

```
import numpy as np
```

```
# gathering the consumer key and secret

api_key = "I8R8K78dKM5Cj5lRkrFACFYV"
api_key_secret = "3VoGqgZRPzP4E2BDeAeOXDPdLlXZW6LiJ4HvY0z4uWAAKS1LNMlt"
access_token = "797894924564263830-8EB8SgmQu2WH8gTCUp0zTpwk0U3Ja"
access_token_secret = "Idx6GMS47McOm8BDM57bkAlf5t7QPoEdlYdyPwBlm"
```

```
# searching keyword name

search_words = "#Budget2022"

# searching number of tweets

number_posts = 1000
```

```
# using the OAuth for accessing tweets

auth = tw.OAuthHandler(api_key, api_key_secret)
auth.set_access_token(access_token, access_token_secret)
tw.API(auth, wait_on_rate_limit=True)
```

```
# downloading the tweets

tweets = tw.Cursor(api.search_tweets, q=search_words,
                    lang='en').items(number_posts)

# list to store downloaded data

tweets_text = []

# downloading different attributes

for tweet in tweets:
    tweets_text.append([
        tweet.user.screen_name, tweet.user.name, tweet.user.description,
        tweet.text, tweet.user.location, tweet.user.followers_count,
        tweet.source, tweet.created_at, tweet.user.friends_count
    ])
```

## 2.STORING THE DOWNLOADED DATA

```
# storing tweet data

tweet_data = pd.DataFrame(tweets_text,
                           columns=[
                               "User Name", "Name", "Description (Status)",
                               "Tweets Posted", "Location", "No of Followers", "Device Used for Tweet", "Date of Tweet Creation", "No of Friends"
                           ])

# starting few tweet data

tweet_data
```

|     | User_Name       | Name                | Description (Status)                            | Tweets Posted   | Location                     | No of Followers | Device Used for Tweet | Date of Tweet Creation       | No of Friends |
|-----|-----------------|---------------------|---|---|------------------------------|-----------------|-----------------------|------------------------------|---------------|
| 0   | AdarshPallav    | Adarsh Pallav       |   | RT @satishacharya:<br>Blueprint for 25 years! #Bu_              |                              | 7               | Twitter for Android   | 2022-02-07<br>07:11:43+00:00 | 38            |
| 1   | bnsidia         | Business Standard   | Latest news on the economy, companies, markets. | #BSMorningShow   Will #RBI go for a hike in ..                  | India                        | 2187279         | TweetDeck             | 2022-02-07<br>07:11:00+00:00 | 430           |
| 2   | bnsidia         | Business Standard   | Latest news on the economy, companies, markets. | #BSMorningShow   Will #RBI go for a hike in rev..               | India                        | 2187279         | Twitter Web App       | 2022-02-07<br>07:10:43+00:00 | 430           |
| 3   | bnsidia         | Business Standard   | Latest news on the economy, companies, markets. | #BSMorningShow   Business Standard's @anuproyt..                | India                        | 2187279         | TweetDeck             | 2022-02-07<br>07:10:00+00:00 | 430           |
| 4   | orfonline       | ORF                 | Non-partisan independent analyses on security.. | Tax concessions for corporate co-operatives are ..              | India                        | 106444          | TweetDeck             | 2022-02-07<br>07:10:00+00:00 | 136           |
| --  | --              | --                  | --  | --  | --                           | --              | --                    | --                           | --            |
| 995 | VikramB12singh  | विक्रम बाबुदुर सिंह | सहोदक अचार्य (Assistant Professor)              | RT @IncomeTaxIndia: Budget Highlights vKey DI..                 | जंजीरबाँसा कस्बिलाल, भरतखण्ड | 467             | Twitter for Android   | 2022-02-06<br>16:54:20+00:00 | 2205          |
| 996 | TheRationalDesi | Rational One        |   | #Budget2022 Congress #BIPhasoDeshBachao #PU_                    |                              | 0               | Twitter Web App       | 2022-02-06<br>16:54:16+00:00 | 58            |
| 997 | VikramB12singh  | विक्रम बाबुदुर सिंह | सहोदक अचार्य (Assistant Professor)              | RT @cbic_india: Key highlights of Union Budget for a hike in .. | जंजीरबाँसा कस्बिलाल, भरतखण्ड | 467             | Twitter for Android   | 2022-02-06<br>16:51:46+00:00 | 2205          |

```

999 Vibhanshu28836256 Vibhanshu
RT @thebubblebust1: My
state on Financial Budg... Seoni, India 79 Twitter for
Android 1651:16+00:00
1651:16+00:00

1000 rows x 9 columns

# collected data information

tweet_data.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1000 entries, 0 to 999
Data columns (total 9 columns):
 #   Column                Non-Null Count  Dtype  ---
  0   User_Name             1000 non-null   object
  1   Name                  1000 non-null   object
  2   Description (Status)  1000 non-null   object
  3   Tweets Posted         1000 non-null   object
  4   Location              1000 non-null   object
  5   No of Followers       1000 non-null   int64
  6   Device Used for Tweet  1000 non-null   object
  7   Date of Tweet Creation 1000 non-null   datetime64[ns, UTC]
  8   No of Friends         1000 non-null   int64
dtypes: datetime64[ns, UTC](1), int64(2), object(6)
memory usage: 70.4+ KB

# saving the data as csv file

tweet_data.to_csv(
'\\DOCUMENTS\\COLLEGE\\CLASSES\\EXPERIMENT_NO_2\\tweet_data_preprocess.csv')

3.PERFORMING PREPROCESSING TASK

# importing the regex library

import re

# Function for removing the urls & special characters

import re

def special_remove(string):
    temp = ""
    for character in string:
        if (character.isalnum() or character == ' '):
            temp += character
    return temp

```

```

for i in range(len(tweet_data.columns)):
    if tweet_data[tweet_data.columns[i]].dtype == np.object_:
        tweet_data[tweet_data.columns[i]].apply(lambda x: re.sub(
            r"https://.*?[/\?&#]", '', x, flags=re.MULTILINE))
        tweet_data[tweet_data.columns[i]] = tweet_data[
            tweet_data.columns[i]].apply(special_remove)

# after the removal of url's & special characters starting few data
tweet_data.head()

```

|   | User Name    | Name              | Description (Status)                              | Tweets Posted                                      | Location | No of Followers | Device Used for Tweet | Date of Tweet Creation   | No of Friends |
|---|--------------|-------------------|---|--|----------|-----------------|-----------------------|--------------------------|---------------|
| 0 | Adarsh Palla | Adarsh Palla      | RT satishacharya Blueprint for 25 years Budget.   |  |          | 7               | Twitter for Android   | 2022-02-07 07:1143+00:00 | 38            |
| 1 | bsindia      | Business Standard | Latest news on the economy companies markets p... | BSMShowingWill Will Rill go for a hike in rever... | India    | 2187279         | TweetDeck             | 2022-02-07 07:1100+00:00 | 430           |
| 2 | bsindia      | Business Standard | Latest news on the economy companies markets p... | BSMShowingWill Will Rill go for a hike in rever... | India    | 2187279         | Twitter Web App       | 2022-02-07 07:1043+00:00 | 430           |
| 3 | bsindia      | Business Standard | Latest news on the economy companies markets p... | BSMShowingWill Will Rill go for a hike in rever... | India    | 2187279         | TweetDeck             | 2022-02-07 07:1000+00:00 | 430           |
| 4 | orfonline    | ORF               | Nonpartisan independent analyses on security s... | Tax concessions for cooperative societies is a...  | India    | 106444          | TweetDeck             | 2022-02-07 07:1000+00:00 | 136           |

### STEP 1: TOKENIZATION

```

# importing the regex library
import re

# installing the libraries
pip install nltk

Requirement already satisfied: nltk in c:\users\h\anaconda3\lib\site-packages (3.6.5)
Requirement already satisfied: click in c:\users\h\anaconda3\lib\site-packages (from nltk) (8.0.3)
Requirement already satisfied: joblib in c:\users\h\anaconda3\lib\site-packages (from nltk) (1.1.0)
Requirement already satisfied: regex>=2021.8.3 in c:\users\h\anaconda3\lib\site-packages (from nltk) (2021.8.3)
Requirement already satisfied: tqdm in c:\users\h\anaconda3\lib\site-packages (from nltk) (4.62.3)
Requirement already satisfied: colorama in c:\users\h\anaconda3\lib\site-packages (from click->nltk) (0.4.4)

# function for removing the url's & special characters

```

```
def url_remove(lst):
    for i in range(len(lst)):
        lst[i] = re.sub(r'\"https?:\\/.*/\"\\n\"',
            '',
            lst[i],
            flags=re.MULTILINE)
    return lst

def special_remove(lst):
    for i in range(len(lst)):
        temp = ''
        for character in lst[i]:
            if character.isalnum():
                temp += character
        lst[i] = temp
    return lst

# importing the libraries

import nltk

# we will perform tokenization for every attribute

tokenized_column_name = [
    'User_Name[Tokenized]',
    'Description [Status][Tokenized]', 'Tweets Posted[Tokenized]',
    'Location[Tokenized]', 'No of Followers',
    'Device Used for Tweet[Tokenized]', 'Date of Tweet Creation[Tokenized]',
    'No of Friends'
]

# creating a dataframe to store tokenized data of each attribute

tokenized_df = pd.DataFrame()
```

```
# performing the tokenized using the nltk.word_tokenize function

for i in range(len(tokenized_column_name)):
    if (tweet_data[tweet_data.columns[i]].dtype == np.object.):
        tokenized_df[tokenized_column_name[i]] = tweet_data.apply(
            lambda row: nltk.word_tokenize(row[tweet_data.columns[i]]), axis=1)

    else:
        tokenized_df[tokenized_column_name[i]] = tweet_data[
            tweet_data.columns[i]]

# after performing tokenization starting few tokenized data
tokenized_df.head()
```

|   | User_Name[Tokenized] | Name[Tokenized]      | Description (Status) [Tokenized]              | Tweets Posted[Tokenized]                          | Location[Tokenized] | No of Followers | Device Used for Tweet[Tokenized] | Date of Tweet Creation[Tokenized] |
|---|----------------------|----------------------|---|---|---------------------|-----------------|----------------------------------|-----------------------------------|
| 0 | (AdarshPallav)       | (Adarsh, Pallav)     | [Latest, news, on the economy, companies, ma, | [RT, satishacharya, Blueprint for, 25,            |                     | 7               | [Twitter, for, Android]          | 2022-07-11T04:30                  |
| 1 | (bsindia)            | (Business, Standard) | [Latest, news, on the economy, companies, ma, | [BSMMorningShow, Will, RBI, go, for, a, hike, I,  | (India)             | 2178279         | [TweetDeck]                      | 2022-07-11T04:30                  |
| 2 | (bsindia)            | (Business, Standard) | [Latest, news, on the economy, companies, ma, | [BSMMorningShow, Will, RBI, go, for, a, hike, in, | (India)             | 2178279         | [Twitter, Web, App]              | 2022-07-10T04:30                  |
| 3 | (bsindia)            | (Business, Standard) | [Latest, news, on the economy, companies, ma, | [BSMMorningShow, Business, Standards,             | (India)             | 2178279         | [TweetDeck]                      | 2022-07-10T04:30                  |

|  |  |  |   |         |        |             |                    |
|--|--|--|---|---------|--------|-------------|--------------------|
|  |  | (Nonpartisan,<br>independent<br>analyses on,<br>secur... | [Tax, concessions,<br>for cooperative<br>societies... | [India] | 106444 | [TweetDeck] | 2022-<br>07/10:00- |
|--|--|--|---|---------|--------|-------------|--------------------|

## STEP 2: STOP WORD REMOVAL

```
# We will perform stop word removal for every attribute

stop_word_removed_column_name = [
    'User_Name(Stop Word Removed)', 'Name(Stop Word Removed)',
    'Description (Status)(Stop Word Removed)',
    'Tweets Retweeted(Stop Word Removed)', 'Location(Stop Word Removed)',
    'No of Followers', 'Device Used for Tweet(Stop Word Removed)',
    'Date of Tweet Creation(Stop Word Removed)', 'No of Friends'
]

# creating a dataframe to store stop word removed data of each attribute

stop_word_removed_df = pd.DataFrame(columns=stop_word_removed_column_name)

# collecting all the 'english' language stop word

eng_stopwords = nltk.corpus.stopwords.words('english')

# stop word list

eng_stopwords

['I',
 'me',
 'my',
 'myself',
 ...]
```

'you're',  
    'you've',  
    'you'll',  
    'you'd',  
    'yours',  
    'yours',  
    'yourself',  
    'yourselves',  
    'he',  
    'him',  
    'his',  
    'himself',  
    'she',  
    'she's',  
    'her',  
    'hers',  
    'herself',  
    'it',  
    'it's',  
    'its',  
    'itself',  
    'they',  
    'them',  
    'their',  
    'theirs',  
    'themselves',  
    'what',  
    'which',  
    'who',  
    'whose',

```

'mightn't',
'mustn',
'mustn't',
'needn',
'needn't',
'shan',
'shan't',
'shouldn',
'shouldn't',
'wasn',
'wasn't',
'weren',
'weren't',
'won',
'won't',
'wouldn',
'wouldn't']

# function for removing stop word

def stop_remove(lst):
    for word in lst:
        if word.lower() in eng_stopwords:
            lst.remove(word)
    return lst

# performing the stop word removal using the stop remove function
```

```

for i in range(len(stop_word_removed_column_name)):
    if (tokenized_df[tokenized_column_name[i]].dtype == np.object_):
        stop_word_removed_df[stop_word_removed_column_name[i]] = tokenized_df[
            tokenized_column_name[i]].copy(deep=True).apply(stop_remove)
    else:
        stop_word_removed_df[stop_word_removed_column_name[i]] = tokenized_df[
            tokenized_column_name[i]]

# after performing stop word removal starting few data

stop_word_removed_df.head()

```

| User Name (Stop Word Removed) | Name (Stop Word Removed) | Description (Status) (Stop Word Removed) | Tweets Posted (Stop Word Removed)                 | Location (Stop Word Removed)                         | No of Followers | Device Used for Tweet (Stop Word Removed) | Date of Tweet Creation (Stop Word Removed) | No of Friends             |     |
|-------------------------------|--------------------------|--|---|--|-----------------|---|--|---------------------------|-----|
| 0                             | [AdarshPallav]           | [Adarsh, Pallav]                         | [RT, satishacharya, Blueprint, 25 years, Budg..]  | [ ]  | 7               | [Twitter, Android]                        | 2022-02-07 07:11:43+00:00                  | 38                        |     |
| 1                             | [bsindia]                | [Business, Standard]                     | [Latest, news, the, economy, companies, market..] | [BSM, MorningShow, RBI, go, a hike, reverse, resp..] | [India]         | 2187279                                   | [TweetDeck]                                | 2022-02-07 07:11:50+00:00 | 430 |
| 2                             | [bsindia]                | [Business, Standard]                     | [Latest, news, the, economy, companies, market..] | [BSM, MorningShow, RBI, go, a hike, reverse, resp..] | [India]         | 2187279                                   | [Twitter, Web, Browser]                    | 2022-02-07 07:11:50+00:00 | 430 |

|   |             |                      |  |   |         |         |             |                              |     |
|---|-------------|----------------------|--|---|---------|---------|-------------|------------------------------|-----|
|   | [bsindia]   | [Business, Standard] | [Latest news, the economy, companies, market...] | [BS]MorningShow, Business, Standards, anuprotive... | [India] | 2187279 | [TweetDeck] | 2022-02-07<br>07:10:00+00:00 | 430 |
| 4 | [orfonline] | [ORF]                | [Nonpartisan, independent analyses, security...] | [Tax, concessions, cooperative societies, a ...]    | [India] | 106444  | [TweetDeck] | 2022-02-07<br>07:10:00+00:00 | 136 |

### PARTS OF SPEECH REMOVAL

**Parts of Speech (POS) Tags Abbreviations Used**

| Tag | Description                              |
|-----|--|
| CC  | Coordinating conjunction                 |
| CD  | Cardinal number                          |
| DT  | Determiner                               |
| EX  | Existential there                        |
| FW  | Foreign word                             |
| IN  | Preposition or subordinating conjunction |
| JJ  | Adjective                                |
| JJR | Adjective, comparative                   |
| JJS | Adjective, superlative                   |

|      |                                       |
|------|---------------------------------------|
| MD   | Modal                                 |
| NN   | Noun, singular or mass                |
| NNS  | Noun, plural                          |
| NNP  | Proper noun, singular                 |
| NNPS | Proper noun, plural                   |
| PD   | Predeterminer                         |
| POS  | Possessive ending                     |
| PP   | Personal pronoun                      |
| PRP  | Possessive pronoun                    |
| RB   | Adverb                                |
| RBR  | Adverb, comparative                   |
| RBS  | Adverb, superlative                   |
| RP   | Particle                              |
| SYM  | Symbol                                |
| TO   | to                                    |
| UH   | Interjection                          |
| VB   | Verb, base form                       |
| VBD  | Verb, past tense                      |
| VBG  | Verb, gerund or present participle    |
| VBN  | Verb, past participle                 |
| VBP  | Verb, non-3rd person singular present |
| VBT  | Verb, 3rd person singular present     |

|     |                       |
|-----|-----------------------|
| WDT | Wh-determiner         |
| WP  | Wh-pronoun            |
| WP  | Possessive wh-pronoun |
| WRB | Wh-adverb             |

  

```

: # list of nltk pos tags
:
nltk.help.upenn_tagset()

$: dollar
$ - $ --$ $S $S $G $K $M $N $Z $S $U $S $U $S
"": closing quotation mark
" ": opening parenthesis
( ): closing parenthesis
[ ]: closing parenthesis
, : comma
- : dash
--: dash
.: sentence terminator
. ! ?
: colon or ellipsis
: ? . ,
CC: conjunction, coordinating
& ^ n and both but either et for less minus neither nor or plus so
therefore times v versus vs whether yet

```

mid-1890 nine-thirty fly-two ten-one-tenth ten million 0.5 one forty-seven 1987 twenty 799 zero two 78-degrees eighty-four IX '60s .025 fifteen 271,124 dozen quintillion DMG,000 ...

DT: determiner  
all an another any both del each either every half la many much nary neither no some such that the them these this those

EX: existential there  
there

FW: foreign word  
gemeinschaft hund ich jeux habebas Haasentiera Her K'ang-i vov  
luthaiw alai je joue objets salutaris fille quibusdam pas trop Monte terran fiche ou corporis ...

IN: preposition or conjunction, subordinating  
astride among uponn whether out inside pro despite on by throughout below within for towards near behind atop around if like until below near into it beside ...

JJ: adjective or numeral, ordinal  
third ill-mannered pre-war regrettable oiled calamitous first separable autoplasmic battery-powered participatory fourth still-to-be-named multilingual multi-disparatory ...

JRn: adjective, comparative  
bleaker braver breezier brierer brighter brisker broader bumper busier calmer cheaper chesher clearer clearer colder colder commoner costlier cozier crierier crunchier cuter ...

JZS: adjective, superlative  
calmest cheapest choicest classiest cleanest cleanest closest commonest corniest costliest crassest creepiest crudest cutest darkest deadliest dearest deepest densest dinkiest ...

LS: list item marker

SP44007 Second Third Three Two 'a b c d first five four one six three two

MD: modal auxiliary  
can cannot could couldn't dare may might must need ought shall should  
shouldn't will would

NN: noun, common, singular or mass  
common-carrier cabbage knuckle-duster Casino aghya shed thermostat  
investment slide humour falloff slick wind hyphen override subhumanity  
machinist ...

NNP: noun, proper, singular  
Norsk Vennedagsgæstecache Ranner Conchita Trumple Ervin Oddi Charly CTOA  
Shannon A.K.C. Meltex Liverpool ...

NNPS: noun, proper, plural  
Americans Americans Amharas Amityvilles Amusements Anarcho-Syndicalists  
Andalusians Andes Andrusas Angeles Animals Anthony Antilles Antiques  
Apache Apaches Apocrypha ...

NNS: noun, common, plural  
undergraduates notches bric-a-brac products bodyguards facets coasts  
disinfectants storehouses designs club fragrances averages  
subjectivists apprehensions messes factory-jobs ...

PDT: pre-determiner  
all both half many quite such sure this

POS: genitive marker  
's

PPR: pronoun, personal  
hers herself him himself himself it itself me myself one oneself ours  
ourselves oneself self she thee their them themselves thy thou thy us

PPRS: pronoun, possessive  
her his mine my our ours their thy your

occasionally unabatingly maddeningly adventurously professedly  
stirring prominently technologically masterfully predominantly  
swiftly fiscally pitilessly ...

MNR: adverb, comparative

very gloomier grander graver greater grimmer harder harsher  
healthier heavier higher however larger later lazier leanlier less-  
perfectly lesser lonelier longer louder lower more ...

RSD: adverb, superlative

best biggest bluntest earliest fastest first furthest hardest  
heartiest highest largest least least most nearest second tightest worst

RP: particle

aboard about across along apart around aside at away back before behind  
by down down ever fast for forth from go high i.e. in into just later  
low more off on open out over per pie raising start teeth that through  
under unto up uppp upon whole with ...

SYM: symbol

\$ % ^ & \* ~ . , | . + , . < - = > B A f% U S U S S R \$ \* \* \* \*

TO: "to" as preposition or infinitive marker -

UH: interjection

Goodbye Goody Gosh Wow Jeepers Gee-oh Sus Hubba Hey Keerleest Cops amen  
hun huddy uh dammit whammo shucks heck anyways whodunnit honey golly  
man baby diddle hush sunavabithe ...

VB: verb, base form

ask assemble assess assign assume attend attention avoid bake balkanize  
bank begin behold believe bend benefit beneil beware bleis blow bomb  
boost break bring broil brush build ...

VBD: verb, past tense

dipped pleaded swiped regummed soaked stidied convened halted registered

```

WPG: speech, present participle ...
WPG: verb, present participle or gerund
    telegraphing stirring focusing angering judging stalling lactating
    hankerin' alleging weeping capping approaching traveling besieging
    encrypting interrupting exclaiming winning ...
VBN: verb, past participle
    multibullded dilapidated aerosolized chained languished panelized used
    expedientized flourished initiated reunified condensed condensed
    unsettled primed dubbed desired ...
VBP: verb, present tense, 3rd person singular
    predominate strap resort sue twist spill cure lengthen brush terminate
    appear tread away gladden obtain comprise detect tease attract
    emphasize mold postpone serve return war ...
VBS: verb, present tense, 3rd person singular
    bases reconstructs marks mixes displeases seals carps weaves snatches
    slumps stretches authorizes amazes amazes pictures amazes stockpiles
    seduces fizzes uses bolsters slanders speaks pleads ...
WDT: WH-determiner
    that that whatever which whichever
WP: WH-pronoun
    that that whatever whatsoever which who whom whosoever
WPP: WH-pronoun, possessive
    whose
WRB: WH-adverb
    how however whenever wherever whereby wherever wherein whereof why
    '': opening quotation mark

# We will perform parts of speech (POS) tagging for every attribute

```

```
pos_tagged_df = pd.DataFrame({
    'User_Name[POS Tagged]', 'Name[POS Tagged]',
    'Description (Status)[POS Tagged]', 'Tweets Posted[POS Tagged]',
    'Location[POS Tagged]', 'No of Followers',
    'Device Used for Tweet[POS Tagged]', 'Date of Tweet Creation[POS Tagged]',
    'No of Friends'
})

# creating a dataframe to store parts of speech (POS) tagged data of each attribute

pos_tag_df = pd.DataFrame()

# performing the parts of speech (POS) tagging using the nltk.pos_tag function

for i in range(len(pos_tag_column_name)):
    if (stop_word_removed_df[stop_word_removed_column_name[i]].dtype ==
        np.object):
        pos_tag_df[pos_tag_column_name[i]] = stop_word_removed_df.apply(
            lambda row: nltk.pos_tag(row[stop_word_removed_column_name[i]]),
            axis=1)
    else:
        pos_tag_df[pos_tag_column_name[i]] = stop_word_removed_df[
            stop_word_removed_column_name[i]]

# after performing tokenization starting few tokenized data

pos_tag_df.head()
```

| User_Name[POS Tagged] | Name[POS Tagged]      | Description (Status)[POS Tagged]  | Tweets Posted[POS Tagged]                                | Location[POS Tagged] | No of Followers | Device Used for Tweet[POS Tagged]      | Date of Tweet Creation[POS Tagged] | No of Friends |
|-----------------------|-----------------------|---|--|----------------------|-----------------|--|------------------------------------|---------------|
| 0                     | [[Adarsh,Pallav, NN]] | [[[Adarsh, NN], [Pallav, NN]]<br>[[satisfacharya, VBD], [Blueprint, ...]]                       | [[[RT, NNP], [satisfacharya, VBD], [Blueprint, ...]]]    | [[[India, NN]]]      | 7               | [[Twitter, NNP], [Android, NN]]        | 2022-02-07<br>07:11:43+00:00       | 38            |
| 1                     | [[[bsindia, NN]]]     | [[[Business, NNP], [Standard, NNP]]<br>[[[Latest, NNP], [news, NN], [the, DT], [econom., ...]]] | [[[BSMorningShow, NNP], [RBI, NNP], [go, VBP], ...]]     | [[[India, NN]]]      | 2178279         | [[TweetDeck, NN]]                      | 2022-02-07<br>07:11:00+00:00       | 430           |
| 2                     | [[[bsindia, NN]]]     | [[[Business, NNP], [Standard, NNP]]<br>[[[Latest, NNP], [news, NN], [the, DT], [econom., ...]]] | [[[BSMorningShow, NNP], [RBI, NNP], [go, VBP], ...]]     | [[[India, NN]]]      | 2178279         | [[Twitter, NNP], [Web, NN], [App, NN]] | 2022-02-07<br>07:10:43+00:00       | 430           |
| 3                     | [[[bsindia, NN]]]     | [[[Business, NNP], [Standard, NNP]]<br>[[[Latest, NNP], [news, NN], [the, DT], [econom., ...]]] | [[[BSMorningShow, NNP], [Business, NNP], [Stand., ...]]] | [[[India, NN]]]      | 2178279         | [[TweetDeck, NN]]                      | 2022-02-07<br>07:10:50+00:00       | 430           |
| 4                     | [[[orfonline, NN]]]   | [[[ORF, NN]]<br>[[[Nonpartisan, NNP], [independent, JJ], [analy., ...]]]                        | [[[Tax, NNP], [concessions, NNS], [cooperative., ...]]]  | [[[India, NN]]]      | 106444          | [[TweetDeck, NN]]                      | 2022-02-07<br>07:10:00+00:00       | 136           |

```
# We will perform stemming & lemmatization for every attribute

stemmed_column_name = [
    'User Name[Stemmed]', 'Name[Stemmed]', 'Description (Status)[Stemmed]',
    'Tweets Posted[Stemmed]', 'Location[Stemmed]', 'No of Followers',
    'Device Used for Tweet[Stemmed]', 'Date of Tweet Creation[Stemmed]',
    'No of Friends'
]

lemmatized_column_name = [
    'User Name[Lemmatized]', 'Name[Lemmatized]',
    'Description (Status)[Lemmatized]', 'Tweets Posted[Lemmatized]',
    'Location[Lemmatized]', 'No of Followers',
    'Device Used for Tweet[Lemmatized]', 'Date of Tweet Creation[Lemmatized]',
    'No of Friends'
]

# creating a dataframe to store stemmed & lemmatized data of each attribute

stem_lemma_df = pd.DataFrame()

# Function for stemming & lemmatizing the word

def stem_word(lst):
    for i in range(len(lst)):
        lst[i] = nltk.stem.PorterStemmer().stem(lst[i])
```

```
def lemma_word(lst):
    for i in range(len(lst)):
        lst[i] = nltk.stem.WordNetLemmatizer().lemmatize(lst[i])
    return lst
```



