用说话者的方言区、性别、说话者标识符、句子类型、句子标识符组成的一个字符串作为标签。

```

1 import nltk
2
3 phonetic = nltk.corpus.timit.phones('dr1-fvmh0/sa1')
4 nltk.corpus.timit.word_times('dr1-fvmh0/sa1')
5
6 timitdict = nltk.corpus.timit.transcription_dict()
7 timitdict['greasy']
8
9 nltk.corpus.timit.spkrinfo('dr1-fvmh0')

```

XML

```

1 merchant_file = nltk.data.find('corpora/shakespeare/merchant.xml')
2 raw = open(merchant_file).read()
3 print(raw[0: 168])
4
5 from xml.etree import ElementTree
6
7 merchant = ElementTree.parse(merchant_file)
8 merchant
9
10 speaker_seq = [s.text for s in
11 merchant.findall('ACT/SCENE/SPEECH/SPEAKER')]
12 speaker_freq = nltk.FreqDist(speaker_seq)
13 top5 = list(speaker_freq.keys())[:5]

```

Toolbox

```

1 from nltk.corpus import toolbox
2
3 lexicon = toolbox.xml('rotokas.dic')
4 # 1、索引访问
5 lexicon[3][0]
6 lexicon[3][0].tag

```



```

7 lexicon[3][0].text
8
9 # 2、路径访问
10 [lexeme.text.lower() for lexeme in lexicon.findall('record/lx')][:10]
11
12 # 格式化
13 html = '<table>\n'
14 for entry in lexicon[70: 80]:
15     lx = entry.findtext('lx')
16     ps = entry.findtext('ps')
17     ge = entry.findtext('ge')
18     html += ' <tr><td>%s</td><td>%s</td><td>%s</td></tr>\n' % (lx, ps, ge)
19
20 html += '</table>'
21 print(html)
22
23 # 为每个条目计算字段的平均个数
24 from nltk.corpus import toolbox
25
26 lexicon = toolbox.xml('rotokas.dic')
27 sum(len(entry) for entry in lexicon) / len(lexicon)

```

OLAC元数据

元数据：关于数据的结构化数据

OLAC：开放语言档案社区