



دانشگاه صنعتی امیرکبیر
(پلی تکنیک تهران)

فاز اول پروژه بازیابی اطلاعات

نام دانشجو: سپهر عسگریان ابیانه

شماره دانشجو: 9531901

اردیبهشت 1399

گزارش پروژه

ابتدا برای هر دو حالت تمامی دیتا ها با یک دیگر مرج شده اند و تشکیل دیتایی با 55000 داده را داده اند.

ابتدا حالت اول به اسم فایل حالت اول که در jupyter پیاده سازی شده است و همچنین نتایج آن نیز در آن فایل قابل مشاهده است .

دیتای اولیه به صورت زیر قابل مشاهده است

```
In [134]: data=df
data['content']

Out[134]: 0      <strong><span dir="RTL">گزارش<span dir="RTL">به گزارش<a href=...
1      <p>کاهش تنوع مدارس، شمار سال‌های اخیر سکندرا...
2      <div class="hideTag">\n <a href="/fa/service/3...
3      <div class="item-text" itemprop="articleBody">...
4      <p style="text-align: justify;">به گزارش<a cl...
5      <div class="body col-xs-36"> \n <a class="ente...
6      <div class="item-text" itemprop="articleBody">...
7      <div class="hideTag">\n <a href="/fa/service/3...
8      <div class="image-news pull-right mrxxl image-...
9      <p>به گزارش ایسنا، پس از اینکه سمیه یزدانی در...
10     <p><img class="image_btn" style="margin: 10px ...
11     <a href="/" class="irinn_link">گزارش یلگاد<span st...
12     <div dir="rtl">\n به گزارش خبرنگار \n <span st...
13     <div class="hideTag">\n <a href="/fa/service/8...
14     <p style="text-align: justify;"><span style="col...
15     <p>ایران به واسطه عدم انجام تمهیدات برجای می‌آید...
16     <p><img class="image_btn" style="margin: 10px ...
17     <div dir="rtl">\n به گزارش کیناک به نقل از <span st...
18     <p>استرس تنها فاکتوری نیست که موجب سفید شدن زو...
19     <img align="left" class="news_corner_image" sr...
20     <p><img class="image_btn" style="margin: 10px ...
21     <p><img class="image_btn" style="margin: 0px 8...
22     <div dir="rtl">\n حساب کاربری زندانیان سیاسی ع...
23     <p><span style="font-size: 14pt;"><img class="...
24     <p><span style="font-size: 12pt;">به گزارش خبرنگار...
25     <div class="hideTag">\n <a href="/fa/service/6...
26     <div class="hideTag">\n <a href="/fa/service/3...
27     <div class="hideTag">\n <a href="/fa/service/6...
28     <div class="hideTag">\n <a href="/fa/service/8...
29     <div class="hideTag">\n <a href="/fa/service/8...
...
7714    <div class="row " id="slideshowDisplay"> \n <d...
7715    <p>به گزارش ایسنا، به نقل از سایت شبکه روسیا ا...
```

1. حالت اول

در حالت اول نرمالایز خیلی ساده همراه با حذف stop words داشتیم به همراه توکنایز کردن بسیار ساده که مثلاً برای واژه

کتاب ها توکنایز به صورت کتاب و ها تشکیل میشد که در اینجا توکنایز اضافی داریم . ابتدا در این بخش با استفاده از

دستور

```
In [83]: regex = re.compile('[a-zA-Z]')
data['content'] = data['content'].map(lambda x: re.sub('[a-zA-Z]', '', x))
```


حال در این مرحله تمامی **stop words** از یک فایل **word** که در اینترنت نیز موجود بود برداشته شده است و خوانده شده

سپس همانطور که در صورت پروژه آمده است **tokenize** بر اساس فاصله ای که موجود است در متن انتخاب شده است و همچنین **stopwords** از متن بیرون کشیده شده اند(فایل **stopwords** استفاده شده در داکيومنت گذاشته شده است) همانند شکل زیر قابل مشاهده است.

```
In [37]: #Read stopwords file .txt
def read_file_stopwords(addr):#read file
    list_result = []
    with open(addr, 'r', encoding="utf_8") as f:
        stopwords = f.read()
        stopwords = stopwords.splitlines()
        for word in stopwords:
            list_result.append(word)

    return list_result
stopword=read_file_stopwords(r'C:\Users\sepehr\Desktop\IR\project\IR-S19-project-data (1)\IR-S19-project-data\project1_phase1
```

```
In [38]: #Tokenization with space

def tokenization(text):
    text = re.split('\W+', text)
    return text

data['content'] = data['content'].apply(lambda x: tokenization(x))
print("tokenization done")

#Removing stop words after Tokenization
def remove_stopwords(text):#remove stop words
    text = [word for word in text if word not in stopword]
    return text

data['content'] = data['content'].apply(lambda x: remove_stopwords(x))
print("remove_stopwords done")

tokenization done
remove_stopwords done
```

حال برای مثال بعد از مراحل بالا به شکل زیر میشود خروجی

```
In [210]: data['content'][3445]

Out[210]: ['حیدر',
            'عباس',
            'زاده',
            'فرمانده',
            'نیروی',
            'انتظامی',
            'استان',
            'خوزستان',
            'امروز',
            'فارس',
            'مرز',
            'تلمچه',
            'اظهار',
            'راستی',
            'تامین',
            'امنیت',
            'زالران',
            'اربعین',
            'جانبی',
            '...']
```

2-حالت دوم :

در حالت دوم نرمالایز کردن با دقت بیشتری انجام میشود همانطور که در صورت پروژه نیز گفته شده است پس از مرج کردن دیتا ابتدا حالت های مختلف کلمه به یک حالت واحد تبدیل شده است

کد آن مانند شکل زیر میباشد

```
def cheak_spell_arabic(text):
    aftery = re.sub("ء", "ئ", text)
    aftera = re.sub("أ", "ا", aftery)
    abfterb = re.sub("ب", "ب", aftera)
    afterp = re.sub("پ", "پ", abfterb)
    aftert = re.sub("ت", "ت", afterp)
    afterc = re.sub("ث", "ث", aftert)
    afterj = re.sub("ج", "ج", afterc)
    afterch = re.sub("ح", "ح", afterj)
    afterkh = re.sub("خ", "خ", afterch)
    afterd = re.sub("د", "د", afterkh)
    afterz = re.sub("ذ", "ذ", afterd)
    afterrr = re.sub("ر", "ر", afterz)
    afterzi = re.sub("ز", "ز", afterrr)
    afterzh = re.sub("ه", "ه", afterzi)
```


فبار ورزشی صادق ورمزبار در گفت وگو با خبرنگار ورزشی خبرگزاری تسنیم، درباره برتری بر صفر تیم فوتبال استقلال مقابل فجر سپاسی شیراز در مرحله یک هشتم نهایی اظهار داشت این مسابقه نتیجه خوبی برای استقلال در پی داشت و صعود به مرحله بعدی جام حذفی، نکته بسیار خوبی بود، اما تیم از لحاظ فنی خیلی قوی نبود و هنوز به مرین بیشتری نیاز دارد استقلال باید بیشتر از این کار کند تا در لیگ برتر موفق باشد، وی ضمن بیان این جمله تصریح کرد در نیمه اول چیزی از استقلال ندیدیم در نیمه دوم بکیت توپ و میدان را به خوبی در اختیار داشتیم، با این حال نتوانستیم به خوبی موقعیت های گلزنی طراحی کنیم معلوم نبود که تیم می خواهد از چه راهی به گل برسد، از جناحین یا پاس در صق این حملات سردرگم باعث شد استقلال زمان را از دست بدهد استراماچونی یک قدم رو به جلو برداشتیم اما هنوز خیلی کار داریم پیشکوت باشگاه است سخ به این پرسش که آیا آبی پوشان به بازیکنی جدید نیاز دارند یا خیر، خاطرنشان کرد استقلال به هافبکی طراح نیاز دارد تا به بازی نظم دهد و حملات تیم را طراحی کند این به زمان احتیاج دارد تا هماهنگ شود و افکار استراماچونی را در تیم بینیم جابه جایی متعدد بازیکنان و کمبود زمان به بازیکنان اجازه نداده که به هماهنگی برسند استقلال بند بازیکنی که بتوانند تیم را راه بیندازند، نیاز دارد ورمزبار با اشاره به تفاوت آبی پوشان با روزهای اوج شان به تسنیم گفت تیم فعلی با استقلال واقعی خیلی فاصله دارد و بداندند پیراهن چه تیمی را پوشیده اند استقلال میلیون ها هوادار دارد که آنها منتظر بازی های خوب این بازیکنان هستند این طرفداران با بازی های استقلال زندگی می کنند و باد این مسائل را بدانند وی در پاسخ به این پرسش که آیا عملکرد شاگردان استراماچونی تا اینجای فصل، امینوارکننده بوده است یا خیر، بیان کرد شکی در این نیست که استقلال ی فصل تاکنون، سیری صعودی طی کرده است پس از اینکه در یکی دو بازی عملکرد متوسطی داشتیم، کم کم به درجه خوبی از اعتمادبه نفس رسیدیم و توانستیم سه بازی متو استقلال برای رسیدن به تیم آرمانی فاصله زیادی دارد، اما این سه پیروزی می تواند نقطه امید برای تیم و هواداران باشد مربی پیشین تیم فوتبال استقلال درباره شانس موفق در هفته های آتی لیگ برتر فوتبال اظهار داشت استقلال در لیگ برتر، بازی های سختی پیش رو دارد که فکر می کنم اگر پیش از این دیدارها به هماهنگی لازم برسد، می ففیت تیم امینوار بود، اما اگر این هماهنگی تا بازی های سخت هم ایجاد نشود و روند پیروزی های استقلال متوقف شود، تیم دچار استرس می شود که این برای استقلال خطرا،

حال در این قسمت علاوه بر خواندن فایل stopwords فایلی به نام عبارت ها نیز ذخیره شده است که از گیت هاب فایل ورودی آن را پیدا کرده ام. این فایل شامل عباراتی مثل فیه مافیه مع ذلک علی ای حال و غیره میباشد که همه در توکن کردن به یک شکل میباشد پس بنابراین برای اینکار این عبارات خوانده شده است و همانند شکل زیر مشکل آن همدل شده است

```
def _fix_phrases(text):
    for phrase in phrases:
        text = text.replace(phrase, iget_whole_phrase(phrase))
        text = _fix_white_space(text)
    return text

def iget_whole_phrase(ph):
    ph = ph.replace(" ", '\u200c')
    return ph

def _fix_white_space(text):
    text = re.sub("\s+", " ", text)
    text = text.strip()
    return text
```

که در اینجا برای هر کدام از آن ها بینشان اگر فاصله موجود میباشد پاک شده است و برای مثال علی ای حال به صورت جدا آمده که همانطور که می بینید به صورت سر هم نوشته شده است

همچنین این عبارت

```
data['content'] = data['content'].apply(lambda x: fix_phrases(x))
```

سپس برای اینکه علائم جمع توکن اضافی همانند حالت اول به حساب نگیرد کد زیر نوشته شده است که علاوه بر علائم جمع پسوند ها و بیوند های مهم دیگر نیز دقت شده که در صورت مشاهده با نیم فاصله بیابند

برای مثال در شکل زیر این تابع برای رشته مشخص شده پیوند ها و پیوندها را درست کرده است

حال مانند قبل **tokenize** و **stopwords** ها هندل شده اند لازم به ذکر است که در **TOKENIZE** علائم پیشوندی و پسوندی به صورت کامل درست جدا شده اند و مزیت آن این است که در مرحله ریشه یابی با استفاده از کتابخانه **hazm** به خوبی ریشه یابی و **lemmatize** انجام میشود

سیس با استفاده از کتابخانه hazm ریشه یابی و lemmatize انجام دشه است همانطور که در شکل زیر میبینید

موارد خواسته شده بدست آمده شده است


```

from hazm import *
lemmatizer = Lemmatizer()
print(lemmatizer.lemmatize('گفت').split('#')[1])
print(lemmatizer.lemmatize('گو'))
print(lemmatizer.lemmatize('رود'))
print(lemmatizer.lemmatize('رو'))
print(lemmatizer.lemmatize('خواه'))
print(lemmatizer.lemmatize('سیاس'))
print(lemmatizer.lemmatize('هنر'))
print(lemmatizer.lemmatize('شریف'))
print(lemmatizer.lemmatize('دوست'))
print(lemmatizer.lemmatize('یاد'))
print(lemmatizer.lemmatize('توان'))
print(lemmatizer.lemmatize('شنو'))
print(lemmatizer.lemmatize('کرد'))
print(lemmatizer.lemmatize('ساز'))
print(lemmatizer.lemmatize('دان'))

```

گو
گو
رود
رو
خواه
سیاس
هنر
شریف
دوست
یاد
توان
شنو
کرد##کن
ساز
دان

```

stemmer = Stemmer()
stemmer.stem('کتاب•ها')

```

'کتاب'

این عملیات مانند عملیات ریشه یابی برای انگلیسی است با این تفاوت که library آن hazm میباشد

```
stemmer = Stemmer()
lemmatiz = Lemmatizer()

print("done3")
def stemming(text):
    text = [stemmer.stem(word) for word in text]
    return text

data['Tweet_stemmed'] = data['content'].apply(lambda x: stemming(x))

print("done4")
def lemmatizer(text):
    #print("new Line")
    s="#"
    listof=[]
    for word in text:
        lemi= lemmatiz.lemmatize(word)
        if(s in lemi):
            lemi=lemi.split("#")[1]

        listof.append(str(lemi))
    return listof

data['Tweet_lemmatized'] = data['Tweet_stemmed'].apply(lambda x: lemmatizer(x))

done3
done4
```

نمونه خروجی :

In [208]: data['Tweet_lemmatized'][0]

```
Out[208]: ['گزار',
'خبرگزار',
'شیس',
'نقل',
'آژانس',
'خبر',
'باخ',
'قول',
'اردو',
'اٹل',
'جنوب',
'افغانستان',
'نشر',
'خبرنامه',
'اعلا',
'منسوبین',
'قول',
'اردو',
'ده',
'کیلوگرا',
'مواد',
'متفجره',
'ده',
'جنگ',
'افزار',
'سبک',
'سنگین',
'سو',
'جنگجو',
'گر',
'وه',
'طالب',
'یک',
'...']
```

بخش سوم:

در این بخش به پیاده سازی heap و zeap میپردازیم اما در ابتدا تابعی نوشته شده تحت عنوان create_posting_list

```
#Creating posting List using Dictioanry

def create_postngs_list(data):
    frequency={}
    diction={}

    for i in range(len(data)) :
        for words in data[i]:
            if (words not in frequency):
                frequency[words] = 0
                frequency[words] = frequency[words] + 1
            else:
                frequency[words] = frequency[words] + 1

            if words not in diction.keys():
                diction[words] = set()
                diction[words].add(i)
            else:
                diction[words].add(i)
    return frequency, diction

frequency5,diction_sample5=create_postngs_list(dat_forposting_sample5)
```

که طبق آنچه در صورت پروژه آمده بود به صورت dic ذخیره شده است تا سرعت بالا باشد در این جا پستینگ لیست مورد نظر ساخته میشود به این صورت که اگر ورد در آن نباشد آن ترم اضافه گشته و اگر موجود باشد شماره ترم آن ذخیره گشته است .

برای محاسبه freq کد زیر مورد استفاده قرار گرفت

```
#Zipf law
import seaborn as sns
import matplotlib.pyplot as plt
def frequency(diction):
    listoffrequency=[]
    #print value
    lengths = [len(v) for v in diction.values()]
    return lengths

list_frequency_first_sample=frequency(diction_sample5)
```

که در آن برای هر کلمه تعداد تکرار آن محاسبه شده است

برای مثال برای لیست 15000 تایی و 5000 تایی در کل به تعداد زیر توکن موجود میباشد

```
] : ▶ print(sum(list_frequency_first_sample),sum(list_frequency_secound_sample))
634906 1894815
```

```
04]: ▶ term_freq_df
```

```
rt[104]: [41808,
22139,
2073,
10089,
352,
12363,
536,
971,
782,
3,
4933,
1118,
913,
13,
13952,
5,
6165,
2,
2621,
```

برای کل دیتاست

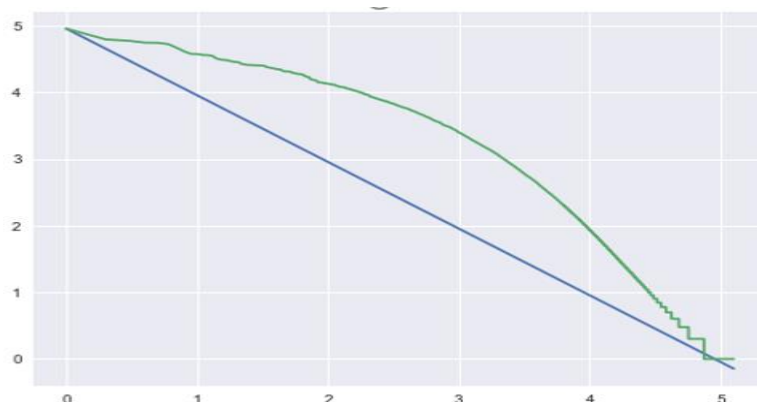
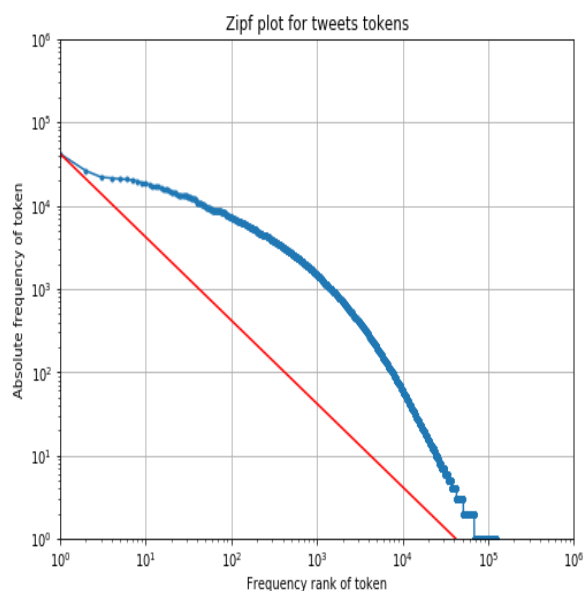
در کل برای هر دو قانون حالت دوم بهتر از حالت اول است بخاطر اینکه بهتر نرمالیز شده و بهتر tokenize انجام شده و همچنین ریشه یابی شده است این باعث میشود که تعداد توکن های منحصر به فرد کمتر شود که بسیار برای ما بهتر است و میتوان دید که هر دو قانون درست میباشند

قانون Zipf یک قانون در مورد توزیع فراوانی کلمات در زبان است (یا در مجموعه ای که به اندازه کافی بزرگ است به طوری که نماینده زبان باشد). برای نشان دادن قانون Zipf اجازه دهید فرض کنیم که ما یک مجموعه داریم و اجازه می دهیم در مجموعه کلمات منحصر به فرد (دایره لغات) باشد.

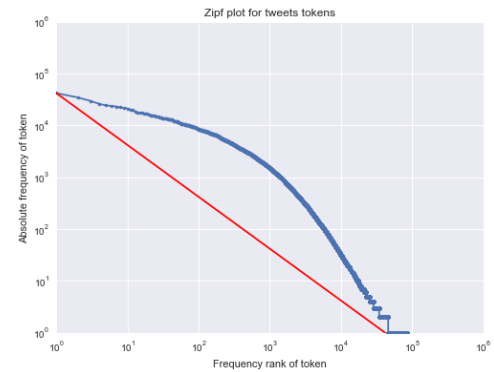
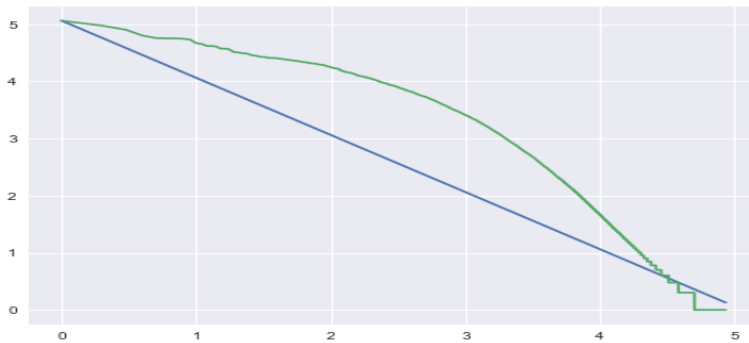
برای پیاده سازی آن همانطور که در کد نیز مشاهده میشود شکل زیر بدست آمد

نمودار ها با seaborn کشیده شده اند

نمودار حالت 1



نمودار در حالت دوم



در این بخش نیز heap پیاده سازی شده است

```
import math as m
def Heap_law(list_sample_5,list_sample_10):
    print(m.log(list_sample_5/list_sample_10,10))
    print(m.log(1/3,10)/m.log(list_sample_5/list_sample_10,10))

    b=(m.log(1/3,10))/(m.log(list_sample_5/list_sample_10,10))
    print(b)

    t =15000/(pow(sumof,b))
    return b,t

b,k=Heap_law(sum(list(frequency5.values())),sum(list(frequency15.values())))
print(b,k)
```

که در اینجا دو زیرمجموعه 5000 تایی و 15000 تایی به صورت رندوم با تابع

```
import random
def random_sample(x,data):
    p=random.sample(range(len(data)), x)
    listofsample=[]
    for i in range(len(p)):
        print(p[i])
        listofsample.append(data[p[i]])
    return listofsample

list_sample_5=random_sample(5000,data['content'])
print(" first random sample done ")
list_sample_10=random_sample(15000,data['content'])
print(" secound random sample done ")
dat_forposting_sample5=pd.Series(list_sample_5)
dat_forposting_sample10=pd.Series(list_sample_10)
```

```
34694
40433
```

ساخته شده که با توجه به عدد رندوم انتخاب میشود که چه داکيومنت هایی برای posting list انتخاب شوند

```
import math as m
def Heap_law(list_sample_5,list_sample_10):
    print(m.log(list_sample_5/list_sample_10,10))
    print(m.log(1/3,10)/m.log(list_sample_5/list_sample_10,10))

    b=(m.log(1/3,10))/(m.log(list_sample_5/list_sample_10,10))
    print(b)

    t =15000/(pow(sumof,b))
    return b,t

b,k=Heap_law(sum(list_frequency_first_sample()),sum(list_frequency_secound_sample))
print(b,k)
```

در اینجا نیز مقدار بهینه b و k با توجه به دو معادله دو مجهول بدست می آید

