

南京邮电大学

实 验 报 告

(2023 / 2024 学年第 1 学期)

课程名称 _____ 自然语言处理 _____

实验名称 _____ 语料预处理和语料库构建与应用 _____

实验日期 _____ 2023/12/1 _____

指导教师 _____ 朱博 _____

学生姓名 _____ 单家俊 _____ 学号 _____ B21080526 _____

专 业 _____ 人工智能 _____ 年级 _____ 大三 _____

实 验 报 告

一、实验目的和要求

- 1、掌握语料库的构建方法。
- 2、掌握语料库的预处理技术。
- 3、掌握语料库的简单分析方法

二、实验环境(实验设备)

PC、PyCharm

三、实验原理及内容

- 1、构建金庸武侠小说集语料库

(1) 收集金庸武侠小说的 txt 文本文件，完成数据采集和预处理，获取保存的文件列表。

使用的核心函数有：nltk.corpus 中的 PlaintextCorpusReader

数据采集和预处理代码：

```

1  # -*- coding: gbk -*-
2
3  import jieba
4  import nltk
5  import re
6  from nltk.book import *
7  from nltk.corpus import PlaintextCorpusReader
8  import matplotlib.pyplot as plt
9
10
11  corpus_root = 'D:/QQ/1413679561/FileRecv/data'
12  # file_pattern = r'*.txt'
13  corpus_reader = PlaintextCorpusReader(corpus_root, fileids='*')
14  file_list = corpus_reader.fileids()
15  print(file_list)
16
17  # 加载停用词列表
18  stop_words = []
19  path = 'E:/NLP/stopword.txt'
20  for line in open(path, encoding='utf8'):
21      # print(line)
22      line = line.strip()
23      stop_words.append(line)
24
25
26  # 去除停用词
27  result = []
28  with open(file=r'D:/QQ/1413679561/FileRecv\data\金庸-神雕侠侣.txt', mode='r', encoding='gbk') as f:
29      content = f.read()
30      cleaned_data = ''.join(re.findall(pattern='[\u4e00-\u9fa5]', content))
31      content_list = jieba.lcut(cleaned_data, cut_all=False)
32      for word in content_list:
33          if word not in stop_words:
34              result.append(word)
35  text = ''.join(result)
36  # print(text)
37
38  with open(file='金庸-神雕侠侣1.txt', mode='w') as f:
39      f.write(text)
40
41  with open('金庸-神雕侠侣1.txt', 'r', encoding='gbk') as f:
42      str = f.read()
43      # len(set(str))
44      # len(str)/len(set(str))
45
46  # cleaned_data = ''.join(re.findall('[\u4e00-\u9fa5]', str))
47  wordlist = jieba.lcut(str)
48  text = nltk.Text(wordlist)
49  print(text)

```

结果:

```

['金庸-书剑恩仇录.txt', '金庸-侠客行.txt', '金庸-倚天屠龙记.txt', '金庸-天龙八部.txt', '金庸-射雕英雄传.txt', '金
庸-白马啸西风.txt', '金庸-碧血剑.txt', '金庸-神雕侠侣.txt', '金庸-笑傲江湖.txt', '金庸-越女剑.txt', '金庸-连城诀
.txt', '金庸-雪山飞狐.txt', '金庸-飞狐外传.txt', '金庸-鸳鸯刀.txt', '金庸-鹿鼎记.txt']
Building prefix dict from the default dictionary ...
Loading model from cache C:\Users\ADan\AppData\Local\Temp\jieba.cache
Loading model cost 0.573 seconds.
Prefix dict has been built successfully.
<Text: 全本 全集 精校 小说 更 资源 下载 声明...>
杨过 小龙女 赵志敬 周伯通 还 笑 武三通 武敦儒 洪七公 耶律齐

```

(2) 选其中一本小说，利用 NLTK 中的基本函数，在你感兴趣的小说中，搜索相似词语、指定词语、搭配词语、查询文本词汇频数分配等，画出词汇分布离散图。（参考教材 21 页-22 页（1）-（11）内容，每条都要做）

内容 1:

```
# (1)
text.similar(word='李莫愁',num=10)
print('\n')
杨过 小龙女 赵志敬 周伯通 还 笑 武三通 武敦儒 洪七公 耶律齐
```

内容 2:

```
# (2)
text.concordance(word='侠',width=30,lines=3)
print('\n')
Displaying 3 of 54 matches:
江湖 上 赫赫有名 神雕 侠 四川人 问道 叫作 神雕
四川人 问道 叫作 神雕 侠 汉子 道 这位 大侠 行侠
小 功劳 那天 晚上 神雕 侠 突然 来到 临安 带领 伙
```

内容 3:

```
# (3)
text.collocations()
print('\n')
潇湘子 尹克西；丘处机 王处一；潇湘子 尼摩星；甄志丙 赵志敬；七十二路 空明拳；鞠躬尽瘁 死而后已；王志坦 宋德方；武敦儒
武修文；全真教 第三代；老顽童 周伯通；李志常 王志坦；大力神 史季强；马真人 丘真人；生死相许 天南地北；小龙女 微微一笑；
郝大通 孙不二；广宁子 郝大通；为国为民 侠之大者；马道长 丘道长；倒大出 意料之外
```

内容 4:

```
60 # (4)
61 text.common_contexts(['杨过','小龙女'])
62 print('\n')
去_道 罢_道 听_说 问_道 是不是_道 听_道 师父_道 干什么_道 说_道 没有_道 知道_道 好_道 人_道 出来_道 道_道
此时_已 走_道 死_道 只_一人 见_脸上
```

内容 5:

```

63 # (5)
64 print(len(str))
65 print('\n')

```

700054

内容 6:

```

66 # (6)
67 print(set(str))
68 print('\n')

```

{'根', '血', '迹', '猫', '埋', '绍', '漱', '湍', '牢', '临', '淘', '邀', '拳', '抑', '忸', '桌', '盟', '仄', '传', '服', '险', '批', '淋', '梨', '廿', '揖', '虬', '前', '履', '嚼', '焉', '族', '内', '砥', '秒', '宇', '话', '髟', '稔', '机', '是', '除', '寒', '池', '程', '哈', '蜈', '介', '瑶', '伍', '袋', '式', '拿', '噪', '响', '目', '尔', '滢', '却', '稻', '脾', '仗', '骄', '诚', '絮', '健', '咋', '旋', '玉', '吱', '典', '然', '要', '居', '肋', '孽', '跽', '缪', '玠', '预', '赎', '软', '炆', '傍', '靠', '蔡', '测', '纶', '闲', '竖', '套', '咽', '单', '疏', '德', '筋', '减', '污', '岩', '搭', '迥', '翔', '冒', '凰', '演', '惟', '产', '骑', '喝', '脐', '陪', '堆', '拌', '斗', '漉', '疹', '霞', '剖', '咫', '倜', '蕃', '梅', '角', '凸', '捷', '著', '骨', '姨', '够', '控', '赶', '越', '眶', '胎', '繁', '兆', '疤', '贝', '逢', '赅', '伯', '冢', '搓', '搓', '癩', '敝', '斌', '多', '套', '洵', '灶', '膝', '溪', '撮', '窠', '汤', '珑', '议', '射', '床', '犖', '微', '敖', '佳', '川', '抖', '症', '跃', '磊', '绌', '灵', '茧', '薛', '括', '巷', '翌', '予', '哀', '樊', '馈', '缈', '蛰', '餐', '功', '们', '雨', '涕', '芥', '缝', '霏', '体', '沿', '良', '酷', '虚', '勒', '郡', '佬', '掌', '簋', '鸣', '聘', '寡', '座', '蟹', '渗', '圃', '艺', '约', '嗣', '瞋', '冢', '扼', '色', '折', '逢', '硬', '茅', '喉', '恹', '响', '叶', '允', '煊', '崕', '蹇', '忡', '头', '耻', '刳', '仔',

内容 7:

```

69 # (7)
70 sorted(set(str))
71 print(f'词汇表大小:{len(set(str))}')
72 print(f'每个词平均使用次数:{len(str)/len(set(str))}')
73 print('\n')

```

词汇表大小:4014

每个词平均使用次数:174.40308918784254

内容 8:

```

74 # (8)
75 print('杨过一词出现的次数: ', end='')
76 print(str.count('杨过'))
77 fdist = FreqDist(str)
78 print('出现频率最高的:')
79 print(fdist.most_common(30))

```

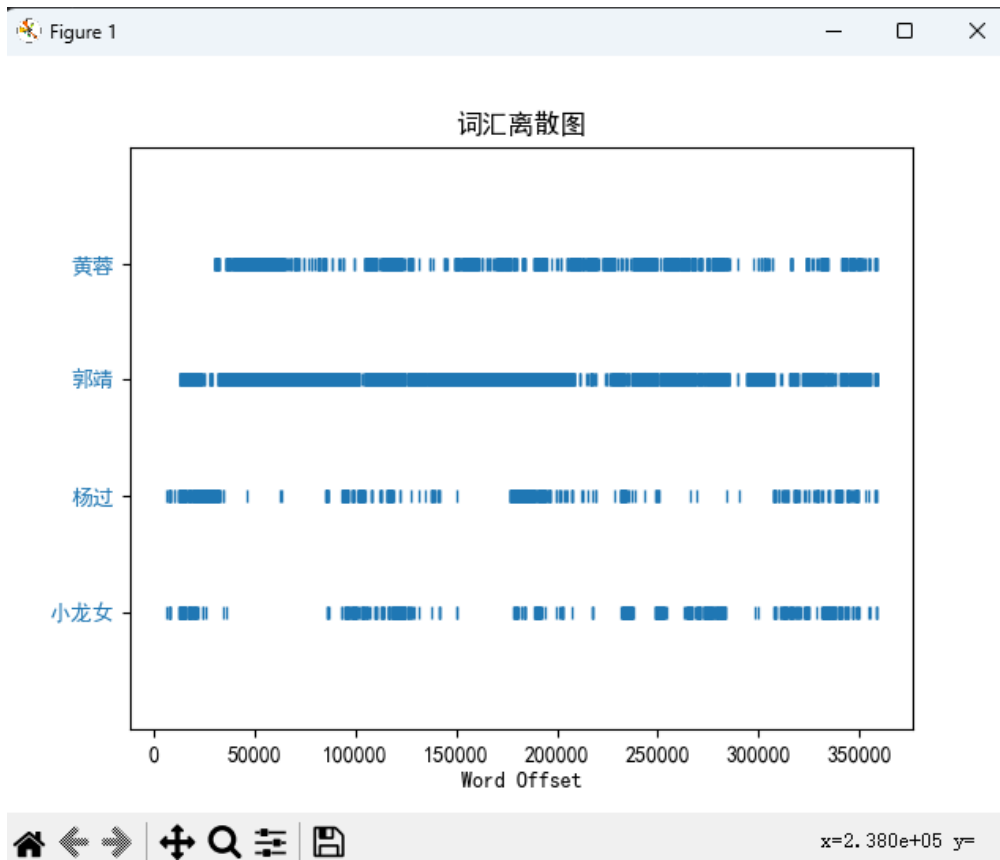
出现频率最高的:

内容 9:

['一', '丁', '七', '万', '父', '三', '上', '下', '不', '与', '丐', '丑', '专', '且', '世', '丘', '丙', '业', '丛', '东', '丝', '丞', '丢', '两', '严', '丧', '个', '丫', '习', '中', '丰', '串', '临', '丸', '丹', '为', '主', '丽', '举', '乃', '久', '么', '义', '之', '乌', '乍', '乎', '乏', '乐', '兵', '兵', '乔', '乘', '乘', '乙', '九', '乞', '也', '习', '乡', '书', '买', '乱', '乳', '乾', '了', '予', '争', '事', '二', '于', '亏', '云', '互', '五', '井', '亚', '些', '巫', '亡', '亢', '交', '亥', '亦', '产', '亨', '亩', '享', '京', '享', '高', '亲', '裹', '人', '什', '仁', '汀', '仄', '仅', '仆', '仇', '今', '介', '仍', '从', '仑', '仓', '仔', '他', '仗', '付', '仙', '仞', '代', '令', '以', '仪', '仆', '们', '仰', '仲', '价', '任', '份', '仿', '企', '伊', '伍', '伎', '伏', '伐', '休', '众', '优', '伙', '会', '伞', '伟', '传', '伤', '伙', '伦', '伦', '伪', '仁', '伯', '估', '伴', '伶', '伸', '伺', '似', '伽', '佃', '但', '位', '低', '住', '佐', '佑', '体', '何', '佗', '余', '佛', '作', '候', '你', '佩', '佬', '伴', '佳', '佻', '使', '佻', '侃', '侄', '例', '侍', '侏', '供', '依', '侠', '侶', '傣', '伙', '侧', '佝', '侏', '侮', '侯', '侵', '便', '促', '俊', '俏', '俐', '俗', '俘', '埋', '保', '侯', '信', '铸', '钐', '俩', '修', '俯', '俱', '俾', '倍', '候', '倒', '偃', '倘', '候', '伺', '伺', '借', '倡', '倦', '偃', '倩', '倪', '傣', '值', '倾', '假', '偈', '偈', '偈', '偏', '偕', '做', '停', '健', '德', '偷', '傚', '傀', '俑', '傣', '侯', '傣', '催', '傣', '傣', '像', '僧', '僧', '僧', '僧', '僧', '僧', '儿', '刀', '允', '元', '兄', '充', '兆', '先', '光', '克', '免',

内容 10:

```
82 # (10)
83 plt.rcParams['font.sans-serif'] = 'SimHei'
84 words = ['小龙女', '杨过', '郭靖', '黄蓉']
85 nltk.draw.dispersion.dispersion_plot(text, words, title='词汇离散图')
```

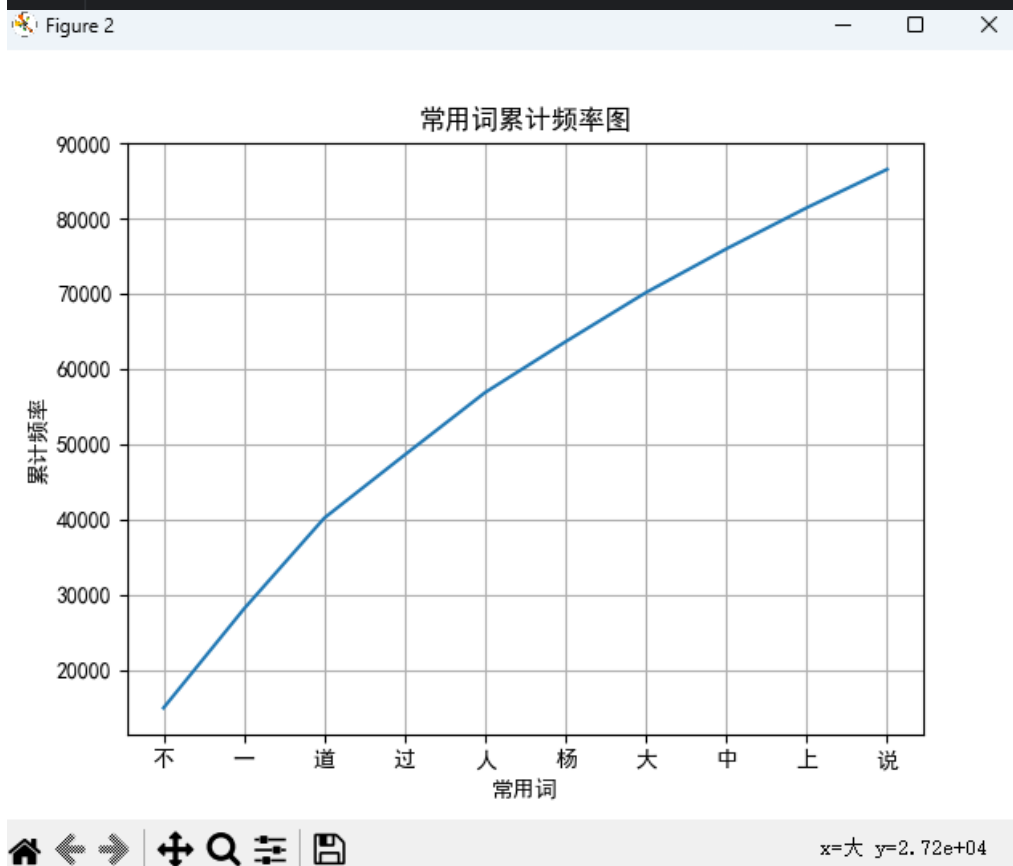


内容 11:

```

82 # (10)
83 plt.rcParams['font.sans-serif'] = 'SimHei'
84 words = ['小龙女', '杨过', '郭靖', '黄蓉']
85 nltk.draw.dispersion.dispersion_plot(text, words, title='词汇离散图')
86 # (11)
87 fig = plt.figure()
88 plt.grid()
89 fdist1 = dict(fdist)
90 fdist1 = sorted(fdist.items(), key= lambda x: x[1], reverse=True)
91 x = []
92 y = []
93 for i in range(10):
94     x.append(fdist1[i][0])
95     y.append(fdist1[i][1])
96 t = 0
97 for i in range(len(y)):
98     y[i] = y[i] + t
99     t = y[i]
100 plt.plot(*args: x, y)
101 plt.title('常用词累计频率图')
102 plt.ylabel('累计频率')
103 plt.xlabel('常用词')
104 plt.show()

```



2、金庸武侠小说集语料库分析

(1) 选其中一本小说进行分析。

预处理代码（去除停用词、去除非正文段落、增添特殊停用词）：

```
3 import jieba
4 import nltk
5 import re
6 from nltk.book import *
7 # 加载停用词列表
8 stop_words = []
9 path = 'E:/NLP/stopword.txt'
10 for line in open(path, encoding='utf8'):
11     # print(line)
12     line = line.strip()
13     stop_words.append(line)
14 # 增添特殊停用词
15 added_stopword = ['具体来说', '再次', '分期分批', '切莫', '到目前为止', '猛然间', '$']
16 stop_words.extend(added_stopword)
17 # 去除非正文段落
18 with open(file=r'D:\QQ\1413679561\FileRecv\data\金庸-天龙八部.txt', mode='r', encoding="gbk") as f:
19     content_list = f.readlines()
20     content_list = content_list[52:]
21     content_list = content_list[:-40]
22 with open(file='金庸-天龙八部1.txt', mode='w') as f:
23     content_str = ''.join(content_list)
24     f.write(content_str)
25 # 去除停用词
26 result = []
27 chinese_result = []
28 with open(file='金庸-天龙八部1.txt', mode='r', encoding="gbk") as f:
29     content = f.read()
30     cleaned_data = ''.join(re.findall(pattern='[\u4e00-\u9fa5]', content))
31     content_list = jieba.lcut(content, cut_all=False)
32     chinese_content_list = jieba.lcut(cleaned_data, cut_all=False)
33     for word in content_list:
34         if word not in stop_words:
35             result.append(word)
36     for word in chinese_content_list:
37         if word not in stop_words:
38             chinese_result.append(word)
39 str = ''.join(result)
40 chinese_str = ''.join(chinese_result)
41
```

(2) 统计小说中总用词量和平均每个词的使用次数

```
46 # (2)
47 print(f'总用词量: {len(set(str))}')
48 print(f'平均每个词的使用次数{len(str)/len(set(str))}')
总用词量: 4192
平均每个词的使用次数211.27862595419847
```

(3) 自行指定 3 个词语，查看小说中的使用次数

```
49 # (3)
50 print(f'乔峰一次的使用次数: {str.count("乔峰")}')
51 print(f'段誉一次的使用次数: {str.count("段誉")}')
52 print(f'虚竹一次的使用次数: {str.count("虚竹")}')
```

乔峰一次的使用次数: 1243
段誉一次的使用次数: 3534
虚竹一次的使用次数: 1671

(4) 统计并输出前 30 个高频字符和对应出现次数

```
53 # (4)
54 fdist = FreqDist(str)
55 print(fdist.most_common(30))
```

[('u3000', 22150), ('不', 19690), ('道', 15684), ('一', 14460), ('人', 12617), ('n', 11222), ('大', 9519), ('说', 7124), ('中', 6936), ('子', 6612), ('上', 6529), ('来', 6528), ('下', 6263), ('之', 6240), ('段', 5733), ('手', 5465), ('去', 5151), ('心', 5021), ('出', 5019), ('便', 5016), ('... ', 4968), ('见', 4615), ('声', 4409), ('是', 4366), ('身', 4228), ('到', 3753), ('只', 3714), ('得', 3680), ('誉', 3658), ('好', 3578)]

(5) 使用 jieba 对小说正文前十句话进行分词

```
56 # (5)
57 # print(content_str.split('。')[0:10])
58 tensentences = content_str.split('。')[0:10]
59 for sentence in tensentences:
60     splited_sentence = jieba.lcut(sentence)
61     print(splited_sentence)
```

['一', ' ', '青衫', '磊落', '险峰', '行', 'n', 'n', 'u3000', 'u3000', '青光闪', '动', ' ', ' ', '一柄', '青钢剑', '候地', '刺', '出', ' ', ' ', '指向', '中年', '汉子', '左肩', ' ', ' ', '使剑', '少年', '不待', '剑', '招用', '老', ' ', ' ', '腕抖', '剑', '斜', ' ', ' ', '剑锋', '已削', '向', '那', '汉子', '右颈']
['那', '中年', '汉子', '竖剑', '挡格', ' ', ' ', '铮', '的', '一', '声响', ' ', ' ', '双剑相', '击', ' ', ' ', '嗡嗡', '做声', ' ', ' ', '震声', '未绝', ' ', ' ', '双刃剑', '光', '霍霍', ' ', ' ', '已', '拆', '了', '三招']
['中年', '汉子', '长剑', '猛地', '击落', ' ', ' ', '直', '斩', '少年', '顶门']
['那', '少年', '避', '向', '右侧', ' ', ' ', '左手', '剑诀', '斜引', ' ', ' ', '青钢剑', '疾刺', '那', '汉子', '大腿']
['n', 'u3000', 'u3000', '两人', '剑法', '迅捷', ' ', ' ', '全力', '相搏']
['n', 'u3000', 'u3000', '练武', '厅', '东边', '坐', '着', '二人']
['上首', '是', '个', '四十左右', '的', '中年', '道姑', ' ', ' ', '铁青', '着', '脸', ' ', ' ', '嘴唇', '紧闭']
['下首', '是', '个', '五十余岁', '的', '老者', ' ', ' ', '右手', '捻', '着', '长须', ' ', ' ', '神情', '甚', '是', '得意']
['两人', '的', ' ', '座位', '相距', '一丈', '有余', ' ', ' ', '身后', '各站', '着', '二十余名', '男女', '弟子']
['西边', '一排', '椅子', '上', '坐', '着', '十余位', '宾客']

(6) 统计并输出前 30 个高频词和对应出现次数

```
62 # (6)
63 fdist = FreqDist(chinese_str)
64 print(fdist.most_common(30))
```

[('不', 19692), ('道', 15684), ('一', 14412), ('人', 12617), ('大', 9519), ('说', 7124), ('中', 6936), ('来', 6656), ('子', 6612), ('之', 6566), ('上', 6529), ('下', 6263), ('段', 5733), ('手', 5465), ('去', 5151), ('心', 5021), ('出', 5019), ('便', 5016), ('见', 4612), ('声', 4409), ('是', 4303), ('身', 4228), ('到', 3742), ('只', 3714), ('誉', 3658), ('得', 3605), ('好', 3577), ('无', 3546), ('头', 3499), ('老', 3457)]

(7) 与同你选择同样一本小说的同学, 对比前述统计结果, 说明结果相

同或者不同的原因。

不同的原因：

- 1.增添的特殊停用词不同
- 2.处理数据的时候选择的是否是只含有中文词的文本
- 3.jieba 分词的时候选用的模式不同

五、指导教师评语

成绩：

批阅人：

日期：