

# تمرین اول درس پردازش زبان طبیعی

نام تهیهکننده: ملیکا نوبختیان

نسخه: ۱

# ۱- بخش تئوري

## 1-1- سوال اول

#### ۱-۱-۱ كتابخانه NLTK

کتابخانه NLTK یا natural language toolkit یک کتابخانه ضروری برای کاربردهای پردازش زبان طبیعی tokenization و classification, stemming, tagging, parsing, semnatic reasoning و است که کارهایی مانند NLTK را می توان ابزار اصلی در کاربردهای NLP و ML دانست. هر چند این را به زبان python انجام می دهد. NLTK را می توان ابزار اصلی در کاربردهای و NLP و است. اگر از این کتابخانه کاربردهای زیادی دارد اما باید قبول کرد استفاده از آن برای پردازش زبان طبیعی دشوار است. اگر از این کتابخانه به درستی استفاده نشود بسیار کند خواهد بود و در کاربردهای production جایی نخواهد داشت. پس باید به روش استفاده از آن دقت شود.

#### ۱-۱-۲ کتابخانه Gensim

Gensim یک کتابخانه به زبان python است که به طور ویژه به مشخص کردن Gensim یک کتابخانه می تواند پیکرههای دو سند از طریق ابزارهای topic modeling و vector space modeling می پردازد. این کتابخانه می تواند پیکرههای متنی بزرگ را با کمک بهینگی data streaming و الگوریتمهای incremental مدیریت کند. ویژگی بسیار مهم و برتر این کتابخانه که نسبت به بقیه نمونهها بی نظیر است memory usage optimization و processing speed آن است.

## ۳-۱-۱- کتابخانه spaCy

این کتابخانه نسبت به بقیه کتابخانهها جدیدتر است و برای کاربردهای production طراحی شدهاست. این کتابخانه سریعترین syntactic parser که در حال حاضر موجود است را ارائه می دهد. به علاوه چون این کتابخانه به زبان Cython نوشته شده است بسیار پرسرعت و بهینه است. از معایب این کتابخانه این است که برخلاف کتابخانههای دیگر تعداد کمتری از زبانها را پشتیبانی می کند.

## **1-1-** سوال دوم

الف) ( $\{9,13\}[9-0]989+//\{9,11\}[9-0]989[0-9][9,11]$ ) نكته: شماره عادى را ۱۳ شماره در نظر گرفتم و براساس آن حداقل را مشخص كردم

ب)

 $\begin{tabular}{ll} $$ $$ (0[1-9]|[12][0-9]|3[01])-(0[1-9]|1[012])-(19|20)\d $$ $$ ((http)\:\/\/)?[a-zA-Z0-9\.\?\:@\-_=\#]+\.(ir|org) (= -2.4.20-9).\] $$$ 

#### $([0-9]{3})[A-Z]([0-9]{3})IR([0-9]{2})$ (3

#### **٦-٣- سوال سوم**

در زبانهایی مانند زبان چینی که بین کلمات فاصله وجود ندارد، tokenize کردن متون و پیداکردن کلمات متن دشوار خواهد بود زیرا بدون فاصله عملا هیچ حرفی را از بقیه نمی توان جدا دانست. اگر زبانهایی مانند انگلیسی نیز این فاصله را حذف کنیم باز هم تشخیص کلمات دشوار خواهد بود. به این مثال دقت کنید:

This is insane درست است که شاید مغز ما بتواند جمله This is insane را تشخیص دهد، اما کامپیوترها چطور به این موضوع پی خواهند برد؟ در اینجا word segmentation با استفاده از الگوریتم خواهند برد؟ در اینجا می کند. این الگوریتم در قدمهای زیر کار را انجام می دهد:

- ۱) با اولین کاراکتر رشته دادهشده کار را آغاز میکند.
- ۲) دنبال بزرگترین کلمه موجود در مرجع کلمات که در دست داریم می گردد که با این کاراکتر آغاز شدهباشد.
- ۳) اگر کلمه مشابهی پیدا شد این کلمه به لیست کلمات اضافه می شود و مرز شروع دوباره تعیین می شود ولی اگر کلمه ای پیدا نشد با همان کاراکتر به عنوان یک کلمه رفتار می شود.
  - ۴) دوباره به مرحله ۲ برمی گردیم با این تفاوت که از مرزی که تعیین شدهاست شروع می کنیم.

با مثال Thisisinsane پیش میرویم. از T شروع می کنیم. اگر پیش برویم بزرگ ترین و تنها کلمه ای که با t شروع می شود و در مرجع کلمات وجود دارد t this خواهد بود. حالا از t به کار خود ادامه می دهیم کلمه بعدی که پیدا خواهیم کرد t خواهد بود. باز هم به t می می می می خوریم t و t insane اما چون insane طولانی تر است به عنوان کلمه موجود انتخاب خواهد شد. پس کلمات این جمله شامل t this, is خواهد بود.

# ۲- بخش عملی

# 1-**۲-** سوال اول

برای این بخش و تشخیص اسامی سه الگو کلی داشتیم: اسامی که با Dr. شروع می شوند، اسامی که با regex می شوند و اسامی که با Doctor شروع می شوند و اسامی که با M. D. پایان می یابند. برای چک کردن این سه الگو به صورت زیر های زیر را در نظر گرفتم:

برای اینکه مواردی که در تمرین گفته شده بود را چک کنم، نمونهها را در یک فایل txt ریختم و برای خواندن نمونهها کافی است آدرس فایل مورد نظر را به آن بدهیم که برای من این فایل doctor\_names.txt خواهد بود. خروجی به این صورت خواهد بود که هر الگو چاپ خواهد شد و سپس مشخص می شود آیا این موارد مطابق regex ها هستند یا نه. قطعه کد زیر این کار را انجام خواهد داد:

خروجی به صورت زیر است:

## **۲-۲-** سوال دوم

برای اینکه نتیجه متدهای word\_tokenize و sent\_tokenize را ببینیم من از یک فایل با متن نمونه به نام دهد. داد تا این کار را روی آن انجام دهد. نام comp.txt استفاده کردم. البته به طور کلی می توان آدرس فایل دلخواه را داد تا این کار را روی آن انجام دهد.

```
file_address = input("Enter File Address: ")
         file1 = open(file_address, 'r')
         text = str(file1.read())
Text: Wil\words = word_tokenize(text)
Matched Pprint("***** Tokenized words with word_tokenize *****")
//////// print("")
Text: Pames = sent_tokenize(text)
Matched Pifor sent in sents:
                                                                                        rn Successfully! True
                                    Text: Dr. Etienne-Emile Baulieu
Matched Pattern Successfully! True
                                      Matched Pattern Successfully! True
                                                                          Text: Doctor Dolittle
                                                                          Matched Pattern Successfully! True
                                      Text: Dr. Karl Thomae
                                      Matched Pattern Successfully! True
Matched Pattern Successfully! True
                                      Text: Dr. Alan D. Lourie
                                      Matched Pattern Successfully! True
                                      Text: Dr. Xiaotong Fei
                                      Matched Pattern Successfully! True
```

قطعه کد زیر فایل مورد نظر را میخواند و کلمات و جملات آن را با استفاده از متدهای گفته شده به دست میآورد و چاپ میکند:

خروجی برای فایل گفتهشده به صورت زیر خواهد بود. ابتدا بخشی از کلماتی که در این فایل موجود است را مشاهده می کنید:

و جملات آن نیز به صورت زیر هستند:

## ٣-٢- سوال سوم

برای این سوال می توانیم از back reference ها در regex استفاده کنیم. با آنها می توانیم یک گروه در برای این سوال می توانیم این گروه با شکل 1, 2, 1 نشان داده می شود. برای حذف حروف تکراری regex او تعدیل که در زیر می بینید استفاده شده است. 2 ( $|w\rangle$ ) نشان می دهد یک کاراکتر آمده و همان کاراکتر که گروه ۲ است دوباره تکرار شده است. برای حذف یک حرف تکراری در هر مرحله از  $|x\rangle$  استفاده شده است. در این جا می بینیم که یک گروه ۲ تکراری حذف شده است. این کار را تا زمانی انجام می دهیم که به کلمه ای برسیم

که در words در nltk موجود باشد:

یک نمونه و جواب را در زیر میبینید:

## ۴-۲- سوال چهارم

Word should be normalized:looooove normal word: love

#### ابتدا فایلهای موردنظر را خواندم و متن آنها را در متغیرهای متناظرشان ذخیره کردم:

```
english_example_Address = "AlbertEinstein.txt"
file1 = open(english_example_Address, 'r')
english_example = str(file1.read())

persian_example_Address = "Shahnameh.txt"
file1 = open(persian_example_Address, 'r')
persian_example = str(file1.read())

english_short_text_Address = "ShortSampleEnglish.txt"
file1 = open(english_short_text_Address, 'r')
english_short = str(file1.read())

persian_short_text_Address = "ShortSamplePersian.txt"
file1 = open(persian_short_text_Address, 'r')
persian_short = str(file1.read())
```

سپس از TreebankWordTokenizer استفاده کردم تا token های هر متن و type آنها را به دست آورم. منظور از type های در اینجا توکنهای منحصر به فردی است که پس از عمل tokenize به وجود می آیند. قطعه کد زیر نحوه کار با این tokenizer است:

```
tt = TreebankWordTokenizer().tokenize(english_example)
print(f'** Token count in AlbertEinstein.txt with Treebankwordtokenizer : {len(tt)}')
print(f'** Token types in AlbertEinstein.txt with Treebankwordtokenizer : {", ".join(list(set(tt)))}')

tt = TreebankWordTokenizer().tokenize(persian_example)
print(f'** Token count in Shahnameh.txt with Treebankwordtokenizer : {len(tt)}')
print(f'** Token types in Shahnameh.txt with Treebankwordtokenizer : {", ".join(list(set(tt)))}')
print("")

tt = TreebankWordTokenizer().tokenize(english_short)
print(f'** Token count in ShortSampleEnglish.txt with Treebankwordtokenizer : {", ".join(list(set(tt)))}')
print(f'** Token types in ShortSampleEnglish.txt with Treebankwordtokenizer : {", ".join(list(set(tt)))}')
print(f'** Token count in ShortSamplePersian.txt with Treebankwordtokenizer : {len(tt)}')
print(f'** Token count in ShortSamplePersian.txt with Treebankwordtokenizer : {", ".join(list(set(tt)))}')
print(f'** Token types in ShortSamplePersian.txt with Treebankwordtokenizer : {", ".join(list(set(tt)))}')
print(f'** Token types in ShortSamplePersian.txt with Treebankwordtokenizer : {", ".join(list(set(tt)))}')
```

برای متن AlbertEinstein.txt خروجی به صورت زیر است. در اینجا تنها بخشی از token type ها قابل مشاهده

```
** Token count in AlbertEinstein.txt with Treebankwordtokenizer : 250
** Token types in AlbertEinstein.txt with Treebankwordtokenizer : Berlin., the, a, Albert, time, Professor, Federal, Kaiser, in, Princeton*., technical, Berlin, position, his,
```

#### برای متن شاهنامه نیز خروجی به شکل زیر است:

هستند:

<sup>\*\*</sup> Token count in Shahnameh.txt with Treebankwordtokenizer : 2207

<sup>,</sup> بهر, باتوت, لرزان, بندهي, گونه, ماهست, خواند, نوبهنو, جويد, پيششبه, ژرفي, هوش, خويشبيگاته, شفود, پرورد : Token types in Shahnameh.txt with Treebankwordtokenizer \*\*

برای متن کوتاه انگلیسی خروجی به شکل زیر است: و برای متن کوتاه فارسی هم خروجی به شکل زیر خواهد بود:

Token count in ShortSampleEnglish.txt with Treebankwordtokenizer : 19

- \*\* Token count in ShortSamplePersian.txt with Treebankwordtokenizer : 16
- خوب ر باشید. ر یک ر کوتاه ر ما ر امیدوارم ر مثن ر تعرین ر ا ر این ر است ر است ر است این ر . ر سلام ۱ ر نعونه : Token types in ShortSamplePersian.txt with Treebankwordtokenizer \*\*

برای به دست آوردن کلمات متن کوتاه انگلیسی با استفاده از RegexpTokenizer از عبارت با قاعده به صورت زیر استفاده کردم که توکنهای متن از جمله کلمات و اعداد را به ما خواهد داد :

```
regexp = RegexpTokenizer('\w+').tokenize(english_short)

print(f'ShortSampleEnglish.txt tokens with Regexp : {", ".join(regexp)}')

و نتیجه آن به صورت زیر خواهد بود:
```

ShortSampleEnglish.txt tokens with Regexp : Hello, Hope, you, re, doing, well, It, is, a, sample, short, text, This, is, our, 1, st, assignment

اگر بخواهیم تنها کلمات چاپ شوند و اعداد در این قسمت نباشند باید از کد زیر استفاده کنیم:

```
regexp = RegexpTokenizer('[a-zA-Z]+').tokenize(english_short)
print(f'ShortSampleEnglish.txt words with Regexp : {", ".join(regexp)}')
```

و نتیجه آن نیز به صورت زیر است:

ShortSampleEnglish.txt words with Regexp : Hello, Hope, you, re, doing, well, It, is, a, sample, short, text, This, is, our, st, assignment

برای اینکه بتوانیم اعداد متن انگلیسی را به دست آوریم به شکل زیر عمل می کنیم:

```
regexp = RegexpTokenizer('[0-9]+').tokenize(english_example)
print(f'AlbertEinstein.txt numbers with Regexp : {", ".join(regexp)}')
```

و نتیجه به صورت زیر است:

AlbertEinstein.txt numbers with Regexp : 14, 1879, 1896, 1901, 1905, 1908, 1909, 1911, 1914, 1914, 1933, 1940, 1945

برای به دست آوردن توکنهای متن کوتاه فارسی هم به همان شکلی که در ابتدا برای تشخیص توکنها متن انگلیسی عمل کردیم پیش میرویم و خروجی آن نیز به شکل زیر است:

سلام, اسيدوارم, خوب, باشيد, اين, يک, متن, کوتاه, نمونه, است, اين, تعرين, ۱, ما, است: ShortSamplePersian.txt tokens with Regexp

نتیجه استفاده از WhiteSpaceTokenizer برای متن کوتاه انگلیسی به شکل زیر است:

ShortSampleEnglish.txt words with WhiteSpaceTokenizer : Hello!, Hope, you're, doing, well., It, is, a, sample, short, text.This, is, our, 1, st, assignment.

این tokenizer به این شکل عمل می کند که می تواند token های یک متن را بدون هیچ فاصله، newline یا tokenizer به دست آورد. به عبارت دیگر برای جداسازی token ها از این موارد استفاده می کند. اگر خواسته باشیم با regexp چنین کاری انجام دهیم می توانیم به شکل زیر عمل کنیم:

هنگامی که در gaps از gaps استفاده می کنیم، این tokenizer به این شکل عمل می کند که سعی می کند با

```
tokenizer = RegexpTokenizer('\s+', gaps=True).tokenize(english_short)
print(f'ShortSampleEnglish.txt words with RegexpTokenizer same as WhiteSpaceTokenizer : {", ".join(tokenizer)}')
```

regex داده شده token های متن را جدا کند. در اینجا سعی می کند با فاصله به هر تعداد این کار را انجام دهد.

ShortSampleEnglish.txt words with RegexpTokenizer like WhiteSpaceTokenizer : Hello!, Hope, you're, doing, well., It, is, a, sample, short, text.This, is, our, 1, st, assignment.

نتیجه استفاده از آن به شکل زیر است:

واضح است که خروجی دو حالت برابر است.

نتیجه استفاده از WordPunctokenizer به شکل زیر است:

ShortSampleEnglish.txt words with WordPunctTokenizer: Hello, !, Hope, you, ', re, doing, well, ., It, is, a, sample, short, text, ., This, is, our, 1, st, assignment,

این tokenizer به این شکل عمل می کند رشتههای کاراکتری الفبایی و هم چنین غیر الفبایی را به عنوان tokenizer قبول می کند. رشتههای الفبایی کلمات معمولی هستند که در متن استفاده می شود اما رشته غیر الفبایی شامل رشته یا کاراکترهایی می شود کاراکترهای دیگر غیر الفبا مانند ,' ,. , ? ,! و ... را در خود دارند. با توجه به توکنهای استخراج شده از متن کوتاه انگلیسی متوجه می شویم که . ,' ,! به عنوان توکن شناخته شده اند در حالی که در token های قبل چنین حالتی نداشتیم و این موارد به قبل یا بعد خود متصل می شدند و به عنوان یک tokenizer شناخته می شدند.

## **۵-۲-** سوال پنجم

اگر PorterStemmer و LancasterStemmer را روی توکنهای گفته شده اعمال کنیم، نتیجه زیر را خواهیم گرفت:

```
*** Original Words ***:
was, Germany, weeks, later, family, moved, Italy

** Stemmed with PorterStemmer **:
wa, germani, week, later, famili, move, itali

** Stemmed with LancasterStemmer **:
was, germany, week, lat, famy, mov, ita
```

همان طور که در تصویر واضح است نتیجه استفاده از این دو stemmer به جز در کلمه weeks متفاوت است. اما حالا به بررسی کلی نتایج این دو stemmer میپردازیم. اگر دقت کنیم میبینیم که بیشتر کلماتی که توسط porter اما stem ، Lancaster شدهاند بیشتر از حالت اولیه خود فاصله گرفتهاند و به ریشه اصلی نزدیک شدهاند. اما stemmer سختگیرتر بیشتر شکل اصلی کلمات را حفظ کردهاست. به عبارت دیگر میتوان گفت stemmer یک stemmer سختگیرتر نسبت به porter است. در حالی که نتیجه porter به چیزی که برای ما آشنا است نزدیک تر است.

می توانیم این طور بگوییم که اگر فشرده سازی حداکثری برای کلماتمان می خواهیم Lancaster گزینه بهتری است در حالی که Porter نتیجه ملموس تر و نزدیک تر به کلمه اصلی را به ما خواهد داد.

اگر از lemmtizer بر روی کلماتی که گفته شدهاست استفاده کنیم نتیجه زیر را خواهیم گرفت:

```
*** Lemmatize given Words ***:
wave, fishing, rock, wa, corpus, better, ate, broken
```

همان طور که مشخص است تنها waves, rocks و corpura به حالت اولیه و درست خود بازگشتهاند و بقیه موارد این چنین نیستند. حتی was هم به اشتباه به wa تبدیل شدهاست که نشان می دهد lemmatizer به اشتباه wasرا به عنوان یک کلمه جمع درنظر گرفته است و s آن را حذف کردهاست.

اما راه حل چیست؟ برای حال این مشکل هنگام lemmatize کردن باید مشخص کنیم نوع هر کلمه چیست برای مثال اسم، فعل، صفت و ... است. در این صورت به نتیجه درست خواهیم رسید. Lemmatizer در حالت عادی همه کلمات اسم در نظر می گیرد:

```
print("*** Lemmatize with Type ***:")
print(f'word : fishing, lemm : {lemmatizer.lemmatize("fishing", pos="v")}')
print(f'word : was, lemm : {lemmatizer.lemmatize("was", pos="v")}')
print(f'word : better, lemm : {lemmatizer.lemmatize("better", pos="a")}')
print(f'word : ate, lemm : {lemmatizer.lemmatize("ate", pos="v")}')
print(f'word : broken, lemm : {lemmatizer.lemmatize("broken", pos="v")}')
```

و نتیجه صحیح lemmatize این کلمات به شکل زیر است:

```
*** Lemmatize with Type ***:
word : fishing, lemm : fish
word : was, lemm : be
word : better, lemm : good
word : ate, lemm : eat
word : broken, lemm : break
```

Fishing, was, ate و broken از نوع verb هستند و Fishing, was, ate

#### ۶-۲- سوال ششم

#### برای بررسی تغییرات به وجودآمده ۳ توییت در نظر گرفتم که محتوای اصلی آنها به شکل زیر است:

Original Tweets:

["RT @Claire\_FOX5: #BREAKING: @GaSecofState's office confirms Floyd County has found 2,600 ballots during audit. Says Sec. Raffensperger waâ $\in$ '"]

//////

['Really disgusting that the failing New York Times allows dishonest writers to totally fabricate stories.']

//////

['What an INCREDIBLE opening by @JesseBWatters. I will put it up later for everyone to see!']

مرحله اول حذف فاصلههای تکراری است. برای این کار از regex استفاده میکنیم و هر جا این الگو تکراری پیدا شد تنها با یک فاصله جایگزین میکنیم:

```
# Step 1:
# Remove Duplicate Spaces between tweet's words
for i, tweet in enumerate(tweets):
    new_tweet = re.sub(r'\s+', ' ', str(tweet))
    tweets[i] = new_tweet
```

#### نتیجه به صورت زیر است:

Step 1: Remove Duplicate Spaces between tweet's words ["RT @Claire\_FOX5: #BREAKING: @GaSecofState's office confirms Floyd County has found 2,600 ballots during audit. Says Sec. Raffensperger waâ $\in$ |"]

['Really disgusting that the failing New York Times allows dishonest writers to totally fabricate stories.'] /////

['What an INCREDIBLE opening by @JesseBWatters. I will put it up later for everyone to see!']

در توییت اول یک فاصله اضافی بین BREAKING# و ...GaS... @ مشاهده می شد که حذف شده است. در مرحله دوم تمام حروف باید به حروف کوچک تبدیل شوند:

```
# Step 2:
# Make text's words Lower case
for i, tweet in enumerate(tweets):
    new_tweet = tweet.lower()
    tweets[i] = new_tweet
```

#### نتیجه اعمال این قسمت به شکل زیر است:

Step 2: Make text's words Lower case

["rt @claire\_fox5: #breaking: @gasecofstate's office confirms floyd county has found 2,600 ballots during audit. says sec. raffensperger wa…"]

['really disgusting that the failing new york times allows dishonest writers to totally fabricate stories.'] /////

['what an incredible opening by @jessebwatters. i will put it up later for everyone to see!']

در مرحله سوم username ها حذف می شوند. برای این قسمت باز هم از regex استفاده می کنیم و الگوهایی که با @ شروع می شوند را حذف می کنیم:

```
# Step 3:
# Remove Handles ( Usernames ) from Tweets
for i, tweet in enumerate(tweets):
    new_tweet = re.sub('@[\w]+', '', str(tweet))
    tweets[i] = new_tweet
```

#### نتیجه اعمال این قدم به شکل زیر است:

```
Step 3: Remove Handles ( Usernames ) from Tweets ["rt : #breaking: 's office confirms floyd county has found 2,600 ballots during audit. says sec. raffensperger waâ\in|"] ///// ['really disgusting that the failing new york times allows dishonest writers to totally fabricate stories.'] ///// ['what an incredible opening by . i will put it up later for everyone to see!']
```

در این مرحله علائم نگارشی و کاراکترهای خاص و اعداد باید از توییتها حذف شوند که این قسمت هم به عهده regex است. هر کلمه یا عبارتی که با کاراکترهای مورد قبول شروع نشود حذف خواهد شد. بعضی از کاراکترها مانند # ها نیز نباید حذف شوند چون به وجود آنها نیاز خواهیم داشت:

```
# Step 4:
# Remove punctuation, special characters and numbers from tweets
for i, tweet in enumerate(tweets):
    new_tweet = re.sub("[^a-z#&%\$]+", '', str(tweet))
    tweets[i] = new_tweet
```

#### نتیجه این قسمت را در شکل زیر میبینیم:

Step 4: Remove punctuation, special characters and numbers from tweets rt #breaking s office confirms floyd county has found ballots during audit says sec raffensperger wa

really disgusting that the failing new york times allows dishonest writers to totally fabricate stories  $^{\prime\prime\prime\prime\prime\prime\prime}$ 

what an incredible opening by i will put it up later for everyone to see برای این کار استفاده کردم. زیرا و WhiteSpaceTokenizer میرویم. من از tokenizer برای این کار استفاده کردم. زیرا فاصله ها را حذف می کند و مواردی که می ماند را به ما می دهد. اگر از tokenizer ای مانند WordPunct استفاده

```
# Step 5:
# Apply Tokenization

for i, tweet in enumerate(tweets):
    new_tweet = WhitespaceTokenizer().tokenize(tweet)
    tweets[i] = new_tweet
```

می کردم، این tokenizer کاراکترهای خاص مانند # که اتفاقا در مراحل بعد به آن نیاز داریم را از کلمات جدا WhiteSpaceTokenizer می کرد که باعث مشکل می شد. برای همین بهترین انتخاب این قسمت به انتخاب من

نتیجه این قسمت را مشاهده می کنید:

```
# Step 6:
# Removing Stop words
stop_words = set(stopwords.words('english'))
for i, tweet in enumerate(tweets):
    new_tweet = [w for w in tweet if not w in stop_words]
    tweets[i] = new_tweet
```

#### نتیجه اجرا این قسمت را میبینیم:

```
Step 6: Removing Stop words
['rt', '#breaking', 'office', 'confirms', 'floyd', 'county', 'found',
'ballots', 'audit', 'says', 'sec', 'raffensperger', 'wa']
/////
['really', 'disgusting', 'failing', 'new', 'york', 'times', 'allows',
'dishonest', 'writers', 'totally', 'fabricate', 'stories']
/////
```

['incredible', 'opening', 'put', 'later', 'everyone', 'see'] در این مرحله باید کلمات با طول کمتر از ۳ را حذف کنیم. اکثر این کلمات هم مانند stopword ها تغییر خاصی در معنی ایجاد نمی کنند و اضافی هستند. در اینجا مواردی مانند rt که در توییت هم موجود است حذف می شوند. با این حال ممکن است بعضی از کلمات با معنی مانند فعل go نیز در این فرآیند حذف شوند که مطابق میل ما نست:

```
# Step 7:
# Remove Words with len less than 3
for i, tweet in enumerate(tweets):
    new_tweet = [w for w in tweet if len(w) > 2]
    tweets[i] = new_tweet
```

#### نتیجه را میبینیم:

مرحله بعد stemming است. به توضیح بیشتر Lancaster و Lancaster میپردازیم. بیشتر کلماتی که توسط stemming میشوند بیشتر از حالت اولیه خود می گیرند و به ریشه اصلی نزدیک میشوند اما تشخیص اصل آنها دشوار میشود. اما porter بیشتر شکل اصلی کلمات را حفظ می کند. به عبارت دیگر می توان گفت stemmer یک stemmer سختگیرتر نسبت به porter است. در حالی که نتیجه porter به چیزی که برای ما آشنا است نزدیک تر است. می توانیم این طور بگوییم که اگر فشرده سازی حداکثری برای کلماتمان می خواهیم است نزدیک تر است در حالی که Porter نیجه ملموس تر و نزدیک تر به کلمه اصلی را به ما خواهد داد. Lancaster شبت به porter سریع تر است. در میان این دو snowball قرار می گیرد که به سخت گیری Lancaster نیست اما مانند porter نیز آنقدر به شکل اصلی کلمات وفادار نیست.

```
# Step 8:
# Stemming
porter = PorterStemmer()
for i, tweet in enumerate(tweets):
    new_tweet = [porter.stem(w) for w in tweet]
    tweets[i] = new_tweet
```

#### نتیجه آن را میبینیم:

#### و WordCloud توییتها به شکل زیر است:

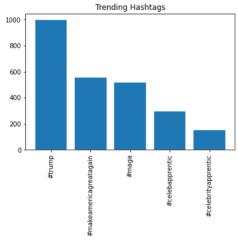


برای شناسایی ترندها یا همان هشتگها به سراغ regex میرویم تا الگوهایی که با # شروع میشوند را شناسایی کنیم و تعداد آنها را مشخص کرده و سپس یک نمودار برای آن رسم کنیم:

```
string_of_All = " ".join(all_words)
trend_reg = re.compile(r"#\w+")
trends = Counter(trend_reg.findall(string_of_All))
trends = trends.most_common(5)
print(trends)
plt.bar([title for title, freq in trends], [freq for title, freq in trends])
plt.title("Trending Hashtags")
plt.xticks(rotation=90)
plt.show()
```

#### نتیجه به شکل زیر است:

[('#trump', 995), ('#makeamericagreatagain', 556), ('#maga', 518), ('#celebapprentic', 295), ('#celebrityapprentic', 151)]



Activate Wir