

nl_ex

April 19, 2016

1 自然语言处理实验报告

专业：智能科学
班级：1301
学号：0909111024
姓名：张日月明

1.1 实验 1

怎样载入自己的英文语料库（obama.txt），在自己的语料库中找出频率小于8，排名前5的词和其频率。（使用nltk的英文分

```
In [7]: import nltk
```

```
text = open('obama.txt', 'rb').read()
tokens = nltk.word_tokenize(text)
freqdist = nltk.FreqDist(tokens).items()
answer = [k for k in freqdist if k[1]<8]
answer.sort(lambda x,y:cmp(y[1],x[1]))
print answer[0:5]
```

```
[('go', 7), ('work', 7), ('learn', 7), ('give', 7), ('something', 7)]
```

1.2 实验 2

写程序处理布朗语料库，找到一下的答案：

- (1) 哪些名词常以他们复数形式而不是它们的单数形式出现？（只考虑常规的复数形式，-s后缀形式的）。

```
In [8]: print('(1):')
wsj = nltk.corpus.treebank.tagged_words()
word_tag_fd = nltk.FreqDist(wsj)
#print([word + '/' + tag for (word, tag) in word_tag_fd if tag.startswith('NNS')])
print([word for (word, tag) in word_tag_fd if tag.startswith('NNS')])
```

(1):

'shipments', u'BIBBSes', u'famalkess', u'qfzanes', u'sgadsunst', u'markshers', u'generalizajons', u'Expetislonerpolitest,

- (2) 选择布朗语料库的不同部分（其他目录），计数包含wh的词，如：what，when，where，who和why。

```

In [9]: print('(2):')
        news_text = nltk.corpus.brown.words(categories='news')
        freqdist = nltk.FreqDist([w.lower() for w in news_text])
        demands = [w.lower() for w in freqdist if w.startswith('wh')]
        for k in demands:
            print(k + ': ', freqdist[k])

(2):
(u'whose:', 22)
(u'whirling:', 1)
(u'whiz:', 2)
(u'whip:', 2)
(u'wherever:', 1)
(u'where:', 59)
(u'whitfield:', 1)
(u'what's:', 1)
(u'wheeler:', 2)
(u'wheeled:', 2)
(u'whopping:', 1)
(u'whelan:', 1)
(u'whereby:', 3)
(u'whichever:', 1)
(u'whipped:', 2)
(u'wheels:', 1)
(u'wheat:', 1)
(u'whatever:', 2)
(u'while:', 55)
(u'whisking:', 1)
(u'white-clad:', 1)
(u'whitey:', 1)
(u'whites:', 2)
(u'whitney:', 1)
(u'whipple:', 1)
(u'wholly:', 1)
(u'wheel:', 4)
(u'when:', 169)
(u'whiplash:', 1)
(u'white:', 57)
(u'whee:', 1)
(u'wholly-owned:', 1)
(u'wholesale:', 1)
(u'what:', 95)
(u'whom:', 8)
(u'whims:', 1)
(u'which:', 245)
(u'who:', 268)
(u'why:', 14)
(u'whole:', 11)
(u'whether:', 18)

```

1.3 实验3

输出brown文本集名词后面接的词性

```
In [10]: brown_lrnd_tagged = nltk.corpus.brown.tagged_words(categories='learned')
tags = [b[1] for (a, b) in nltk.bigrams(brown_lrnd_tagged) if a[1].startswith('N')] #a[0] == '
freqdist = nltk.FreqDist(tags)
freqdist.tabulate()
```

IN	.	,	NN	CC	NNS	BEZ	MD	CS	VBN	VBZ
11913	4998	4968	3725	2577	1737	1336	1080	827	732	

1.4 实验4

句法分析演示

运行以下代码即会弹出 GUI 界面查看演示

```
nltk.app.srparser()
nltk.app.rdparsr()
```

1.5 实验5

中文分词

实验数据：

- word_freq_list.txt 分词词典
- pku_test.txt 未经过分词的文档文件
- pku_test_gold.txt 经过分词的文档文件

```
In [11]: import jieba
#import sys
#reload(sys)
#sys.setdefaultencoding('utf-8')

text = open('pku_test.txt', 'rb').read()

jieba.load_userdict('word_freq_list.txt')
seg_list = jieba.cut(text)

#open('pku_test_seg.txt', 'wb').write(' '.join(seg_list))
print(' '.join(seg_list))
```

--

1

2001

21

20

20

20

"

12 31

2001 1 1

12 31

4

- - - 2001

2000

40

12 31

12 30

55.6

11.1

25

2000

12 31

30
20

3000

28

" "

"

"

"

... .. "

12 31

20

" " - - -

21

WTO

20

- - -

21

2001 1 1

20

0

[illegible]

11	8					14		3	8	9	14		
11	9	200		5	"	"						841	1200
12	27												
12	30				1								
2000	12	31											
	21					20							
	10					10	1993	9					
1993		"	"			7	"	"			15		
				"	"							6	
									"	"			
12	31			31									
12	19			89		24							
19				"				"			.		
12	30							.		30			
12	31			20									
28					3								
								2000		"	"	9	
			"	"									"
12	30												
12	31						2001						
	31						10					2020	
2001											"	"	
	"						"						
"											"		

[illegible]

[illegible]

2000

12 30 " " " " 250
1999

" " IT
21 " "
IT 2000 23
90 " "
1915 3 4 47
- - -

" AM21B " 37 1 49

" " " happynewyear . txt . vbs "
" " 60
5 1932 90 10 1916 1920
" " 60 6
300 " 1 5

" " " " " "

1 1 1
2001 " "
" "

A D

" "

.

1 1

1

" "

" "

" "

-

A B

" " " " " "

"

" "

" "

" "

1 1

1 1
12 31 21 23 21 - - - " " 100 2.5
2000 12 31 8 30 1
9 4000 "
" 3000 2001 9 2100 21
12 " 2001 "

1 1

1 1 15
8 50 15
6 15
1 1 15 " " 3
" " " " " "
600
1 1 500
" "

12 31
31
31
31
30

2001

|| ||

" "

2000

21

$$\begin{array}{r} 12 \quad 31 \\ 31 \\ \cdot \end{array}$$

" "

2001

|| ||

1 1

10

1 1

12 31 31

12 31 31 242 338

$$\begin{array}{ccccccc} 1 & & 1 & & & & 1 \end{array}$$

12 30 16

12 28 1947

1932	1954	1963	30		1969
		1998		"	"
		21			

29
12 30

4

1	2000	12	31	
2	2000	12	31	
3	2001	1	1	2001
4	2000	12	31	

2001	42	20
------	----	----

1

12 31

12 31 " "

30

|| || ||

1 1

1 1 21 21

12 31 2000

12 30

50

2000

" "

" 1 1 " "

" - - - " "

100 2008 6000

1 1 " 250

1 1 " "

1998 3 28 56

160 7.8 39 6.3 2000

3 10 1.8

1 1

1 2 20 - 1 3 20

1 2 3 4 - 6 6 - 7

9 °C 0 °C

4 °C 10 °C

9 °C 0 °C

10 °C 19 °C

5 °C 6 °C

6 °C 14 °C

10 °C 5 °C

3 °C 5 °C

17 °C 5 °C

5 °C 1 °C

25 °C 13 °C

1 °C 8 °C

12 °C 4 °C

3 °C 10 °C

23 °C 17 °C

4 °C 8 °C

30 °C 23 °C

12 ℃ 23 ℃
 4 ℃ 12 ℃
 15 ℃ 20 ℃
 0 ℃ 9 ℃
 13 ℃ 24 ℃
 4 ℃ 9 ℃
 7 ℃ 10 ℃
 7 ℃ 10 ℃
 16 ℃ 21 ℃
 3 ℃ 7 ℃
 4 ℃ 7 ℃
 7 ℃ 17 ℃
 24 ℃ 34 ℃
 6 ℃ 10 ℃
 17 ℃ 28 ℃
 1 ℃ 6 ℃
 10 ℃ 20 ℃
 6 ℃ 4 ℃
 11 ℃ 21 ℃
 14 ℃ 3 ℃
 5 ℃ 2 ℃
 8 ℃ 4 ℃
 2 ℃ 6 ℃
 11 ℃ 4 ℃
 2 ℃ 5 ℃
 14 ℃ 21 ℃
 3 ℃ 6 ℃
 16 ℃ 20 ℃
 2 ℃ 5 ℃

3 ℃ 5 ℃
 2 ℃ 16 ℃
 5 ℃ 12 ℃
 2 ℃ 6 ℃
 3 ℃ 7 ℃
 2 ℃ 6 ℃
 13 ℃ 5 ℃
 35 ℃ 25 ℃
 7 ℃ 10 ℃
 4 ℃ 8 ℃
 - - -

2001 1 1 " "
 9 30 " " " " " "
 20 " 20 20 30 " " 20
 " " " "
 17 " 20 " - - -
 " 20 " - - -

[illegible]

1 1 " " 5 55 - - -
175 60 . " " 15 - - -
5 55 " " 8 150 - 300
2000 12 28 5000 19 85
75
1 1
2001 2008
2000 12 31 7 1 25 46
" " 12 31
" " www . in paku . go . jp 100 1 50
" " 202
IT 1 IT 90
1 1
7.58 1 1.45 - - 12 31 6 2000
1
1.3 1963 153
90
" "
2000
20 18 40 " 21 " 8500 12 15 22
2 600 1700
2001 573 20 2030 450 2000 280
2000 2480 " " 2000 6
80 19 2000 2000 27 2000 2001
3 4 14

2000 14 97

6.6 12 31 . -

6.6 1 3 66

20 10 - - -

1 2 " "

3

" " - - -

2000 12 31 .

160 50

5

"

" "

1

12 20 2001 2

" " 1999

20

160 50 30

4

1963 4

118 7000 4500

1996 1996 40

1999 50 3

" "

" "

" "

1890

1910

1840

" "

1900

1868

1949

" "

" "

" "

" "

"

" "

1996

2001

" "

- - -

100733

010 - 65092341

caibian3 peopledaily com cn

300

3

" - - -

" "

" "

" "

- - -

2000 11 22 1999 11 · 16

1999 11 16 10 25

" "

4 24

" "

" "

" " " " " "

" " " " " "

" " " " " "

162 2465 85 93

39

162

188 3402

206

" " "

1

" "

1

" "

1

- - -

1

1958

1

2001

2000

— — —

1

11

11

11

11

2000

11

11

11

11

2000
2000

5

13

6

1

1

1

2000 12 31

1

— — —

11

11

11

11

— — —

2000 9 7

100 1900 8

9

11

100

11

“ ”

1949	10	1
------	----	---

A	C
1	1
2	2
3	3
4	4
5	5
6	6
7	7
8	8
9	9
10	10
11	11
12	12
13	13
14	14
15	15
16	16
17	17
18	18
19	19
20	20
21	21
22	22
23	23
24	24
25	25
26	26
27	27
28	28
29	29
30	30
31	31
32	32
33	33
34	34
35	35
36	36
37	37
38	38
39	39
40	40
41	41
42	42
43	43
44	44
45	45
46	46
47	47
48	48
49	49
50	50
51	51
52	52
53	53
54	54
55	55
56	56
57	57
58	58
59	59
60	60
61	61
62	62
63	63
64	64
65	65
66	66
67	67
68	68
69	69
70	70
71	71
72	72
73	73
74	74
75	75
76	76
77	77
78	78
79	79
80	80
81	81
82	82
83	83
84	84
85	85
86	86
87	87
88	88
89	89
90	90
91	91
92	92
93	93
94	94
95	95
96	96
97	97
98	98
99	99
100	100

11

11

19

1955

50

1971 10

11

" 80

26

" "

1 2

1991 200 1000 2000 3 4

800

" 3 5 " " 10 1 " ... "

1 1

" "

" " - -

300 16

300.5 16

11 20 13.6

244 18.3

2000 244.8 1999 18.3 " "

25 3 18.8

2000 1141.6 24.6 449.9 12.3 497.9

" "

682

2000 682.6

1996 300 1998 1999 400 500 2000

600 17.6

2000 600.3 90 17.6 10 224

22

1994

1 2 WTO 8 2 3 100

" "

" "

24

10 - - -

171

1 2 " " " " 5 6250 3200

1 2 1 10 " "

" " 3 " "

45

1

2

3 " 1 "

" "

1 2

[illegible]

1 42 " "

1

6000 " 1999 12 31 8 1 200 10

27 1 1

1 1

8 31 300

5 8 130 200

" "

1 1 " "

1996 " "

1940

4 " " 50 2

- - -

10

1996 3 21 6

3

5 2000 100 1

3 10

10 2300 6000 1

8 5000

" " 50

5 12

683

" " 3000

10 " " 4

105 20 3.4 161 7

1 2

GDP

2000 1 - 11 2000 " "

1812 1999 22.0 1228 2000 2 .. 2

1 2 5.5 207 138 76

4 5200

2000 9 2001 9

9957.5 1.3 317.2 11.1 32.8

2000 12 28 75

			7000	8		183	810		
1	2			500					
1	2								
		1							
2000	12	20				2000	12	11	88
			1928		1929	1935			
					1952		1956		19
"			"						
1	2								
1	2								
		"	"						
				17	"	"			
1	2	9					20		
					20			20	2001
"	"					"	"		
			"	"					
							5		10
1	3	20	- 1	4	20				
						3	4		4
10 °C	2 °C								
4 °C	10 °C								
7 °C	1 °C								
9 °C	18 °C								
5 °C	4 °C								
8 °C	13 °C								
8 °C	6 °C								
6 °C	3 °C								
	16 °C	6 °C							
	5 °C	1 °C							
27 °C	16 °C								
2 °C	4 °C								

12 ℃ 4 ℃
 3 ℃ 9 ℃
 29 ℃ 20 ℃
 3 ℃ 7 ℃
 30 ℃ 22 ℃
 14 ℃ 22 ℃
 6 ℃ 11 ℃
 12 ℃ 16 ℃
 0 ℃ 7 ℃
 17 ℃ 25 ℃
 5 ℃ 9 ℃
 8 ℃ 10 ℃
 8 ℃ 12 ℃
 16 ℃ 21 ℃
 3 ℃ 6 ℃
 4 ℃ 9 ℃
 7 ℃ 18 ℃
 24 ℃ 34 ℃
 6 ℃ 9 ℃
 17 ℃ 30 ℃
 2 ℃ 8 ℃
 11 ℃ 29 ℃
 5 ℃ 3 ℃
 12 ℃ 22 ℃
 13 ℃ 4 ℃
 4 ℃ 1 ℃
 11 ℃ 2 ℃
 1 ℃ 5 ℃
 10 ℃ 4 ℃
 4 ℃ 6 ℃
 14 ℃ 23 ℃
 6 ℃ 10 ℃
 16 ℃ 20 ℃
 5 ℃ 1 ℃

 2 ℃ 7 ℃
 4 ℃ 13 ℃
 4 ℃ 14 ℃
 2 ℃ 7 ℃
 3 ℃ 9 ℃
 4 ℃ 12 ℃
 10 ℃ 4 ℃
 30 ℃ 19 ℃
 6 ℃ 11 ℃
 4 ℃ 8 ℃

		2000	12	30			8.9
1							
1							
	1				1952	679	679
1							
	1		1	1998	87599		

- - -

" "

1 2

19.8 7.1

12

3

9

28

50

3

"

"

"

"

"

" - - -

50

20

" "

"

"

"

"

1 2

1 11

5

260

64

24

42

13

7

"

"

"

"

"

"

"

"

"

"

"

"

"

31

30

440

"

"

"

"

"

"

"

"

"

"

"

" - -

"

"

"

"

"

"

"

[illegible]

[illegible]

A

1 2 A " "

2001 1 1 " "

9 300

1 2 2000 28 16 15

2000 18 110 1 2 15 4

14 2 30 22 8 21 14 3 2 6 5 3 3 3 4

2000 1408 1042 2000 1378 865

www . people . com . cn sports

3000

4000 1000

" " 17

20

1 2 " "

3

4

2000 - 2001 " " 2000 12 30

2000 12 24

2000 12 28 - - - 3

2003 7.4 217

4

1 1 2001

1 2000 12 30 2001 1 1 " "

2

3

4

- - -

" "

" "

- - -

" "

" "

" " "

- - -

- - -

" " - - -

" "

10

" " " "

" "

" "

.

1

2000 11 5

"

"

" "

400

380

3000

1996

1997

1000

8000

5000

700

1

1998 11 15

1999 8 2

8 21

1999 10 9

1999 11 1

11 10

1999 12 8

1983

100

5000

2000 1 21

2

7000

5

"

"

" "

"

" "

" "

" "

- - -

50

84

1998

- - -

... ..

1998

79

15

65

20

" "

" "

" "

" "

39