

امیر علی بلباسی، ۹۸۳۱۱۰۹

دو خبر اول را در فایل دیگری برای تست آوردم:

[illegible]

فایل کلمات پرتکرار که در کتابخانه استفاده شده (hazm) موجود است:

```
stopwords.dat x
stopwords.dat
1
2 در
3 به
4 از
5 که
6 این
7 را
8 با
9 است
10 برای
11 آن
12 یک
13 خود
14 تا
15 کرد
16 بر
17 هم
18 نیز
19 گفت
20 می‌شود
21 وی
22 شد
23 دارد
24 ما
25 اما
26 با
27
```

In [1]:

```
%pip install hazm

from __future__ import unicode_literals
from hazm import *
import codecs
import json
Requirement already satisfied: hazm in d:\applications\anaconda\lib\site-
packages (0.7.0)
Requirement already satisfied: nltk==3.3 in
d:\applications\anaconda\lib\site-packages (from hazm) (3.3)
Requirement already satisfied: six in d:\applications\anaconda\lib\site-
packages (from nltk==3.3->hazm) (1.16.0)
Note: you may need to restart the kernel to use updated packages.
```

In [2]:

```
f = open('foo.json')
#f = open('IR_data_news_12k.json')
data = json.load(f)
f.close()

contents = []
for i in data:
    tmp = data[i]['content']
    contents.append(tmp[0:len(tmp)-16])
print('#1')
for i in contents:
    print(i)

normalizer = Normalizer()
for i in range(len(contents)):
    contents[i] = normalizer.normalize(contents[i])
print('#2')
for i in contents:
    print(i)

tokens = []
for i in range(len(contents)):
    tokens.append(word_tokenize(contents[i]))
print('#3')
for i in tokens:
    print(i)
    print()

lemmatizer = Lemmatizer()
for i in range(len(tokens)):
```

```

        for j in range(len(tokens[i])):
            tokens[i][j] = lemmatizer.lemmatize(tokens[i][j])
print('#4')
for i in tokens:
    print(i)
    print()

"""
stemmer = Stemmer()
for i in range(len(tokens)):
    for j in range(len(tokens[i])):
        tokens[i][j] = stemmer.stem(tokens[i][j])
print('#5')
for i in tokens:
    print(i)
    print()
"""

f = codecs.open('stopwords.dat', encoding='utf-8')
stopWords = []
for l in f.readlines():
    stopWords.append(l.strip('\n'))
f.close()
print('A#')
print(stopWords)
print()

for i in range(len(stopWords)):
    stopWords[i] = normalizer.normalize(stopWords[i])
print('B#')
print(stopWords)
print()

for i in range(len(stopWords)):
    stopWords[i] = lemmatizer.lemmatize(stopWords[i])
print('C#')
print(stopWords)
print()

tmpTokens = []
for i in range(len(tokens)):
    tmpTokensDoc = []
    for j in range(len(tokens[i])):
        if(tokens[i][j] not in stopWords): tmpTokensDoc.append(tokens[i][j])
    tmpTokens.append(tmpTokensDoc)

```

```

tokens = tmpTokens.copy()
print('#5')
for i in tokens:
    print(i)
    print()
#1

```

در نامه ای رسمی به (AFC) به گزارش خبرگزاری فارس، کنفدراسیون فوتبال آسیا فدراسیون فوتبال ایران و باشگاه گیتی پسند زمان قرعه کشی جام باشگاه های فوتسال آسیا را رسماً اعلام کرد. بر این اساس ۲۵ فروردین ماه ۱۴۰۱ مراسم قرعه کشی جام باشگاه های فوتسال آسیا در مالزی برگزار می شود. باشگاه گیتی پسند بعنوان قهرمان فوتسال ایران در سال ۱۴۰۰ به این مسابقات راه پیدا کرده است. پیش از این گیتی پسند تجربه ۳ دوره حضور در جام باشگاه های فوتسال آسیا را داشته که هر سه دوره به فینال مسابقات راه پیدا کرده و یک عنوان قهرمانی و دو مقام دومی بدست آورده است.

به گزارش خبرگزاری فارس، سید حمید سجادی در حاشیه مراسم گرامیداشت روز جوان در جمع خبرنگاران در رابطه با عرضه سهام سرخابی ها در بورس اظهار داشت: منتظر طی روند هستیم و بعداً اطلاع رسانی خواهیم کرد. وی در مورد حضور تماشاگران در مسابقات فوتبال اظهار داشت: حضور تماشاگران در لیگ برتر فوتبال تابع نظر فدراسیون است.، سازمان لیگ و ستاد ملی مبارزه با کرونا است.

#2

در نامه ای رسمی به (AFC) به گزارش خبرگزاری فارس، کنفدراسیون فوتبال آسیا فدراسیون فوتبال ایران و باشگاه گیتی پسند زمان قرعه کشی جام باشگاه های فوتسال آسیا را رسماً اعلام کرد. بر این اساس ۲۵ فروردین ماه ۱۴۰۱ مراسم قرعه کشی جام باشگاه های فوتسال آسیا در مالزی برگزار می شود. باشگاه گیتی پسند بعنوان قهرمان فوتسال ایران در سال ۱۴۰۰ به این مسابقات راه پیدا کرده است. پیش از این گیتی پسند تجربه ۳ دوره حضور در جام باشگاه های فوتسال آسیا را داشته که هر سه دوره به فینال مسابقات راه پیدا کرده و یک عنوان قهرمانی و دو مقام دومی بدست آورده است.

به گزارش خبرگزاری فارس، سید حمید سجادی در حاشیه مراسم گرامیداشت روز جوان در جمع خبرنگاران در رابطه با عرضه سهام سرخابی ها در بورس اظهار داشت: منتظر طی روند هستیم و بعداً اطلاع رسانی خواهیم کرد. وی در مورد حضور تماشاگران در مسابقات فوتبال اظهار داشت: حضور تماشاگران در لیگ برتر فوتبال تابع نظر فدراسیون، سازمان لیگ و ستاد ملی مبارزه با کرونا است.

#3

[('آسیا', 'فوتبال', 'کنفدراسیون', 'فارس', 'خبرگزاری', 'گزارش', 'به', 'فوتبال', 'فدراسیون', 'به', 'رسمی', 'ای\۲00c\نامه', 'در', 'AFC', 'جام', 'کشی', 'قرعه', 'زمان', 'پسند', 'گیتی', 'باشگاه', 'و', 'ایران', 'بر', 'کرد', 'اعلام', 'رسمی', 'را', 'آسیا', 'فوتسال', 'های\۲00c\باشگاه', 'جام', 'کشی', 'قرعه', 'مراسم', '۱۴۰۱', 'ماه', 'فروردین', '۲۵', 'اساس', 'این', 'شود\۲00c\می', 'برگزار', 'مالزی', 'در', 'آسیا', 'فوتسال', 'های\۲00c\باشگاه', 'در', 'ایران', 'فوتسال', 'قهرمان', 'بعنوان', 'پسند', 'گیتی', 'باشگاه', 'پیش', 'کرده\_است', 'پیدا', 'راه', 'مسابقات', 'این', 'به', '۱۴۰۰', 'سال']

'جام', 'در', 'حضور', 'دوره', '۳', 'تجربه', 'پسند', 'گیتی', 'این', 'از',  
'دوره', 'سه', 'هر', 'که', 'داشته', 'را', 'آسیا', 'فوتسال', 'هایu200c', 'باشگاه',  
'عنوان', 'یک', 'و', 'کرده', 'پیدا', 'راه', 'مسابقات', 'فینال', 'به',  
[.], 'آورده\_است', 'بدست', 'دومی', 'مقام', 'دو', 'و', 'قهرمانی'

'در', 'سجادی', 'حمید', 'سید', 'ا', 'فارس', 'خبرگزاری', 'گزارش', 'به',  
'در', 'خبرنگاران', 'جمع', 'در', 'جوان', 'روز', 'گرامیداشت', 'مراسم', 'حاشیه',  
'اظهار', 'بورس', 'در', 'هاu200c', 'سرخابی', 'سهام', 'عرضه', 'با', 'رابطه',  
'رسانی', 'اطلاع', 'بعدا', 'و', 'هستیم', 'روند', 'طی', 'منتظر', 'ا', 'داشت',  
'مسابقات', 'در', 'تماشاگران', 'حضور', 'مورد', 'در', 'وی', 'ا', 'خواهیم\_کرد',  
'برتر', 'لیگ', 'در', 'تماشاگران', 'حضور', 'ا', 'داشت', 'اظهار', 'فوتبال',  
'ستاد', 'و', 'لیگ', 'سازمان', 'ا', 'فدراسیون', 'نظر', 'تابع', 'فوتبال',  
[.], 'است', 'کرونا', 'با', 'مبارزه', 'ملی'

#4

('), 'آسیا', 'فوتبال', 'کنفدراسیون', 'ا', 'فارس', 'خبرگزاری', 'گزارش', 'به',  
'و', 'ایران', 'فوتبال', 'فدراسیون', 'به', 'رسم', 'نامه', 'در', 'ا', 'AFC',  
'باشگاه', 'جام', 'کشید#کش', 'قرعه', 'زمان', 'پسند', 'گیتی', 'باشگاه',  
'اساس', 'این', 'بر', 'ا', 'کرد#کن', 'اعلام', 'رسم', 'را', 'آسیا', 'فوتسال',  
'باشگاه', 'جام', 'کشید#کش', 'قرعه', 'مراسم', '۱۴۰۱', 'ماه', 'فروردین', '۲۵',  
'گیتی', 'باشگاه', 'ا', 'شد#شو', 'برگزار', 'مالزی', 'در', 'آسیا', 'فوتسال',  
'به', '۱۴۰۰', 'سال', 'در', 'ایران', 'فوتسال', 'قهرمان', 'بعنوان', 'پسند',  
'گیتی', 'این', 'از', 'پیش', 'ا', 'کرد#کن', 'پیدا', 'راه', 'مسابقات', 'این',  
'فوتسال', 'باشگاه', 'جام', 'در', 'حضور', 'دوره', '۳', 'تجربه', 'پسند',  
'مسابقات', 'فینال', 'به', 'دوره', 'سه', 'هر', 'که', 'داشته', 'را', 'آسیا',  
'مقام', 'دو', 'و', 'قهرمان', 'عنوان', 'یک', 'و', 'کرده', 'پیدا', 'راه',  
[.], 'آورد#آور', 'بدست', 'دوم'

'در', 'سجاد', 'حمید', 'سید', 'ا', 'فارس', 'خبرگزاری', 'گزارش', 'به',  
'در', 'خبرنگار', 'جمع', 'در', 'جوان', 'روز', 'گرامیداشت', 'مراسم', 'حاشیه',  
'اظهار', 'بورس', 'در', 'هاu200c', 'سرخابی', 'سهام', 'عرضه', 'با', 'رابطه',  
'اطلاع', 'بعدا', 'و', 'هست', 'روند', 'طی', 'منتظر', 'ا', 'داشت',  
'در', 'تماشاگر', 'حضور', 'مورد', 'در', 'وی', 'ا', 'کرد#کن', 'رساند#رسان',  
'لیگ', 'در', 'تماشاگر', 'حضور', 'ا', 'داشت', 'اظهار', 'فوتبال', 'مسابقات',  
'و', 'لیگ', 'سازمان', 'ا', 'فدراسیون', 'نظر', 'تابع', 'فوتبال', 'برتر',  
[.], 'است', 'کرونا', 'با', 'مبارزه', 'ملی', 'ستاد'

A#

'یک', 'آن', 'برای', 'است', 'با', 'را', 'این', 'که', 'از', 'به', 'در', 'و',  
'شد', 'وی', 'شودu200c', 'می', 'گفت', 'نیز', 'هم', 'بر', 'کرد', 'تا', 'خود',  
'دیگر', 'او', 'بود', 'آنها', 'هر', 'باید', 'شده', 'یا', 'اما', 'ما', 'دارد',  
'شده\_است', 'پیش', 'بین', 'وجود', 'کند', 'شود', 'کندu200c', 'می', 'مورد', 'دو',  
'نیست', 'کنند', 'من', 'هستند', 'حال', 'یکی', 'همه', 'اگر', 'نظر', 'پس',  
'هایی', 'افزود', 'همین', 'کنندu200c', 'می', 'بخش', 'می', 'بی', 'چه', 'باشد',  
'داشت', 'سه', 'سیار', 'بیشتر', 'داد', 'روی', 'همچنین', 'راه', 'دارند'

'ولی', 'جدید', 'بعد', 'شدن', 'اینکه', 'میان', 'هیچ', 'تنها', 'سوی', 'چند',  
 'کرده\_است', 'نه', 'اول', 'دهد\u200cمی', 'کردند', 'برخی', 'کردن', 'حتی',  
 'بار', 'درباره', 'تمام', 'افراد', 'طور', 'چنین', 'شما', 'بیش', 'نسبت',  
 'طی', 'بزرگ', 'دوم', 'ندارد', 'چون', 'کرده', 'تواند\u200cمی', 'بسیاری',  
 'خواهد\_شد', 'دیگری', 'گوید\u200cمی', 'آنان', 'البته', 'بدون', 'همان', 'حدود',  
 'قبل', 'ویژه', 'کل', 'وارد', 'توان\u200cمی', 'رشد', 'یعنی', 'قابل', 'کنیم',  
 'چرا', 'بوده\_است', 'سازی', 'لازم', 'هنوز', 'گذاری', 'نیاز', 'براساس',  
 'پیدا', 'تغییر', 'حالی', 'جای', 'کم', 'گرفت', 'وقتی', 'شوند\u200cمی',  
 'رو', 'بیان', 'آیا', 'تعداد', 'زیادی', 'فقط', 'مدت', 'باعث', 'تحت', 'اکنون',  
 'برابر', 'دهد', 'جاری', 'بلکه', 'نوع', 'بودن', 'کرده\_اند', 'عدم', 'شدند',  
 'آقای', 'خصوص', 'شاید', 'گیری', 'زیر', 'امر', 'مربوط', 'اخیر', 'بوده', 'مهم',  
 'کنید', 'سایر', 'سوم', 'اولین', 'کنار', 'فکر', 'بودند', 'کننده', 'اثر',  
 'مثل', 'پی', 'دارای', 'حل', 'ممکن', 'گیرد\u200cمی', 'باز', 'مانند', 'ضمن',  
 'آنچه', 'امکان', 'طول', 'موجب', 'کسی', 'منظور', 'دور', 'اجرا', 'رسد\u200cمی',  
 'رسید', 'تاکنون', 'گونه', 'علاوه', 'خیلی', 'جمع', 'شوند', 'گفته', 'تعیین',  
 'طرف', 'خواهد\_بود', 'داشته\_باشد', 'چهار', 'علت', 'شده\_اند', 'گرفته', 'ساله',  
 'جریان', 'نزدیک', 'توانند\u200cمی', 'مشخص', 'زیرا', 'مناسب', 'تبدیل', 'تهیه',  
 'ریزی', 'پنج', 'بالا', 'نخستین', 'یافت', 'دهند\u200cمی', 'بنابراین', 'روند',  
 'خوب', 'خوبی', 'خاص', 'شده\_بود', 'ترتیب', 'بیشتری', 'نخست', 'چیزی', 'عالی',  
 'جدی', 'دادن', 'آخرین', 'دهند', 'رود\u200cمی', 'غیر', 'کامل', 'فرد', 'شروع',  
 'حد', 'داده\_است', 'بهرتر', 'تمامی', 'باشند', 'بخشی', 'گیرد', 'شامل', 'بهرترین',  
 'دانست', 'باشد\u200cمی', 'علیه', 'داریم', 'کرد\u200cمی', 'کسانی', 'نبود',  
 'دچار', 'گرفته\_است', 'آنجا', 'ایشان', 'شد\u200cمی', 'دهه', 'داشتند', 'ناشی',  
 'برداری', 'اند', 'هستیم', 'بعضی', 'داده', 'آنکه', 'لحاظ', 'آید\u200cمی',  
 'وگو', 'اش', 'آمد', 'همیشه', 'سهم', 'نشست', 'کنیم\u200cمی', 'نباید',  
 'رفت', 'چگونه', 'نوعی', 'خواهد\_کرد', 'جا', 'طبق', 'حداقل', 'کنم\u200cمی',  
 'کافی', 'کلی', 'شمار', 'بندی', 'سعی', 'ندارند', 'روش', 'فوق', 'هنگام',  
 'آورد', 'پشت', 'چیز', 'داشته\_است', 'کوچک', 'سمت', 'زیاد', 'همچنان', 'مواجه',  
 'نیمه', 'عهده', 'کردند\u200cمی', 'دادند', 'های\u200cسال', 'روبه', 'حالا',  
 'سپس', 'کنم', 'جز', 'آمده\_است', 'یکدیگر', 'بروز', 'سی', 'دیگران', 'جایی',  
 'رسیدن', 'شود\u200cنمی', 'صرف', 'شان', 'یافته', 'همواره', 'خودش', 'کنندگان',  
 'کردم', 'نحوه', 'باره', 'کرده\_بود', 'داشته', 'ساز', 'متر', 'یابد', 'چهارم',  
 'سراسر', 'متفاوت', 'کمی', 'پخش', 'محسوب', 'داشته\_باشند', 'شخصی', 'تو',  
 'خطر', 'همچون', 'ع', 'فردی', 'گروهی', 'آمده', 'نظیر', 'داشتن', 'کاملاً',  
 'دار', 'بیرون', 'متاسفانه', 'آوری', 'عین', 'سبب', 'دسته', 'کدام', 'خویش',  
 'نیستند', 'درون', 'سالهای', 'گویند\u200cمی', 'افراد', 'شش', 'ابتدا',  
 'دوباره', 'اغلب', 'جمعی', 'گاه', 'خاطرنشان', 'پر', 'یافته\_است',  
 '[اینجا', 'گردد', 'زاده', 'لذا', 'یابد\u200cمی']

B#

'یک', 'آن', 'برای', 'است', 'با', 'را', 'این', 'که', 'از', 'به', 'در', 'و',  
 'شد', 'وی', 'شود\u200cمی', 'گفت', 'نیز', 'هم', 'بر', 'کرد', 'تا', 'خود',  
 'دیگر', 'او', 'بود', 'آنها', 'هر', 'باید', 'شده', 'یا', 'اما', 'ما', 'دارد',  
 'شده\_است', 'پیش', 'بین', 'وجود', 'کند', 'شود', 'کند\u200cمی', 'مورد', 'دو',  
 'نیست', 'کنند', 'من', 'هستند', 'حال', 'یکی', 'همه', 'اگر', 'نظر', 'پس'

'هایِ', 'افزود', 'همین', 'کنند\u200cمی', 'بخش', 'می', 'بی', 'چه', 'باشد',  
 'داشت', 'سه', 'بسیار', 'بیشتر', 'داد', 'روی', 'همچنین', 'راه', 'دارند',  
 'ولی', 'جدید', 'بعد', 'شدن', 'اینکه', 'میان', 'هیچ', 'تنها', 'سوی', 'چند',  
 'کرده\_است', 'نه', 'اول', 'دهد\u200cمی', 'کردند', 'برخی', 'کردن', 'حتی',  
 'بار', 'درباره', 'تمام', 'افراد', 'طور', 'چنین', 'شما', 'بیش', 'نسبت',  
 'طی', 'بزرگ', 'دوم', 'ندارد', 'چون', 'کرده', 'تواند\u200cمی', 'بسیاری',  
 'خواهد\_شد', 'دیگری', 'گوید\u200cمی', 'آنان', 'البته', 'بدون', 'همان', 'حدود',  
 'قبل', 'ویژه', 'کل', 'وارد', 'توان\u200cمی', 'رشد', 'یعنی', 'قابل', 'کنیم',  
 'چرا', 'بوده\_است', 'سازی', 'لازم', 'هنوز', 'گذاری', 'نیاز', 'براساس',  
 'پیدا', 'تغییر', 'حالی', 'جای', 'کم', 'گرفت', 'وقتی', 'شوند\u200cمی',  
 'رو', 'بیان', 'آیا', 'تعداد', 'زیادی', 'فقط', 'مدت', 'باعث', 'تحت', 'اکنون',  
 'برابر', 'دهد', 'جاری', 'بلکه', 'نوع', 'بودن', 'کرده\_اند', 'عدم', 'شدند',  
 'آقای', 'خصوص', 'شاید', 'گیری', 'زیر', 'امر', 'مربوط', 'اخیر', 'بوده', 'مهم',  
 'کنید', 'سایر', 'سوم', 'اولین', 'کنار', 'فکر', 'بودند', 'کننده', 'اثر',  
 'مثل', 'پی', 'دارای', 'حل', 'ممکن', 'گیرد\u200cمی', 'باز', 'مانند', 'ضمن',  
 'آنچه', 'امکان', 'طول', 'موجب', 'کسی', 'منظور', 'دور', 'اجرا', 'رسد\u200cمی',  
 'رسید', 'تاکنون', 'گونه', 'علاوه', 'خیلی', 'جمع', 'شوند', 'گفته', 'تعیین',  
 'طرف', 'خواهد\_بود', 'داشته\_باشد', 'چهار', 'علت', 'شده\_اند', 'گرفته', 'ساله',  
 'جریان', 'نزدیک', 'توانند\u200cمی', 'مشخص', 'زیرا', 'مناسب', 'تبدیل', 'تهیه',  
 'ریزی', 'پنج', 'بالا', 'نخستین', 'یافت', 'دهند\u200cمی', 'بنابراین', 'روند',  
 'خوب', 'خوبی', 'خاص', 'شده\_بود', 'ترتیب', 'بیشتری', 'نخست', 'چیزی', 'عالی',  
 'جدی', 'دادن', 'آخرین', 'دهند', 'رود\u200cمی', 'غیر', 'کامل', 'فرد', 'شروع',  
 'حد', 'داده\_است', 'بهتر', 'تمامی', 'باشند', 'بخشی', 'گیرد', 'شامل', 'بهترین',  
 'دانست', 'باشد\u200cمی', 'علیه', 'داریم', 'کرد\u200cمی', 'کسانی', 'نبود',  
 'دچار', 'گرفته\_است', 'آنجا', 'ایشان', 'شد\u200cمی', 'دهه', 'داشتند', 'ناشی',  
 'برداری', 'اند', 'هستیم', 'بعضی', 'داده', 'آنکه', 'لحاظ', 'آید\u200cمی',  
 'وگو', 'اش', 'آمد', 'همیشه', 'سهم', 'نشست', 'کنیم\u200cمی', 'نباید',  
 'رفت', 'چگونه', 'نوعی', 'خواهد\_کرد', 'جا', 'طبق', 'حداقل', 'کنم\u200cمی',  
 'کافی', 'کلی', 'شمار', 'بندی', 'سعی', 'ندارند', 'روش', 'فوق', 'هنگام',  
 'آورد', 'پشت', 'چیز', 'داشته\_است', 'کوچک', 'سمت', 'زیاد', 'همچنان', 'مواجه',  
 'نیمه', 'عهده', 'کردند\u200cمی', 'دادند', 'های\u200cسال', 'روبه', 'حالا',  
 'سپس', 'کنم', 'جز', 'آمده\_است', 'یکدیگر', 'بروز', 'سی', 'دیگران', 'جایی',  
 'رسیدن', 'شود\u200cنمی', 'صرف', 'شان', 'یافته', 'همواره', 'خودش', 'کنندگان',  
 'کردم', 'نحوه', 'باره', 'کرده\_بود', 'داشته', 'ساز', 'متر', 'یابد', 'چهارم',  
 'سراسر', 'متفاوت', 'کمی', 'پخش', 'محسوب', 'داشته\_باشند', 'شخصی', 'تو',  
 'خطر', 'همچون', 'ع', 'فردی', 'گروهی', 'آمده', 'نظیر', 'داشتن', 'کاملاً',  
 'دار', 'بیرون', 'متاسفانه', 'آوری', 'عین', 'سبب', 'دسته', 'کدام', 'خویش',  
 'نیستند', 'درون', 'سالهای', 'گویند\u200cمی', 'افراد', 'شش', 'ابتدا',  
 'دوباره', 'اغلب', 'جمعی', 'گاه', 'خاطرنشان', 'پر', 'یافته\_است',  
 '[اینجا', 'گردد', 'زاده', 'لذا', 'یابد\u200cمی']

C#

'یک', 'آن', 'برای', 'است#', 'با', 'را', 'این', 'که', 'از', 'به', 'در', 'و',  
 'شد#شو', 'وی', 'شد#شو', 'گفت#گو', 'نیز', 'هم', 'بر', 'کرد#کن', 'تا', 'خود',  
 'او', 'بود#باش', 'آن', 'هر', 'باید', 'شده', 'یا', 'اما', 'ما', 'داشت#دار'

, 'پیش', 'بین', 'وجود', 'کند', 'شد#شو', 'کرد#کن', 'مورد', 'دو', 'دیگر', 'کرد#کن', 'من', 'هست#', 'حال', 'یک', 'همه', 'اگر', 'نظر', 'پس', 'شد#شو', 'افزود#افزا', 'همین', 'کرد#کن', 'بخش', 'می', 'بی', 'چه', 'بود#باش', 'نیست', 'سه', 'بسیار', 'بیشتر', 'داد', 'روی', 'همچنین', 'راه', 'داشت#دار', 'های', 'جدید', 'بعد', 'شدن', 'اینکه', 'میان', 'هیچ', 'تنها', 'سو', 'چند', 'داشت', 'کرد#کن', 'نه', 'اول', 'داد#ده', 'کرد#کن', 'برخی', 'کردن', 'حتی', 'ولی', 'بار', 'درباره', 'تمام', 'افراد', 'طور', 'چنین', 'شما', 'بیش', 'نسبت', 'طی', 'بزرگ', 'دوم', 'داشت#دار', 'چون', 'کرده', 'توانست#توان', 'بسیار', 'کرد#کن', 'شد#شو', 'دیگر', 'گفت#گو', 'آن', 'البته', 'بدون', 'همان', 'حدود', 'براساس', 'قبل', 'ویژه', 'کل', 'وارد', 'توان\200cمی', 'رشد', 'یعنی', 'قابل', 'شد#شو', 'چرا', 'بود#باش', 'سازی', 'لازم', 'هنوز', 'گذاشت#گذار', 'نیاز', 'تحت', 'اکنون', 'پیدا', 'تغییر', 'حال', 'جای', 'کم', 'گرفت#گیر', 'وقت', 'عدم', 'شد#شو', 'رو', 'بیان', 'آیا', 'تعداد', 'زیاد', 'فقط', 'مدت', 'باعث', 'بوده', 'مهم', 'برابر', 'داد#ده', 'جاری', 'بلکه', 'نوع', 'بودن', 'کرد#کن', 'اثر', 'آقا', 'خصوص', 'شاید', 'گرفت#گیر', 'زیر', 'امر', 'مربوط', 'اخیر', 'ضمن', 'کرد#کن', 'سایر', 'سوم', 'اولین', 'کنار', 'فکر', 'بود#باش', 'کننده', 'رسید#رس', 'مثل', 'پی', 'دارا', 'حل', 'ممکن', 'گرفت#گیر', 'باز', 'مانند', 'تعیین', 'آنچه', 'امکان', 'طول', 'موجب', 'کس', 'منظور', 'دور', 'اجرا', 'ساله', 'رسید', 'تاکنون', 'گونه', 'علاوه', 'خیلی', 'جمع', 'شد#شو', 'گفته', 'تهیه', 'طرف', 'بود#باش', 'داشت#دار', 'چهار', 'علت', 'شد#شو', 'گرفته', 'روند', 'جریان', 'نزدیک', 'توانست#توان', 'مشخص', 'زیرا', 'مناسب', 'تبدیل', 'عالی', 'ریخت#ریز', 'پنج', 'بالا', 'نخستین', 'یافت#یاب', 'داد#ده', 'بنابراین', 'شروع', 'خوب', 'خوبی', 'خاص', 'شد#شو', 'ترتیب', 'بیشتر', 'نخست', 'چیز', 'بهترین', 'جدی', 'دادن', 'آخرین', 'داد#ده', 'رفت#رو', 'غیر', 'کامل', 'فرد', 'حد', 'داد#ده', 'بهتر', 'تمام', 'بود#باش', 'بخشی', 'گرفت#گیر', 'شامل', 'دانست#دان', 'بود#باش', 'علیه', 'داشت#دار', 'کرد#کن', 'کسانی', 'بود#باش', 'دچار', 'گرفت#گیر', 'آنجا', 'ایشان', 'شد#شو', 'دهه', 'داشت#دار', 'ناشی', 'نباید', 'بردار', 'اند', 'هست#', 'بعضی', 'داده', 'آنکه', 'لحاظ', 'آمد#آ', 'حداقل', 'کرد#کن', 'وگو', 'اش', 'آمد#آ', 'همیشه', 'سهم', 'نشست', 'کرد#کن', 'روش', 'فوق', 'هنگام', 'رفت#رو', 'چگونه', 'نوع', 'کرد#کن', 'جا', 'طبق', 'همچنان', 'مواجه', 'کافی', 'کلی', 'شمار', 'بست#بند', 'سعی', 'داشت#دار', 'روبه', 'حالا', 'آورد#آور', 'پشت', 'چیز', 'داشت#دار', 'کوچک', 'سمت', 'زیاد', 'بروز', 'سی', 'دیگر', 'جا', 'نیمه', 'عهده', 'کرد#کن', 'داد#ده', 'سال', 'همواره', 'خودش', 'کنندگان', 'سپس', 'کرد#کن', 'جز', 'آمد#آ', 'یکدیگر', 'ساز', 'متر', 'یافت#یاب', 'چهارم', 'رسیدن', 'شد#شو', 'صرف', 'شان', 'یافته', 'داشت#دار', 'شخص', 'تو', 'کرد#کن', 'نحوه', 'باره', 'کرد#کن', 'داشته', 'آمده', 'نظیر', 'داشتن', 'کاملا', 'سراسر', 'متفاوت', 'کم', 'پخش', 'محسوب', 'عین', 'سبب', 'دسته', 'کدام', 'خویش', 'خطر', 'همچون', 'ع', 'فرد', 'گروهی', 'گفت#گو', 'افراد', 'شش', 'ابتدا', 'دار', 'بیرون', 'متاسفانه', 'آورد#آور', 'اغلب', 'جمع', 'گاه', 'خاطر نشان', 'پر', 'یافت#یاب', 'هست#', 'درون', 'سال', '[اینجا', 'گشت#گرد', 'زاده', 'لذا', 'یافت#یاب', 'دوباره']

#5

, 'آسیا', 'فوتبال', 'کنفدراسیون', 'ء', 'فارس', 'خبرگزاری', 'گزارش', 'گیتی', 'باشگاه', 'ایران', 'فوتبال', 'فدراسیون', 'رسم', 'نامه', 'AFC', ')



'رسم', 'آسیا', 'فوتسال', 'باشگاه', 'جام', 'کشید#کش', 'قرعه', 'زمان', 'پسند',  
'قرعه', 'مراسم', '۱۴۰۱', 'ماه', 'فروردین', '۲۵', 'اساس', '.', 'اعلام',  
'', 'برگزار', 'مالزی', 'آسیا', 'فوتسال', 'باشگاه', 'جام', 'کشید#کش',  
'۱۴۰۰', 'ایران', 'فوتسال', 'قهرمان', 'بعنوان', 'پسند', 'گیتی', 'باشگاه',  
'جام', 'حضور', 'دوره', '۳', 'تجربه', 'پسند', 'گیتی', '.', 'مسابقات',  
'قهرمان', 'عنوان', 'مسابقات', 'فینال', 'دوره', 'آسیا', 'فوتسال', 'باشگاه',  
[',', 'بدست', 'مقام']

'مراسم', 'حاشیه', 'سجاد', 'حمید', 'سید', '،', '،', 'فارس', 'خبرگزاری', 'گزارش',  
'سهام', 'عرضه', 'رابطه', 'خبرنگار', 'جوان', 'روز', 'گرامیداشت',  
'اطلاع', 'بعدا', 'منتظر', ':', '،', 'اظهار', 'بورس', 'ها200c\u2013', 'سرخابی',  
'،', 'اظهار', 'فوتبال', 'مسابقات', 'تماشاگر', 'حضور', '.', 'رساند#رسان',  
'،', 'فدراسیون', 'تابع', 'فوتبال', 'برتر', 'لیگ', 'تماشاگر', 'حضور',  
[',', 'کرونا', 'مبارزه', 'ملی', 'ستاد', 'لیگ', 'سازمان']

همانطور که در مثال مشاهده می کنید در مرحله اول فایل تست را لود کردیم سپس محتوا اخبار را استخراج کرده و پس از آن متن محتوا را نرمال می کنیم (در هر مرحله هم چاپ می کنیم تا تغییرات مشخص شود. مثلا با نرمال سازی space های بین "ها" و "خود کلمه" تبدیل به نیم فاصله می شود). در مرحله بعدی توکن ها را استخراج کرده و در یک لیست از لیست ها ذخیره می کنیم. در مرحله بعد ریشه کلمات را جایگزین فرم های مختلف آن می کنیم (مثلا "باشگاه های" تبدیل به "باشگاه" می شود). سپس همین مراحل را روی stop word ها تکرار می کنیم تا مشابه هم شوند. سر آخر نیز با حلقه روی لیست توکن ها، آن توکن هایی که جزو کلمات پر تکرار هستند را حذف می کنیم (مانند توکن "در").

دلیل پردازش نرمال سازی یکدست کردن کل متون است زیرا مثلا ممکن است نویسنده در بعضی جاها نیم فاصله را رعایت کرده باشد و در بعضی جاها نکرده باشد. دلیل استخراج توکن ها هم که واضح است زیرا برای ساخت شاخص مکانی لازم اند. بعضی کلمات در جاهای مختلف یک یا چند سند با اشکال متفاوت ظاهر می شوند ولی همه آنها به یک مفهوم اشاره می کنند و یک ریشه مشترک دارند. به منظور یک جستجو دقیق تر، لازم است تمام این اشکال به یک شکل واحد در آورده شوند. به همین جهت پردازش ریشه یابی را انجام می دهیم. شرکت دادن کلمات پر تکرار فارسی مانند "از" در عملیات جستجو و بازیابی اسناد، دقت ما را به مراتب پایین می آورد (چون این کلمات در عموم اسناد وجود دارند). به همین جهت با پیش پردازش، این دست کلمات را از لیست توکن های استخراج شده حذف می کنیم.