



上海外国语大学
SHANGHAI INTERNATIONAL STUDIES UNIVERSITY

Computational Journalism

Lecture 7: Big Data Analysis and Visualization
Ting Wang

Outlines

- Word Cloud using Python
- Chinese Word Segmentation
- Chinese Word Cloud



data visualization with word cloud

Word Cloud using Python

Word Cloud using Python

Case Description

Motivations:

- To measure a news objectively
- To obtain new information efficiently

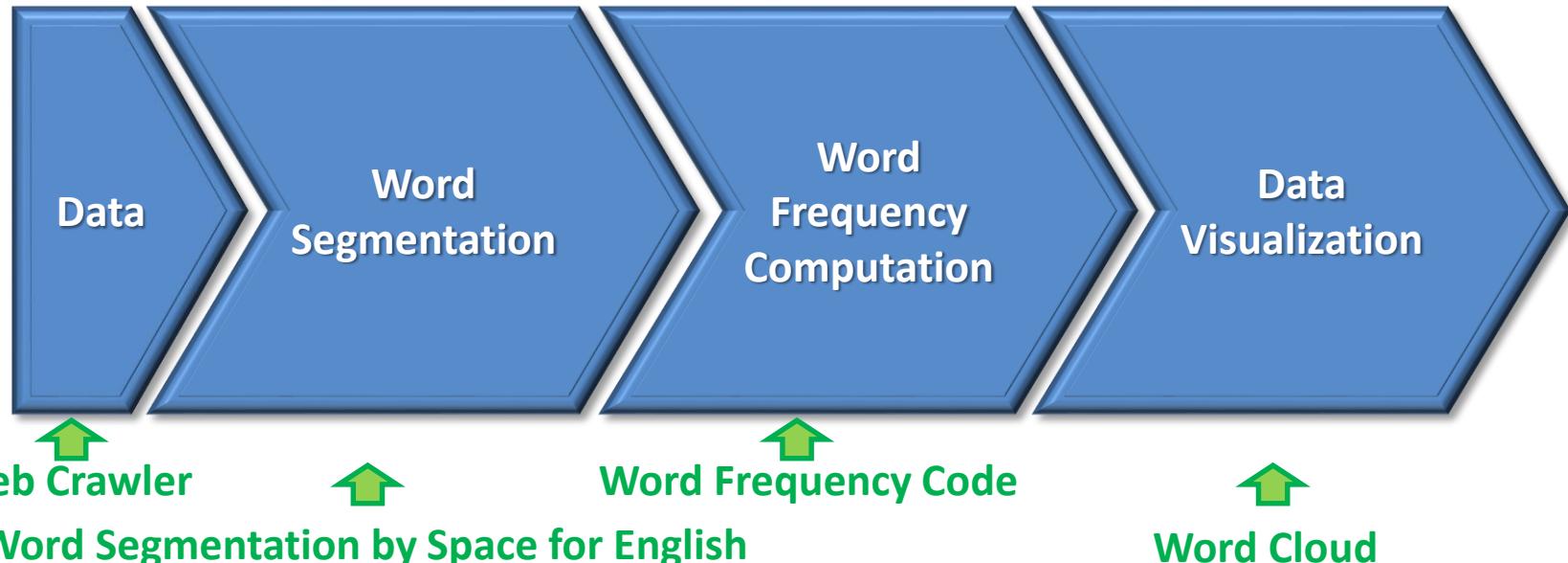


Methodologies:

- Describe a news report by quantitative method
- Technical integration by computer science, statistics and journalism

Word Cloud using Python

Now, we have data, how to mining it?



Word Cloud using Python

EXAMPLE 1:
Word Cloud for English News



Word Cloud using Python

Preparation: Take a news for an example.

<https://edition.cnn.com/2018/12/22/politics/shutdown-mattis-whitaker-trump/index.html>

① https://edition.cnn.com/2018/12/22/politics/shutdown-mattis-whitaker-trump/index.html

CNN politics 45 CONGRESS SUPREME COURT 2018 ELECTION RESULTS f t i S

The Washington nightmare before Christmas: a government in chaos

Analysis by Stephen Collinson, CNN
Updated 2209 GMT (0609 HKT) December 22, 2018



BREAKING NEWS
FEDERAL GOVERNMENT PARTIALLY SHUTDOWN AT MIDNIGHT

12:00 AM ET

NEWS & BUZZ

Lawmaker calls Trump official 'liar,' storms out

Bloomberg 'is not going to be intimidated' by Trump, says his...

Ad closed by Google

Word Cloud using Python

- Step 1-1: Build a Web Crawler, and extract the text from this page.

```
import urllib.request
response = urllib.request.urlopen('https://edition.cnn.com/2018/12/22/politics/shutdown-mattis-whitaker-trump/index.html')
HTMLText = response.read()

print(HTMLText)
```

The screenshot shows a Jupyter Notebook interface with a single code cell containing the provided Python code. The code uses the `urllib` module to open a URL and read the content. The output of the code is a large block of HTML text, which is displayed in the notebook's output area. The notebook has tabs for 'Python Console', 'Terminal', 'Run', 'TODO', and 'Event Log'.

```
C:\Program Files\Python36\python.exe "D:/SISU/Courses/2018计算机新闻学（研）/Lecture 7/EnglishWordCloud.py"
... (The output is a very long string of HTML code, likely the content of the CNN news article about the shutdown.)
```

Word Cloud using Python

- Step 1-2: Build a Web Crawler, and extract the text from this page.

```
import urllib.request
from bs4 import BeautifulSoup

response = urllib.request.urlopen('https://edition.cnn.com/2018/12/22/politics/shutdown-mattis-whitaker-trump/index.html')
HTMLText = response.read()
BSobj = BeautifulSoup(HTMLText, "html.parser")
Content = BSobj.find("section", {"class": "zn zn-body-text zn-body zn--idx-0 zn--ordinary zn-has-multiple-containers zn-has-48-containers"})

print(Content.get_text())
```

The screenshot shows a CNN news article titled "FEDERAL GOVERNMENT PARTIALLY SHUTDOWN AT MIDNIGHT". The browser's developer tools are open, with the "Elements" tab selected. The "body-text" section of the page is highlighted, and its class attributes are visible: "zn zn-body-text zn-body zn--idx-0 zn--ordinary zn-has-multiple-containers zn-has-48-containers". The main content of the page discusses the political situation and the shutdown.

EnglishWordCloud2

C:\Program Files\Python36\python.exe "D:/SISU/Courses/2018计算新闻学(研)/Lecture 7/EnglishWordCloud2.py"

(CNN) - President Donald Trump is precipitating chaos and seeking to wield unrestrained power as America enters a holiday period overshadowed by political pandemonium orchestrated by the disruptor-in-chief. For the third time this year, Congress is paralyzed, unable to prevent a shutdown that sent thousands of federal employees home for Christmas unsure about their upcoming paychecks. Trump is polling advisers on whether he has the power to fire Federal Reserve Chairman Jerome Powell following sell-offs on Wall Street that have taken away one of his favorite measures of his own job performance -- soaring stock markets. The revelation came days after the President announced a snap withdrawal of US troops in Syria against the advice of his advisers and without consulting allies. The move provoked the resignations of his most admired Cabinet officer, Defense Secretary James Mattis, who penned a devastating critique of Trump's "America First" world view, and a day later, of Trump's special envoy in the ISIS fight, Brett McGurk. A surprise announcement that Justice Ruth Bader Ginsburg underwent surgery for cancerous growths on her lung added a frenetic mood in Washington, as the Supreme Court dealt a blow to Trump by knocking back his new restrictions on asylum seekers who cross the southern border. Read MoreThe sense of things slipping out of control on multiple levels in the political world more unsettled and on edge than at any other time in Trump's tumultuous presidency. Even Republicans who have rarely dared to cross the President fumed that Trump appeared to navigate himself into a no-win situation with the government shutdown, and border wall funding prospects that will only worsen once Democrats take control of the House in a couple of weeks. "We are pretty much flying here without an instruction book," said Republican Sen. Roy Blunt. Trump, who bowed to a right-wing revolt and forced the fight by digging in on a dispute over funding for his border wall, addressed the crisis by tweeting a picture of himself signing already-passed legislation that included the naming of post offices -- while also complaining that he was staying in Washington instead of heading out on his 16-day Florida golf vacation as planned. That was after one senator, Democrat Brian Schatz, revealed a holiday period overshadowed by political pandemonium orchestrated by the disruptor-in-chief. For the third time this year, Congress is paralyzed, unable to prevent a shutdown that sent thousands of federal employees home for Christmas unsure about their upcoming paychecks. Trump is calling advisers on whether he has the power to fire Federal Reserve Chairman Jerome Powell following sell-offs on Wall Street that have taken away one of his favorite measures of his own job performance -- soaring stock markets.

上海外国语大学
SHANGHAI INTERNATIONAL STUDIES UNIVERSITY

Word Cloud using Python

Step2: Do you remember the word frequency code in Lecture 3? Revise it to a function.

```
with open('text_flu.txt', 'rU') as f:  
    s_array = str(f.read())
```

INPUT

```
s_array = s_array.replace(',', ' ')  
s_array = s_array.replace('.', ' ')  
s_array = s_array.replace('\'', ' ')  
s_array = s_array.replace('\"', ' ')  
s_array = s_array.replace(';', ' ')  
s_array = s_array.replace(';', ' ')  
s_array = s_array.replace('?', ' ')  
s_array = s_array.replace('\'', ' ')  
s_array = s_array.replace('\"', ' ')  
s_array = s_array.replace('!', ' ')  
s_array = s_array.replace('@', ' ')  
s_array = s_array.replace('#', ' ')  
s_array = s_array.replace('%', ' ')  
s_array = s_array.replace('$', ' ')  
s_array = s_array.replace('\"', ' ')  
s_array = s_array.replace('&', ' ')  
s_array = s_array.replace('*', ' ')  
s_array = s_array.replace('(', ' ')  
s_array = s_array.replace(')', ' ')  
s_array = s_array.replace('\'', ' ')  
s_array = s_array.replace('\"', ' ')
```

```
s_array = s_array.replace('_', ' ')  
s_array = s_array.replace('=', ' ')  
s_array = s_array.replace('{', ' ')  
s_array = s_array.replace('}', ' ')  
s_array = s_array.replace('[', ' ')  
s_array = s_array.replace(']', ' ')  
s_array = s_array.replace('|', ' ')  
s_array = s_array.replace('\'', ' ')  
s_array = s_array.replace('<', ' ')  
s_array = s_array.replace('>', ' ')  
s_array = s_array.split()
```

OUTPUT

```
word_dict = {}  
  
for i in range(len(s_array)):  
    if s_array[i] not in word_dict:  
        word_dict[s_array[i]] = 1  
    else:  
        word_dict[s_array[i]] += 1  
  
with open('result_flu_test.txt', 'w') as f:  
    for word, number in word_dict.items():  
        result = word + "\t" + str(number) + "\n"  
        f.write(result)
```

Make sure that you know what is input, and what is output.



Word Cloud using Python

- Revised function

```
def WordFrequency(Text):  
    s_array = str(Text)  
    s_array = s_array.replace('，',',')  
    s_array = s_array.replace('。',',')  
    s_array = s_array.replace('\'',',')  
    s_array = s_array.replace('\'',',')  
    s_array = s_array.replace('：',',')  
    s_array = s_array.replace('；',',')  
    s_array = s_array.replace('？',',')  
    s_array = s_array.replace('～',',')  
    s_array = s_array.replace('！',',')  
    s_array = s_array.replace('@',',')  
    s_array = s_array.replace('#',',')  
    s_array = s_array.replace('%',',')  
    s_array = s_array.replace('$',',')  
    s_array = s_array.replace('^',',')
```



啦啦啦啦啦啦啦啦

You still have to revise one place in your main code,
but I do NOT want to tell you,
could you try it by yourself?

```
s_array = s_array.replace('～', ',')  
s_array = s_array.replace('&', ',')  
ay = s_array.replace('*', ',')  
ay = s_array.replace('(', ',')  
ay = s_array.replace(')', ',')  
ay = s_array.replace('-', ',')  
ay = s_array.replace('+', ',')  
ay = s_array.replace('_', ',')  
ay = s_array.replace('=', ',')  
ay = s_array.replace('{', ',')  
ay = s_array.replace('}', ',')  
ay = s_array.replace('[', ',')  
ay = s_array.replace(']', ',')  
ay = s_array.replace('|', ',')  
ay = s_array.replace('\\\\', ',')  
ay = s_array.replace('<', ',')  
ay = s_array.replace('>', ',')  
ay = s_array.split()  
dict = {}  
    in range(len(s_array)):  
if s_array[i] not in word_dict:  
    word_dict[s_array[i]] = 1  
else:  
    word_dict[s_array[i]] += 1  
return word_dict
```

Word Cloud using Python

Results

Run: EnglishWordCloud3 x

```
C:\Program Files\Python36\python.exe "D:/SISU/Courses/2018计算新闻学（研）/Lecture 7/EnglishWordCloud3.py"
```

The word cloud visualization is displayed in the main window, showing the most frequent words from the processed text. The words are represented by their size and color, indicating frequency and possibly sentiment.

Python Console Terminal Run TODO

Event Log



上海外国语大学

SHANGHAI INTERNATIONAL STUDIES UNIVERSITY

Word Cloud using Python

Before the most important step:

1. *Build a folder called ‘static’ in your computer.*
2. *Install some packages for word cloud.*



Word Cloud using Python

Data Visualization using Python

- Necessity:
 - NumPy (Computing Package)
 - Scipy (Scientific Computing Package)
 - Pillow(Image)
 - Matplotlib (Diagram Package)
 - wordcloud (Word Cloud Package)
- Some packages also need some other required packages

Installation Sequence



Word Cloud using Python

Step 4: employ the word cloud code.

```
1 import urllib.request
2 from bs4 import BeautifulSoup
3 from DataWordCloud import word_cloud_generate
```

```
50 response = urllib.request.urlopen('https://edition.cnn.com/2018/12/22/politics/shutdown-mattis-whitaker-trump/index.html')
51 HTMLText = response.read()
52 BSobj = BeautifulSoup(HTMLText, "html.parser")
53 Content = BSobj.find("section", {"class": "zn zn-body-text zn-body zn--idx-0 zn--ordinary zn-has-multiple-containers zn-has-48-containers"})
54
55 # print(WordFrequency(Content.get_text()))
56
57 picurl=word_cloud_generate(WordFrequency(Content.get_text()))
58
59 print("word cloud saved in "+picurl)
```



Word Cloud using Python

Results:

Disadvantages:

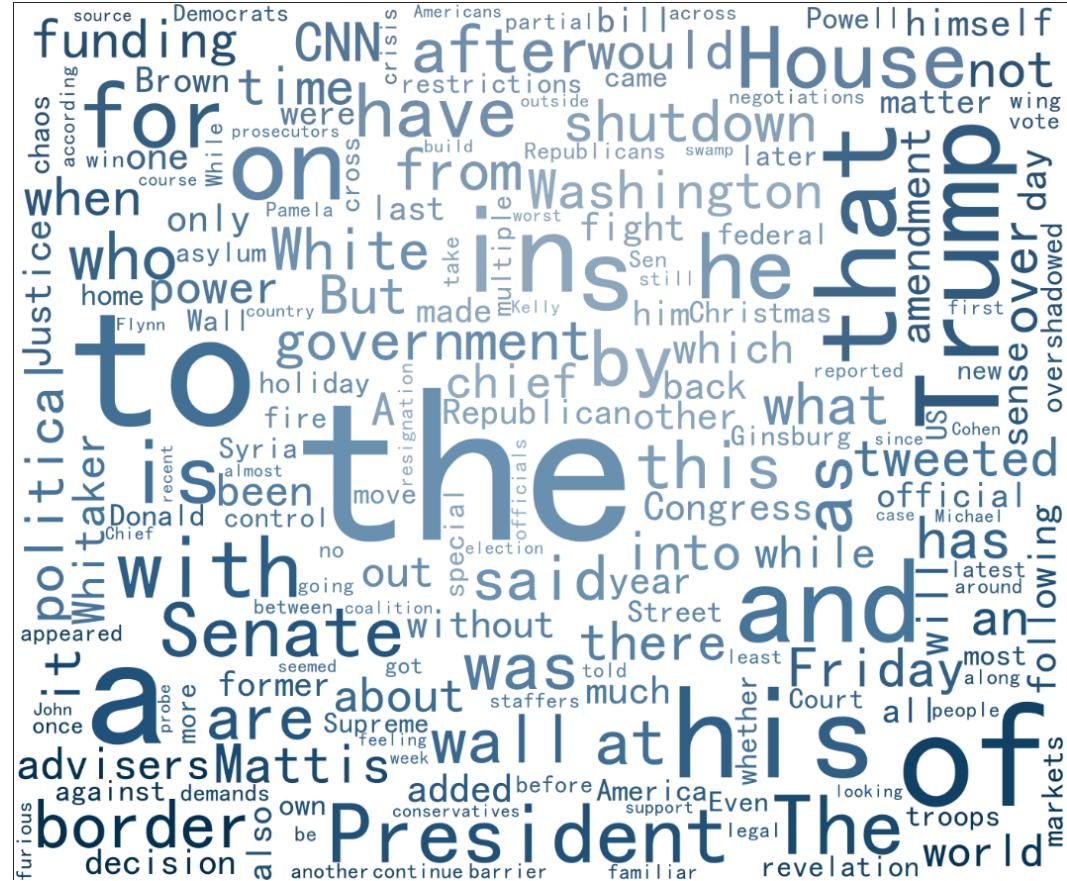
1. Too many words

→to set a threshold

2. Lots of unmeaning words

→to employ a list for stop words

<http://www.lextek.com/manuals/onix/stopwords1.html>



Word Cloud using Python

Optimization

Get the Word Dict

Loop words in Word Dict:

if word threshold > n:

 if word not a stop word:

 save it in Selected Word Dict

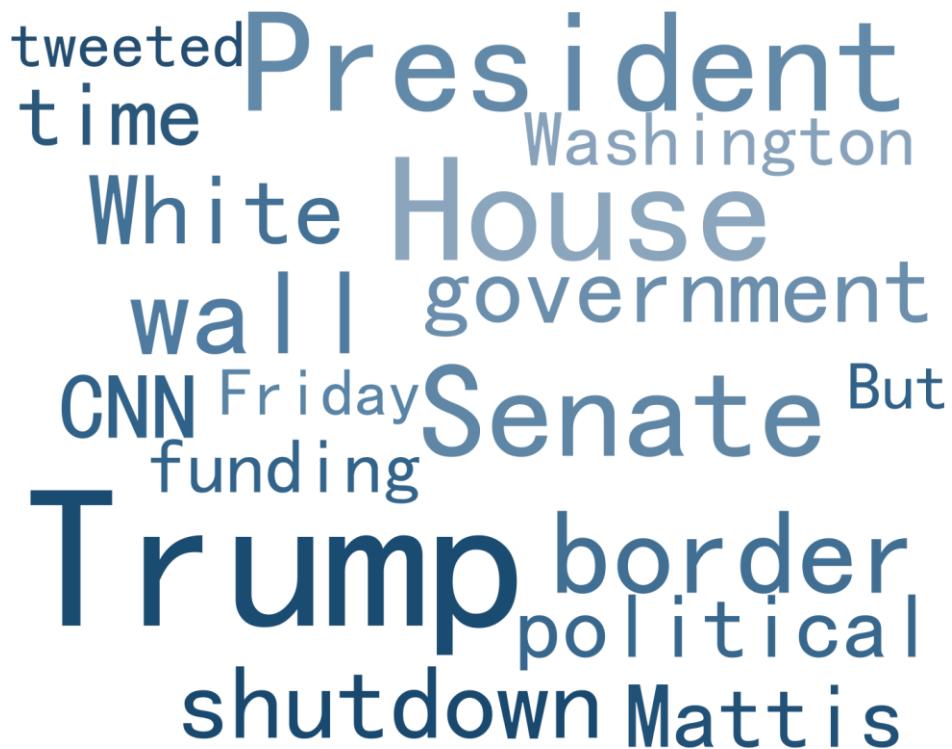
 else:

 continue

else:

 continue

**Know the content without reading:
President Trump, White House, and
Washington Government shutdown
this Friday!**



Word Cloud using Python



Ask A Question



We use ‘ ’(a space) to divide words one by one in English, how to segment words in Chinese or Japanese?





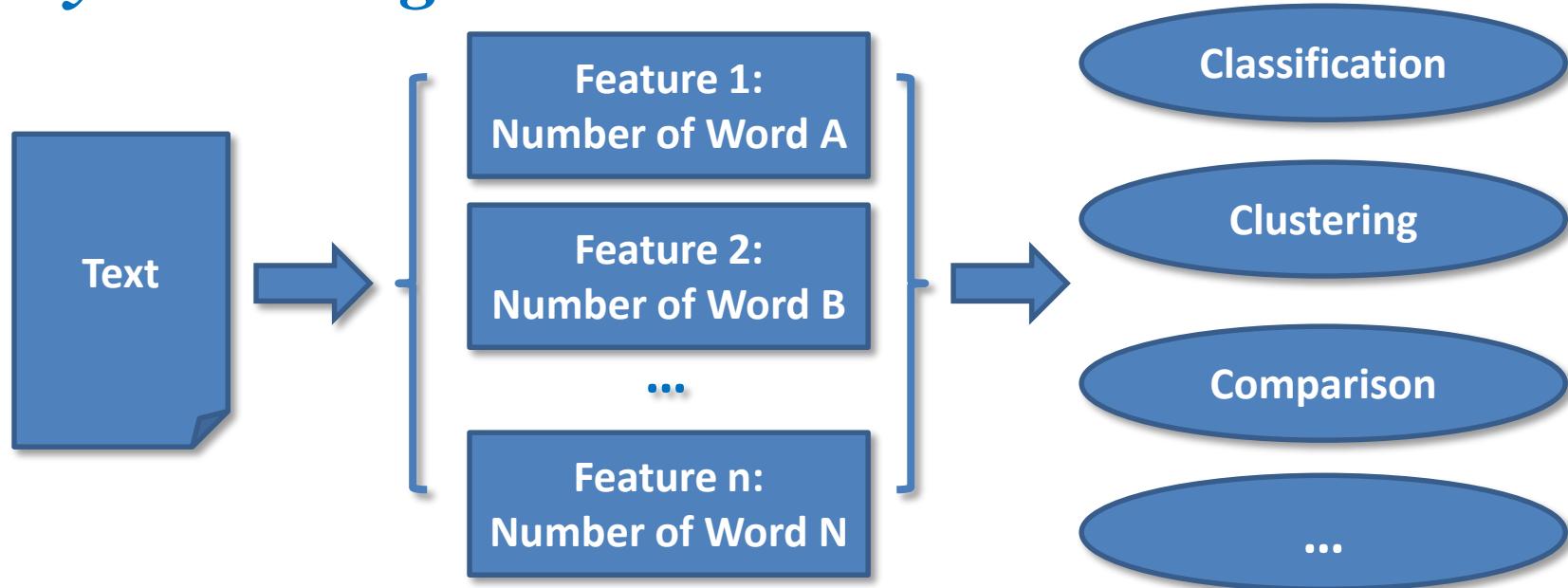
上海外国语大学
SHANGHAI INTERNATIONAL STUDIES UNIVERSITY

the first step for Chinese information processing

Chinese Word Segmentation

Chinese Word Segmentation

Why Word Segmentation?



However, it is difficult to extract words from Chinese text.

Chinese Word Segmentation

Difficulties: Disambiguation

乒乓球拍卖完了

乒乓|球拍|卖完了

乒乓球|拍卖|完了

一脸懵逼



- Chinese Word Segmentation 分词
 - Forward Max. matching method 正向最大匹配
 - Backward Max. matching method 逆向最大匹配
 - Statistical matching method 统计学方法

Chinese Word Segmentation

Forward Max. matching method, FMM

正向最大匹配

准备工作：需要分词词典D

设MaxLen表示最大词长度

算法：

1. 从生语料N中取长度为MaxLen的字串str, 令Len= MaxLen
2. 把str与D中的词相匹配
3. 若匹配成功, 则认为str为词, N中去掉str(指针前移Len个单位), 返回1
4. 若匹配不成功,
 - ◆ 若Len>1则Len--, 从生语料N中取长度为Len的字串str返回2;
 - ◆ 否则, 得到单字词, N中去掉str(指针前移1个单位), 返回1

若4中得到的单字不是词, 则要进行未登录词处理

若待切分的语料字串长度小于MaxLen, 则取str为待切分语料



Chinese Word Segmentation

Backward Max. matching method, BMM

逆向最大匹配

1. Similar to FMM, but the text is scanned from the right side
2. Often jointly use with FMM

Chinese Word Segmentation

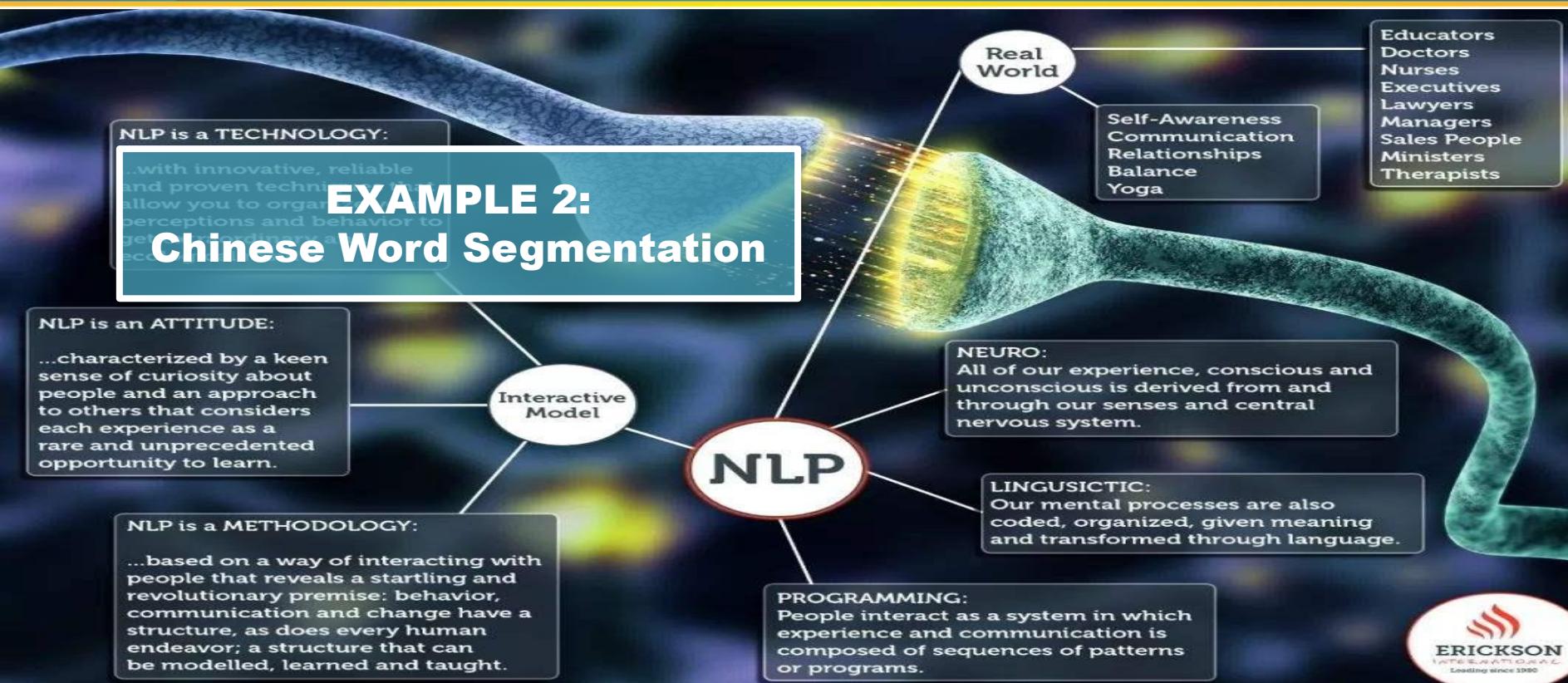
• Statistical Matching Method

FMM and BMM

```
Begin initialize Path←{}, AmbiguousString, SubString←{}  
    While (AmbiguousString.Length>0)  
    {  
        //只考虑以当前HMM第一个状态开始的匹配序列  
        SubString←以AmbiguousString中的第一个字为基准，取出所有可能的匹配字符串  
        Foreach SubString  
        {  
            //提供当前情况下所有的概率，为判断歧义作参考  
            计算当前每一种可能情况的概率P(SubString) //unigram, bigram, trigram with smoothing  
        }  
        //选择概率最大的SubString添加到Path  
        将argmax(P(SubString))添加到Path  
        //准备考察除去最大概率的SubString后的AmbiguousString，从HMM序列首部开始，除去所有的匹配状态  
        AmbiguousString.Remove(0, argmax(P(SubString)).Length)  
    }  
    Return Path  
End
```



Chinese Word Segmentation





新华每日电讯

4.2万 文章 2.2亿 总阅读

查看TA的文章>

2



分享到



创新是改革开放的生命

2018-12-23 01:08

新华社评论员“创新是改革开放的生命。”习近平总书记在庆祝改革开放40周年大会重要讲话中作出这一重要论断,深刻指明创新对于推进改革开放事业的极端重要性,为进一步解放思想、开拓创新注入了强大精神动力。

“苟日新,日日新,又日新”,创新是一个民族进步的灵魂,是一个国家兴旺发达的不竭动力,也是中华民族最深沉的民族禀赋。改革开放40年来,党带领人民大胆地试、勇敢地改,闯出了一片新天地,探索出前无古人的中国特色社会主义道路。从实行家庭联产承包到农村承包地“三权”分置,从兴办经济特区到设立自贸区,改革开放事业就是在不断解放思想、开拓创新中从无到有发展起来的。

创新,是破局开路的利器;创新,是点亮未来的希望。没有创新就没有中国的今天,也就没有中国的明天。推进新时代的改革开放,必须把创新置于更加突出位置,不断推进理论创新、制度创新、科技创新、文化创新等各方面创新,让创新在全社会蔚然成风。

创新源泉的涌流,来自思想解放的破冰。“真理标准大讨论”如一声春雷,打破“两个凡是”的禁锢;建立社会主义市场经济体制的提出,解开束缚人们思想和行动多年的绳索。40年实践历程启示我们,没有思想的大解放,就不会有创新活力的大迸发,就不会有改革开放的大突破。实践发展永无止境,解放思想永无止境。新征程上,面对新形势、新任务、新挑战,更加需要打开

Elements Console Sources

```
<!doctype html>
<html data-log-pv="{"mpc":46}" style="font-size: 79px;">
  <head>...</head>
  <body class="article-page">
    <div class="wrapper-box">
      <header id="main-header" class="error-head">...</header>
      <div class="location area">...</div>
      <div class="area clearfix" id="article-container">
        <div class="column left">...</div>
        <div class="left main">
          <div class="text">
            <div class="text-title">...</div>
            <article class="article" id="mp-editor" style="display: none">
              <p data-role="original-title" style="display: none">原标题: 创新是改革开放的生命</p>
              <p>...</p>
              <p>...</p>
              <p>...</p>
              <p>...</p>
              <p>当今世界,变革创新的潮流滚滚向前。谁排斥变革,谁拒绝创新,谁就会落后于时代,谁就会被历史淘汰。从这个意义上说,抓创新就是抓发展,谋创新就是谋未来。</p>
              <p>...</p>
              <p>惟改革者进,惟创新者强,惟改革创新者胜。新时代的改革开放,呼唤我们弘扬创新精神、汇聚创新力量、推动创新发展,</p>
            </article>
          </div>
        </div>
      </div>
    </div>
  </body>

```

html body div #article-container div div article#mp-editor.article

Styles Event Listeners DOM Breakpoints Properties Accessibility

Filter :hover .cls +



Chinese Word Segmentation

Step 1: get the text

```
import urllib.request  
from bs4 import BeautifulSoup  
  
response = urllib.request.urlopen('http://www.sohu.com/a/283822245_117503')  
HTMLText = response.read()  
BSobj = BeautifulSoup(HTMLText, "html.parser")  
Content = BSobj.find("article", {"class": "article"})  
print(Content.get_text())
```

Run: ChineseWordSegmentation01 ×

"C:\Program Files\Python36\python.exe" "D:/SISU/Courses/2018计算新闻学（研）/Lecture 7/ChineseWordSegmentation01.py"

原标题：创新是改革开放的生命
新华社评论员“创新是改革开放的生命。”习近平总书记在庆祝改革开放40周年大会重要讲话中作出的这一重要论断，深刻指明创新对于推进改革开放事业的极端重要性，为进一步解放思想、开拓创新注入了强大精神动力。
“苟日新、日日新、又日新”，创新是一个民族进步的灵魂，是一个国家兴旺发达的不竭动力，也是中华民族最深沉的民族禀赋。改革开放40年来，党带领人民大胆地试、勇敢地改，闯出了一片新天地，探索出前无古人的中国特色社会主义道路。从实行家庭联产承包到农村承包地“三权”分置，从兴办经济特区到设立自贸区，改革开放事业就是在不断解放思想、开拓创新中从无到有发展起来的。
创新，是破局开路的利器；创新，是点亮未来的希望。没有创新就没有中国的今天，也就没有中国的明天。推进新时代的改革开放，必须把创新置于更加突出位置，不断推进理论创新、制度创新、科技创新、文化创新等各方面创新，让创新在全社会蔚然成风。
创新源泉的涌流，来自思想解放的破冰。“真理标准大讨论”如一声春雷，打破“两个凡是”的禁锢；建立社会主义市场经济体制的提出，解开束缚人们思想和行动多年的绳索。40年实践历程启示我们，没有思想的大解放，就不会有创新活力的大迸发，就不会有改革开放的大突破。实践发展永无止境，解放思想永无止境。新征程上，面对新形势、新任务、新挑战，更加需要打开解放思想这个“总开关”，继续鼓起闯的勇气、迈开试的步子，及时回答时代之问、人民之问，廓清困扰和束缚实践发展的思想迷雾，不断推动实践基础上的理论创新，推动改革开放取得新成就。
当今世界，变革创新的潮流滚滚向前。谁排斥变革，谁拒绝创新，谁就会落后于时代，谁就会被历史淘汰。从这个意义上说，抓创新就是抓发展，谋创新就是谋未来。
“问渠那得清如许？为有源头活水来。”激发创新活力，要靠全面深化改革、扩大开放；随着改革开放不断向前推进，更强劲的创新动能必将喷涌而出。这是一个相互促进、相得益彰的过程。要通过深化各领域改革，破除一切制约创新的思想障碍和制度藩篱，在全社会大力营造勇于创新、鼓励成功、宽容失败的良好氛围，培育并用好各类创新人才，进一步激发全社会创新活力和创造潜能。
惟改革者进，惟创新者强，惟改革创新者胜。新时代的改革开放，呼唤我们弘扬创新精神、汇聚创新力量、推动创新发展，在更加广袤的历史时空创造新的更大奇迹。
新华社北京12月22日电返回搜狐，查看更多 责任编辑：

Chinese Word Segmentation

Step2: Cut words using FMM

Step2-1: build the dictionary of Chinese Words

准备工作：需要分词词典D

设MaxLen表示最大词长度

算法：

1. 从生语料N中取长度为MaxLen的字串str, 令Len= MaxLen
2. 把str与D中的词相匹配
3. 若匹配成功, 则认为str为词, N中去掉str(指针前移Len个单位), 返回1
4. 若匹配不成功,
 - ◆ 若Len>1则Len--, 从生语料N中取长度为Len的字串str返回2;
 - ◆ 否则, 得到单字词, N中去掉str(指针前移1个单位), 返回1

若4中得到的单字不是词, 则要进行未登录词处理

若待切分的语料字串长度小于MaxLen, 则取str为待切分语料

Chinese Word Segmentation

Step2: Cut words using FMM

Step2-1: build the dictionary of Chinese Words

```
def word_map_initial(): #把词汇装入词典
    WordMap=set()
    conn = pymysql.connect(user='root', password='123456', database='nlp', charset="utf8")
    cursor = conn.cursor()
    sqlstr="SELECT DISTINCT WORD_NAME FROM NLP_WORD_MAP"
    cursor.execute(sqlstr)
    rows = cursor.fetchall()
    for row in rows:
        WordMap.add(str(row[0]))
    conn.commit()
    cursor.close()
    conn.close()
    return WordMap
```

Chinese Word Segmentation

Step2: Cut words using FMM

Step2-2: FMM Code

```
def word_seg_fmm(content): #正向匹配
    WordMap = word_map_initial()
    MaxLen=10 #最大词长
    Len=MaxLen #动态切割词长
    Seg_Content="" #返回的切割结果

    while len(content)>0:
        if content[0:Len] in WordMap: #词典中有匹配
            Seg_Content=Seg_Content+content[0:Len]+"|"
            content=content[Len:]
            Len=MaxLen
            #print ("Seg_Content1:" + Seg_Content)
            continue
        else: #词典中无匹配
            Len=Len-1
            if Len==1:#仅剩一个词还没匹配到
                Seg_Content = Seg_Content + content[0:Len] + " | "
                content = content[Len:]
                Len = MaxLen
                #print ("Seg_Content2:" + Seg_Content)
    return Seg_Content[:-1]
```

Chinese Word Segmentation

Step 3: Revise the head and the main body.

```
1 import pymysql
2 import urllib.request
3 from bs4 import BeautifulSoup
4
5 def word_map_initial():
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
```



Chinese Word Segmentation

Results:

```
"C:\Program Files\Python36\python.exe" "D:/SISU/Courses/2018计算新闻学(研)/Lecture 7/ChineseWordSegmentation.py"
```

原	标题	：	创新	是	改革	开放	的	生命																																																																																																																																																																																																																																										
	新华社评论员：“创新”是“改革”“开放”的“生命”。习总书记“庆祝改革开放40周年大会”重要讲话中作出的这一重要论断，深刻指明创新对于推进改革、开放事业的极端重要性，为“进一步解放思想”、“开拓创新”注入了强大精神动力。																																																																																																																																																																																																																																																	
“苟日新，日日新，又日新”，创新是一个民族进步的“灵魂”，是“一个国家兴旺发达的不竭动力”，也是“中华民族最深沉的民族禀赋”。改革开放40年来，党带领人民大胆地试、勇敢地改，“闯出了一片新天地”，探索出前无古人的中国特色社会主义道路。“从实行家庭联产承包到农村承包地‘三权’分置，从兴办经济特区到设立自贸区，改革开放事业就是“在不断解放思想”、“开拓创新”中“从无到有”发展起来的”。																																																																																																																																																																																																																																																		
创新	，	是	破	局	开	路	的	利器	；	创新	，	是	点	亮	未	来	的	希望	。	没有	创新	就	没有	中国	的	今天	，	也	就	没有	中国	的	明天	。	推进	新	时	代	的	改革	开放	，	必须	把	创新	置	于	更	加	突	出	位	置	，	不	断	推	进	理	论	创	新	、	制	度	创	新	、	科	技	创	新	、	文	化	创	新	等	各	方	面	创	新	，	让	创	新	在	全	社	会	蔚	然	成	风	。																																																																																																																																																		
创新	源泉	的	涌流	，	来	自	思想解放	的	破	冰	。	“真	理	标	准	大	讨	论	”如	一	声	春	雷	，	打	破	“两	个	凡	是	”	的	禁	锢	，	建	立	社	会	主	义	市	场	经	济	体	制	的	提	出	，	解	开	束	缚	人	们	思	想	和	行	动	多	年	的	绳	索	。	40	年	实	践	历	程	启	示	我	们	，	没	有	思	想	的	大	解	放	，	就	不	会	有	创	新	活	力	的	大	进	发	，	就	不	会	有	改	革	开	放	的	大	突	破	。	实	践	发	展	永	无	止	境	，	解	放	思	想	永	无	止	境	。	新	征	程	上	，	面	对	新	形	势	、	新	任	务	、	新	挑	战	，	更	加	需	要	打	开	解	放	思	想	这	个	“总	开关	”	，	继	续	鼓	起	闯	的	勇	气	、	迈	开	试	的	步	子	，	及	时	回	答	时	代	之	问	、	人	民	之	问	，	廓	清	困	扰	和	束	缚	实	践	发	展	的	思	想	迷	雾	，	不	断	推	动	实	践	基	础	上	的	理	论	创	新	，	推	动	改	革	开	放	取	得	新	成	就	。	
当	今	世	界	，	变	革	创	新	的	潮	流	滚	滚	向	前	。	谁	排	斥	变	革	，	谁	拒	绝	创	新	，	谁	就	会	落	后	于	时	代	，	谁	就	会	被	历	史	淘	汰	。	从	这	个	意	义	上	说	，	抓	创	新	就	是	抓	发	展	，	谋	创	新	就	是	谋	未	来	。																																																																																																																																																																										
“问	渠	那	得	清	如	许	？	为	有	源	头	活	水	来	。	”	激	发	创	新	活	力	，	要	靠	全	面	深	化	改	革	、	扩	大	开	放	，	随	着	改	革	开	放	不	断	向	前	推	进	，	更	强	劲	的	创	新	动	能	必	将	喷	涌	而	出	。	这	是	一	个	相	互	促	进	、	相	得	益	彰	的	过	程	。	要	通	过	深	化	各	领	域	改	革	，	破	除	一	切	制	约	创	新	的	思	想	障	碍	和	制	度	藩	篱	，	在	全	社	会	大	力	营	造	勇	于	创	新	、	鼓	励	成	功	、	宽	容	失	败	的	良	好	氛	围	，	培	育	并	用	好	各	类	创	新	人	才	，	进	一步	激	发	全	社	会	创	新	活	力	和	创	造	潜	能	。																																																																									
惟	改	革	者	进	，	惟	创	新	者	强	，	惟	改	革	创	新	者	胜	。	新	时	代	的	改	革	开	放	，	呼	唤	我	们	弘	扬	创	新	精	神	、	汇	聚	创	新	力	量	、	推	动	创	新	发	展	，	在	加	更	广	袤	的	历	史	时	空	创	造	新	的	更	大	奇	迹	。																																																																																																																																																																										
新华社	北京	1	2	月	2	2	日	电	返	回	搜	狐	，	查	看	更	多		责	任	编	辑	：																																																																																																																																																																																																																											

Process finished with exit code 0



Chinese Word Segmentation

Chinese Word Segmentation Algorithms

- FMM
- BMM



```
def word_seg_fmm(content): #正向匹配
    MaxLen=10 #最大词长
    Len=MaxLen #动态切割词长
    Seg_Content="" #返回的切割结果

    while len(content)>0:
        if content[0:Len] in WordMap: #词典中有匹配
            Seg_Content=Seg_Content+content[0:Len]+"|"
            content=content[Len:]
            Len=MaxLen
            #print("Seg_Content1:"+Seg_Content)
            continue
        else: #词典中无匹配
            Len=Len-1
            if Len==1:#仅剩一个词还没匹配到
                Seg_Content = Seg_Content + content[0:Len] + "|"
                content = content[Len:]
                Len = MaxLen
                #print("Seg_Content2:" + Seg_Content)
    return Seg_Content[:-1]
```

```
def word_seg_bmm(content): #逆向匹配
    MaxLen=10 #最大词长
    Len=MaxLen #动态切割词长
    Seg_Content="" #返回的切割结果

    while len(content)>0:
        if content[-Len:] in WordMap: #词典中有匹配
            Seg_Content=content[-Len:]+ "|" + Seg_Content
            content=content[:-Len]
            Len=MaxLen
            #print("Seg_Content1:"+Seg_Content)
            continue
        else: #词典中无匹配
            Len=Len-1
            if Len==1:#仅剩一个词还没匹配到
                Seg_Content = content[-Len:] + "|" + Seg_Content
                content = content[:-Len]
                Len = MaxLen
                #print("Seg_Content2:" + Seg_Content)
    return Seg_Content[:-1]
```

Chinese Word Segmentation

Tips for Chinese Word Segmentation

- Initialization is very important
- Segment in the memory (not hard disk or data bases) to accelerate the segmentation speed
- Using “set” to store the dictionary, and “dict” for segmented words in Python
- For Tag Analysis, a precise word segmentation is unnecessary



上海外国语大学
SHANGHAI INTERNATIONAL STUDIES UNIVERSITY

data visualization with word cloud for Chinese text

Chinese Word Cloud

Chinese Word Cloud

Steps:

1. Download a news report
2. Word segmentation
3. Word tag extraction and statistical computing
4. Data visualization and news summarization



Chinese Word Cloud

Database Preparation

- Word Dictionary (required)
- Stop Word Dictionary (required)
- Dictionaries of Terms (optional)
- Word Chains (required if using N-gram)
- Part of Speech (optional)
- Word Sentiment (optional for Sentiment Analysis)



Chinese Word Cloud

Word Frequency Computing for Chinese

- str.split() for all tags
- Discarding One-Char tags
- Discarding Stop-Word tags
- Select tags whose term frequencies are larger than a threshold (for example >2)
- Other statistical computing

Chinese Word Cloud

Chinese Word Frequency

```
def WordFrequency(InputText):
    s_array = str(InputText)
    s_array = s_array.replace('，', ' ')
    s_array = s_array.replace('。', ' ')
    s_array = s_array.replace('\'', ' ')
    s_array = s_array.replace('\\"', ' ')
    s_array = s_array.replace('：“', ' ')
    s_array = s_array.replace('；', ' ')
    s_array = s_array.replace('？', ' ')
    s_array = s_array.replace('～', ' ')
    s_array = s_array.replace('~', ' ')
    s_array = s_array.replace('！', ' ')
    s_array = s_array.replace('@', ' ')
    s_array = s_array.replace('#', ' ')
    s_array = s_array.replace('%', ' ')
    s_array = s_array.replace('$', ' ')
    s_array = s_array.replace('“', ' ')
    s_array = s_array.replace('&', ' ')
    s_array = s_array.replace('*', ' ')
```

```
s_array = s_array.replace('‘', ' ')
s_array = s_array.replace('’', ' ')
s_array = s_array.replace('－', ' ')
s_array = s_array.replace('＋', ' ')
s_array = s_array.replace('＿', ' ')
s_array = s_array.replace('＝', ' ')
s_array = s_array.replace('{', ' ')
s_array = s_array.replace('}', ' ')
s_array = s_array.replace('[', ' ')
s_array = s_array.replace(']', ' ')
s_array = s_array.replace('‘', ' ')
s_array = s_array.replace('’', ' ')
s_array = s_array.replace('＼', ' ')
s_array = s_array.replace('＜', ' ')
s_array = s_array.replace('＞', ' ')
s_array = s_array.replace('。', ' ')
s_array = s_array.replace('“', ' ')
s_array = s_array.replace('”', ' ')
```

```
s_array = s_array.replace('；', ' ')
s_array = s_array.replace('，', ' ')
s_array = s_array.replace('——', ' ')
s_array = s_array.replace('.....', ' ')
s_array = s_array.replace('？', ' ')
s_array = s_array.replace('！', ' ')
s_array = s_array.replace('￥', ' ')
s_array = s_array.replace('、', ' ')
s_array = s_array.replace('《', ' ')
s_array = s_array.replace('》', ' ')
s_array = s_array.replace('\u3000', ' ')
s_array = s_array.replace('\n', ' ')
s_array = s_array.split(' ')
word_dict = {}
for i in range(len(s_array)):
    if s_array[i] not in word_dict:
        word_dict[s_array[i]] = 1
    else:
        word_dict[s_array[i]] += 1
return word_dict
```

Chinese Word Cloud

Results:

解放思想
不断改革
实践
创新
没有
思想
发展
开放
推进

```
1 import urllib.request
2 from bs4 import BeautifulSoup
3 import DataWordCloud
4 from ChineseWordSeg import word_seg_fmm,WordFilter
5
6 def WordFrequency(InputText):
7
8     response = urllib.request.urlopen('http://www.sohu.com/a/283822245_117503')
9     HTMLText = response.read()
10    BObj = BeautifulSoup(HTMLText, "html.parser")
11    Content = BObj.find("article", {"class": "article"})
12
13    # print(WordFrequency(Content.get_text()))
14
15    WordSegResult = word_seg_fmm(Content.get_text())
16
17    WordDict = WordFrequency(WordSegResult)
18
19    SelectedWordDict={}
20    SelectedWordDict = WordFilter(WordDict)
21
22    threshold = 4
23    for key, value in SelectedWordDict.items():
24        if value >= threshold:
25            continue
26        else:
27            SelectedWordDict.pop(key)
28
29    picurl=DataWordCloud.word_cloud_generate(SelectedWordDict)
30
31    print("word cloud saved in "+picurl)
```





上海外国语大学
SHANGHAI INTERNATIONAL STUDIES UNIVERSITY

Tips for your Final Project

Tips for your Final Project

Word seg for NEWS_CONTENT

Word Segmentation

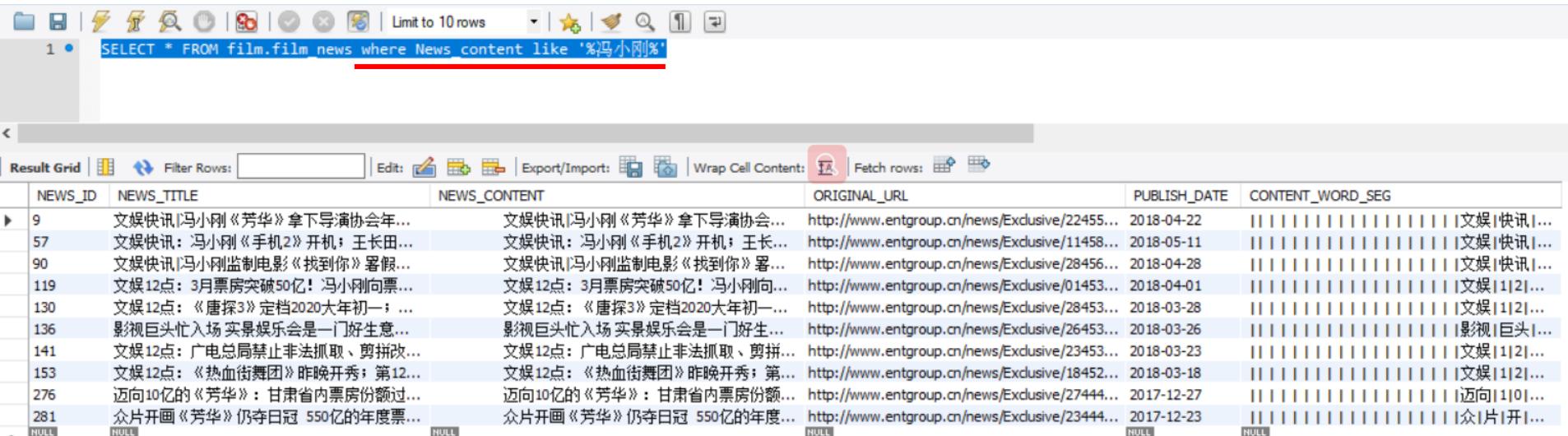
NEWS_ID	NEWS_TITLE	NEWS_CONTENT	ORIGINAL_URL	PUBLISH_DATE	CONTENT_WORD_SEG
1	成立发行公司发布豪华片单 腾讯与阿里在...	成立发行公司发布豪华片单 腾讯...	http://www.entgroup.cn/news/Exclusive/24455...	2018-04-24	成立 发行 ...
2	从泛娱乐到新文创 腾讯的这次升级到底“升”...	从泛娱乐到新文创 腾讯的这次升...	http://www.entgroup.cn/news/Exclusive/24455...	2018-04-24	从 泛 娱乐 ...
3	文娱12点 喜剧片《误入江湖》定档2019大...	文娱12点 喜剧片《误入江湖》定...	http://www.entgroup.cn/news/Exclusive/24455...	2018-04-24	文娱 12 ...
4	艺恩IP中心推荐：小说《失婚情陷冷情大叔》...	艺恩IP中心推荐：小说《失婚情...	http://www.entgroup.cn/news/Exclusive/24455...	2018-04-24	艺 恩 IP ...
5	超600亿的新媒体版权市场 谁是这一波的红...	超600亿的新媒体版权市场 谁是...	http://www.entgroup.cn/news/Exclusive/24455...	2018-04-24	超 600 ...
6	艺恩IP中心推荐：小说《青春是一个人的兵》...	艺恩IP中心推荐：小说《青春是...	http://www.entgroup.cn/news/Exclusive/23455...	2018-04-23	艺 恩 IP ...
7	文娱12点：腾讯影业成立腾影发行公司；电...	文娱12点：腾讯影业成立腾影发...	http://www.entgroup.cn/news/Exclusive/23455...	2018-04-23	文娱 12 ...
8	优酷春秋集发百部新作，你有多少睡眠时间够...	优酷春秋集发百部新作，你有多...	http://www.entgroup.cn/news/Exclusive/23455...	2018-04-23	优 酷 春 ...
9	文娱快讯 冯小刚《芳华》拿下导演协会年...	文娱快讯 冯小刚《芳华》拿下导...	http://www.entgroup.cn/news/Exclusive/22455...	2018-04-22	文娱 快讯 ...
10	艺恩IP中心推荐：小说《无间枭雄》，警察...	艺恩IP中心推荐：小说《无间枭...	http://www.entgroup.cn/news/Exclusive/22455...	2018-04-22	艺 恩 IP ...
11	《最后的锦衣卫》导演杨振豪：在国际舞台...	《最后的锦衣卫》导演杨振豪：...	http://www.entgroup.cn/news/Exclusive/22455...	2018-04-22	《 最后 的 ...
12	文娱快讯 《画皮3》宣布启动电影 网剧；...	文娱快讯 《画皮3》宣布启动电...	http://www.entgroup.cn/news/Exclusive/21455...	2018-04-21	文娱 快讯 ...
13	艺恩IP中心推荐：小说《孤星满月》，人生...	艺恩IP中心推荐：小说《孤星满...	http://www.entgroup.cn/news/Exclusive/21455...	2018-04-21	艺 恩 IP ...
14	艺恩对话 Base FX副总裁谢宁：影视后期需...	艺恩对话 Base FX副总裁谢宁：...	http://www.entgroup.cn/news/Exclusive/21455...	2018-04-21	艺 恩 对话 ...
15	北影节 孙向辉：观众满意度调查引导市场...	北影节 孙向辉：观众满意度调...	http://www.entgroup.cn/news/Exclusive/20455...	2018-04-20	北影 节 ...
16	承华传媒“拿下”《极限特工4》合资公司内...	承华传媒“拿下”《极限特工4》合...	http://www.entgroup.cn/news/Exclusive/20455...	2018-04-20	承华 传媒 ...
17	文娱12点：《头号玩家》延至5月29日上映...	文娱12点：《头号玩家》延至5月...	http://www.entgroup.cn/news/Exclusive/20455...	2018-04-20	文娱 12 ...
18	淘票票公布2018战略规划 发布电影宣发...	淘票票公布2018战略规划 发布...	http://www.entgroup.cn/news/Exclusive/19455...	2018-04-19	淘 票 票 ...
19	张昭：“屌丝”转中产将成为中国电影产业升...	张昭：“屌丝”转中产将成为中国电...	http://www.entgroup.cn/news/Exclusive/19455...	2018-04-19	张 昭：“ ...
20	艺恩扣映由国内电影市场独热 分享艺恩新...	艺恩扣映由国内电影市场独热 分享...	http://www.entgroup.cn/news/Exclusive/19455...	2018-04-19	艺 恩 扣 ...



Tips for your Final Project

Content limitations for Word Cloud

– Condition: Where ...



The screenshot shows a MySQL database interface with the following details:

- Query Bar: `SELECT * FROM film.film_news where News_content like '%冯小刚%'`
- Result Grid Headers: NEWS_ID, NEWS_TITLE, NEWS_CONTENT, ORIGINAL_URL, PUBLISH_DATE, CONTENT_WORD_SEG
- Data Rows (partial):

NEWS_ID	NEWS_TITLE	NEWS_CONTENT	ORIGINAL_URL	PUBLISH_DATE	CONTENT_WORD_SEG
9	文娱快讯 冯小刚《芳华》拿下导演协会年...	文娱快讯 冯小刚《芳华》拿下导演协会...	http://www.entgroup.cn/news/Exclusive/22455...	2018-04-22	文 娱 快 讯 ...
57	文娱快讯：冯小刚《手机2》开机；王长田...	文娱快讯：冯小刚《手机2》开机；王长...	http://www.entgroup.cn/news/Exclusive/11458...	2018-05-11	文 娱 快 讯 ...
90	文娱快讯 冯小刚监制电影《找到你》暑假...	文娱快讯 冯小刚监制电影《找到你》署...	http://www.entgroup.cn/news/Exclusive/28456...	2018-04-28	文 娱 快 讯 ...
119	文娱12点：3月票房突破50亿！冯小刚向票...	文娱12点：3月票房突破50亿！冯小刚向票...	http://www.entgroup.cn/news/Exclusive/01453...	2018-04-01	文 娱 1 2 ...
130	文娱12点：《唐探3》定档2020大年初一；...	文娱12点：《唐探3》定档2020大年初一；...	http://www.entgroup.cn/news/Exclusive/28453...	2018-03-28	文 娱 1 2 ...
136	影视巨头忙入场 实景娱乐会是一门好生意...	影视巨头忙入场 实景娱乐会是一门好生...	http://www.entgroup.cn/news/Exclusive/26453...	2018-03-26	影 视 巨 头 ...
141	文娱12点：广电总局禁止非法抓取、剪拼改...	文娱12点：广电总局禁止非法抓取、剪拼...	http://www.entgroup.cn/news/Exclusive/23453...	2018-03-23	文 娱 1 2 ...
153	文娱12点：《热血街舞团》昨晚开秀；第12...	文娱12点：《热血街舞团》昨晚开秀；第...	http://www.entgroup.cn/news/Exclusive/18452...	2018-03-18	文 娱 1 2 ...
276	迈向10亿的《芳华》：甘肃省内票房份额过...	迈向10亿的《芳华》：甘肃省内票房份额...	http://www.entgroup.cn/news/Exclusive/27444...	2017-12-27	迈向 1 0 ...
281	众片开画《芳华》仍夺日冠 550亿的年度票...	众片开画《芳华》仍夺日冠 550亿的年度票...	http://www.entgroup.cn/news/Exclusive/23444...	2017-12-23	众 片 开 ...





The End of Lecture 7

Thank You

<http://www.wangting.ac.cn>

