

D1A1tweets.csv

In [53]:

```
import numpy as np
import pandas as pd
import re #text cleaning (preprocessing)
import nltk #natural language toolkit, used for preprocessing
import string
from nltk.stem.porter import PorterStemmer
from nltk.tokenize import word_tokenize
from nltk.corpus import stopwords
```

Out [54]:

|   | TWEET_ID           | TWEET   | TWEET_Label |
|---|--------------------|---|-------------|
| 0 | 8030860482123700   | seriously? racist mcdonalds™s sign is obvious...  | false       |
| 1 | 80384559739030000  | hoax: mcdonald's issues official statement on...  | false       |
| 2 | 91728807861426900  | spsa please do not drink any pepsi soda, a wor... | false       |
| 3 | 26595285247209000  | deep-fried left wings demo-crab cakes barack-a... | false       |
| 4 | 273182568298450000 | 42 million dead in bloodiest black friday week... | false       |

In [55]:

```
length_tweet = df['TWEET'].str.len().plot.hist(color = 'blue', figsize = (6,4))
```

In [56]:

```
dfWhole['length'] = df['TWEET'].str.len()
dfWhole['length'].describe()
```

Out [56]:

```
count    2388    0.000000
mean      92.633016
std       23.948301
min        9.000000
50%      77.000000
90%      95.000000
95%      110.000000
max      257.000000
Name: length, dtype: float64
```

In [57]:

```
dfWhole[dfWhole['length']==100]['TWEET'].iloc[0]
```

Out [57]:

```
"highly radioactive water is seeping into the ocean from japan's crippled fukushima nuclear plant URL"
```

In [58]:

```
print(dfWhole)
```

Out [58]:

|                         | TWEET_ID            | TWEET   | TWEET_Label |
|-------------------------|---------------------|---|-------------|
| 0                       | 80808680482123700   | seriously? racist mcdonalds™s sign is obvious...  | 0           |
| 1                       | 80808680482123700   | hoax: mcdonald's issues official statement on...  | 1           |
| 2                       | 91728807861426900   | spsa please do not drink any pepsi soda, a wor... | 2           |
| 3                       | 26595285247209000   | deep-fried left wings demo-crab cakes barack-a... | 3           |
| 4                       | 273182568298450000  | 42 million dead in bloodiest black friday week... | 4           |
| ...                     | ...                 | ...   | ...         |
| 2387                    | 7786882510645370000 | black men may have cause to run from police, a... | ...         |
| 2394                    | 7789497491562459800 | a black man who runs from police merely might...  | ...         |
| 2395                    | 7786384848896628000 | this network of tunnels is from the stone age...  | ...         |
| 2396                    | 7786384848896628000 | well, this is a new one: chelsea clinton impli... | ...         |
| 2397                    | 788882510645370000  | chelsea clinton implies here that marijuana ca... | ...         |
| ...                     | ...                 | ...   | ...         |
| 2393                    | unverified          | 110   | ...         |
| 2394                    | unverified          | 110   | ...         |
| 2395                    | unverified          | 117   | ...         |
| 2396                    | unverified          | 113   | ...         |
| 2397                    | unverified          | 65  | ...         |
| [2388 rows x 4 columns] |                     |   |             |

In [59]:

```
#.....changing the casing of the words.....
df = dfWhole.iloc[:,1:2]
df['TWEET'] = df['TWEET'].astype(str)
df['TWEET_lower'] = df['TWEET'].str.lower()
df.head()
```

Out [59]:

|   | TWEET   | TWEET_lower                                       |
|---|---|---|
| 0 | seriously? racist mcdonalds™s sign is obvious...  | seriously? racist mcdonalds™s sign is obvious...  |
| 1 | hoax: mcdonald's issues official statement on...  | hoax: mcdonald's issues official statement on...  |
| 2 | spsa please do not drink any pepsi soda, a wor... | spsa please do not drink any pepsi soda, a wor... |
| 3 | deep-fried left wings demo-crab cakes barack-a... | deep-fried left wings demo-crab cakes barack-a... |
| 4 | 42 million dead in bloodiest black friday week... | 42 million dead in bloodiest black friday week... |

In [60]:

```
#Remove punctuations
df.drop(['TWEET_lower'], axis=1, inplace=True)
puncremove = string.punctuation
def remove_punctuation(TWEET):
    return TWEET.translate(str.maketrans('', '', puncremove))
df['TWEET_punct'] = df['TWEET'].apply(lambda TWEET: remove_punctuation(TWEET))
df.head()
```

Out [60]:

|   | TWEET   | TWEET_punct  |
|---|---|--|
| 0 | seriously? racist mcdonalds™s sign is obvious...  | seriously? racist mcdonalds™s sign is obvious...   |
| 1 | hoax: mcdonald's issues official statement on...  | hoax mcdonalds issues official statement on ra...  |
| 2 | spsa please do not drink any pepsi soda, a wor... | psa please do not drink any pepsi soda, a worke... |
| 3 | deep-fried left wings demo-crab cakes barack-a... | deepfried left wings democrab cakes barackamol...  |
| 4 | 42 million dead in bloodiest black friday week... | 42 million dead in bloodiest black friday week...  |

In [61]:

```
import numpy as np
import pandas as pd
import re #text cleaning (preprocessing)
import nltk #natural language toolkit, used for preprocessing
import string
from nltk.stem.porter import PorterStemmer
from nltk.tokenize import word_tokenize
from nltk.corpus import stopwords
```

In [62]:

```
stops = set(stopwords.words('english'))
print(stops)
```

Out [62]:

```
{'those', 'mustn', 'ain', 'no', 'or', 'she's', 'me', 'down', 'most', 'had', 'you're', 'it', 'which', 'that'll', 'from', 'having', 'needn', 'wouldn', 'why', 'were', 'isn', 'more', 'w
on't', 'up', 'am', 'just', 'ourselves', 'the', 'mustn't', 'by', 'yourself', 'whom', 'did', 'and', 'what', 'after', 'haven', 're', 't', 'a', 'than', 'them', 'you'll', 'o', 'hasn', 'ha
ven't', 'is', 'was', 'between', 'doesn', 'about', 'once', 'm', 'does', 'for', 'ours', 'hers', 'into', 'couldn't', 'don', 'shan't', 'so', 'didn', 'been', 'out', 'myself', 'again',
'be'one', 'too', 'n't', 'shouldn', 'shouldn't', 'won't', 'theirs', 'hisself', 'his', 'herself', 'he', 'above', 'few', 'aren', 'doesn't', 'your', 'you've', 'do', 'all', 'with',
'hadn't', 'when', 'it's', 'you'd', 'there', 'being', 'mightn', 'doing', 'as', 'other', 'wasn't', 'hadn', 'at', 'each', 'some', 'not', 'any', 'any', 'shan', 'i', 'you', 'isn't', 'its
elf', 'these', 'here', 'y', 'then', 'yours', 'until', 'off', 'how', 'are', 'while', 'mightn't', 'very', 'needn't', 'don't', 'lover', 'or', 'under', 'themselves', 'd', 'its', 'becau
se', 't', 'remove_punctuations', 'her', 'own', 'further', 'didn't', 'to', 'now', 'ha', 'an', 'will', 'against', 'where', 'couldn', 'during', 'has', 'nor', 'both', 'sh', 've', 'only',
'he', 'ain', 'same', 'this', 'yourselves', 'on', 'such', 'should', 'wasn', 'through', 'they', 'aren't', 'but', 'that', 's', 'hasn't', 'wouldn't', 'we', 'have', 'has', 'weren't', 'i
n', 'weren', 'our', 'can', 'who'}
```

In [63]:

```
STOPWORDS = set(stopwords.words('english'))
def remove_stopwords(TWEET):
    """custom function to remove the stopwords"""
    return " ".join([word for word in str(TWEET).split() if word not in STOPWORDS])
df['TWEET_stop'] = df['TWEET_punct'].apply(lambda TWEET: remove_stopwords(TWEET))
df.head()
```

Out [63]:

|   | TWEET   | TWEET_punct  | TWEET_stop  |
|---|---|--|---|
| 0 | seriously? racist mcdonalds™s sign is obvious...  | seriously racist mcdonalds™s sign is obvious...    | seriously racist mcdonalds™s sign obviously hoax  |
| 1 | hoax: mcdonald's issues official statement on...  | hoax mcdonalds issues official statement on ra...  | hoax mcdonalds issues official statement racis... |
| 2 | spsa please do not drink any pepsi soda, a wor... | psa please do not drink any pepsi soda, a worke... | psa please drink pepsi soda worker company put... |
| 3 | deep-fried left wings demo-crab cakes barack-a... | deepfried left wings democrab cakes barackamol...  | deepfried left wings democrab cakes barackamol... |
| 4 | 42 million dead in bloodiest black friday week... | 42 million dead in bloodiest black friday week...  | 42 million dead bloodiest black friday weekend... |

In [64]:

```
from collections import Counter
cnt = Counter()
for TWEET in df['TWEET_stop'].values:
    for word in TWEET.split():
        cnt[word] += 1
cnt.most_common(10)
```

Out [64]:

```
[('URL', 2990),
 ('police', 139),
 ('shot', 111),
 ('says', 102),
 ('ferguson', 92),
 ('breaking', 89),
 ('new', 88),
 ('us', 84),
 ('obama', 82),
 ('paul', 82)]
```

In [65]:

```
freqwords = set([w for (w, wc) in cnt.most_common(10)])
def remove_freqwords(TWEET):
    """custom function to remove the frequent words"""
    return " ".join([word for word in str(TWEET).split() if word not in freqwords])
df['TWEET_stopfreq'] = df['TWEET_stop'].apply(lambda TWEET: remove_freqwords(TWEET))
df.head()
```

Out [65]:

|   | TWEET   | TWEET_punct  | TWEET_stop  | TWEET_stopfreq                                    |
|---|---|--|---|---|
| 0 | seriously? racist mcdonalds™s sign is obvious...  | seriously racist mcdonalds™s sign is obvious...    | seriously racist mcdonalds™s sign obviously hoax  | seriously racist mcdonalds™s sign obviously hoax  |
| 1 | hoax: mcdonald's issues official statement on...  | hoax mcdonalds issues official statement on ra...  | hoax mcdonalds issues official statement racis... | hoax mcdonalds issues official statement racis... |
| 2 | spsa please do not drink any pepsi soda, a wor... | psa please do not drink any pepsi soda, a worke... | psa please drink pepsi soda worker company put... | psa please drink pepsi soda worker company put... |
| 3 | deep-fried left wings demo-crab cakes barack-a... | deepfried left wings democrab cakes barackamol...  | deepfried left wings democrab cakes barackamol... | deepfried left wings democrab cakes barackamol... |
| 4 | 42 million dead in bloodiest black friday week... | 42 million dead in bloodiest black friday week...  | 42 million dead bloodiest black friday weekend... | 42 million dead bloodiest black friday weekend... |

In [66]:

```
df.drop(['TWEET_punct', 'TWEET_stop'], axis=1, inplace=True)
```

In [67]:

```
#stemming (playing play)
stemmer = PorterStemmer()
def stem_words(TWEET):
    return " ".join([stemmer.stem(word) for word in TWEET.split()])
df['TWEET_stemmed'] = df['TWEET'].apply(lambda TWEET: stem_words(TWEET))
df.head(10)
```

Out [67]:

|   | TWEET   | TWEET_stopfreq                                    | TWEET_stemmed                                      |
|---|---|---|--|
| 0 | seriously? racist mcdonalds™s sign is obvious...  | seriously racist mcdonalds™s sign obviously hoax  | seriously? racist mcdonalds™s sign obviously h...  |
| 1 | hoax: mcdonald's issues official statement on...  | hoax mcdonalds issues official statement racis... | hoax: mcdonald' issu official statement on raci... |
| 2 | spsa please do not drink any pepsi soda, a wor... | psa please drink pepsi soda worker company put... | psa pleas do not drink any pepsi soda, a wor...    |
| 3 | deep-fried left wings demo-crab cakes barack-a... | deepfried left wings democrab cakes barackamol... | deep-fri left wing demo-crab cake barack-amol...   |
| 4 | 42 million dead in bloodiest black friday week... | 42 million dead bloodiest black friday weekend... | 42 million dead in bloodiest black friday week...  |
| 5 | 42 million dead in bloodiest black friday week... | 42 million dead bloodiest black friday weekend... | 42 million dead in bloodiest black friday week...  |
| 6 | #prayforchristopher 5k run - well we are walk...  | prayforchristopher 5k run well walking lol sav... | #prayforchristoph 5k run - well we are walk l...   |
| 7 | a photo of black nurses saving the life of a k... | photo black nurses saving life kkk member         | a photo of black nurs save the life of a kkk m...  |
| 8 | a photo of black nurses saving the life of a k... | photo black nurses saving life kkk member         | a photo of black nurs save the life of a kkk m...  |
| 9 | a photo of black nurses saving the life of a k... | photo black nurses saving life kkk member         | a photo of black nurs save the life of a kkk m...  |

In [68]:

```
from nltk.stem.snowball import SnowballStemmer
SnowballStemmer.languages
```

Out [68]:

```
('arabic',
 'danish',
 'dutch',
 'english',
 'finnish',
 'french',
 'german',
 'hungarian',
 'italian',
 'norwegian',
 'porter',
 'portuguese',
 'romanian',
 'russian',
 'spanish',
 'swedish')
```

In [69]:

```
#lemmizer
lemmatizer = WordNetLemmatizer()
def lemmatize_words(TWEET):
    return " ".join([lemmatizer.lemmatize(word) for word in TWEET.split()])
df['TWEET_lemmatized'] = df['TWEET'].apply(lambda TWEET: lemmatize_words(TWEET))
df.head()
```

Out [69]:

|   | TWEET   | TWEET_stopfreq                                    | TWEET_stemmed                                      | TWEET_lemmatized                                   |
|---|---|---|--|--|
| 0 | seriously? racist mcdonalds™s sign is obvious...  | seriously racist mcdonalds™s sign obviously hoax  | seriously? racist mcdonalds™s sign obviously h...  | seriously? racist mcdonalds™s sign obviously hoax  |
| 1 | hoax: mcdonald's issues official statement on...  | hoax mcdonalds issues official statement racis... | hoax: mcdonald' issu official statement on raci... | hoax: mcdonald's issue official statement racis... |
| 2 | spsa please do not drink any pepsi soda, a wor... | psa please drink pepsi soda worker company put... | psa pleas do not drink any pepsi soda, a wor...    | psa please do not drink any pepsi soda, a wor...   |
| 3 | deep-fried left wings demo-crab cakes barack-a... | deepfried left wings democrab cakes barackamol... | deep-fri left wing demo-crab cake barack-amol...   | deep-fried left wings democrab cake barack-amol... |
| 4 | 42 million dead in bloodiest black friday week... | 42 million dead bloodiest black friday weekend... | 42 million dead in bloodiest black friday week...  | 42 million dead in bloodiest black friday week...  |

In [70]:

```
lemmatizer.lemmatize("sleeping")
```

Out [70]:

```
'sleeping'
```

In [71]:

```
lemmatizer.lemmatize("sleeping","v") # v1
```

Out [71]:

```
'sleep'
```

In [72]:

```
print("The word is : stripes")
print("Lemma result for verb : ",lemmatizer.lemmatize("stripes", 'v'))
print("Lemma result for noun : ",lemmatizer.lemmatize("stripes", 'n'))

The word is : stripes
Lemma result for verb : strip
Lemma result for noun : stripe
```

In [73]:

```
def remove_urls(TWEET):
    url_pattern = re.compile(r'(https?://S+|www\..S+)'
    return url_pattern.sub(r'', TWEET)
```

In [74]:

```
s = 'a\b\nc\td'
print(s)

a      b
      c      d
```

In [75]:

```
s = r'a\b\nc\td'
print(s)

a\b\nc\td
```

In [76]:

```
#removal url
TWEET1 = "This is my website, https://www.abc.com, check it out"
remove_urls(TWEET1)
```

Out [76]:

```
'This is my website, check it out'
```

In [77]:

```
TWEET = "Want to learn more. Checkout www.h2o.ai for additional information"
remove_urls(TWEET)
```

Out [77]:

```
'Want to learn more. Checkout for additional information'
```

In [78]:

```
import numpy as np
import pandas as pd
import re # used for preprocessing
import nltk # Natural Language Toolkit, used for preprocessing
import string used for preprocessing
from nltk.stem.porter import PorterStemmer
from nltk.stem import WordNetLemmatizer
from nltk.tokenize import word_tokenize
from nltk.corpus import stopwords

stops = set(stopwords.words('english'))
```

In [79]:

```
dfWhole = pd.read_csv("D:/A1/tweets.csv", nrows=2388)
df = dfWhole[['TWEET']]
dfWhole.head()
```

Out [79]:

|   | TWEET_ID           | TWEET   | TWEET_Label |
|---|--------------------|---|-------------|
| 0 | 8008080482123700   | seriously? racist mcdonalds™s sign is obvious...  | false       |
| 1 | 80084557338030000  | hoax: mcdonald's issues official statement on...  | false       |
| 2 | 917288078614269000 | spsa please do not drink any pepsi soda, a wor... | false       |
| 3 | 265952852472090000 | deep-fried left wings demo-crab cakes barack-a... | false       |
| 4 | 273182568298450000 | 42 million dead in bloodiest black friday week... | false       |

In [80]:

```
#df = dfWhole.iloc[:,1:2]
#df['TWEET'] = df['TWEET'].astype(str)
#df.head()
```

In [81]:

```
# remove all urls
def remove_urls(TWEET):
    url_pattern = re.compile(r'(https?://S+|www\..S+)'
    return url_pattern.sub(r'', TWEET)
# make all text lowercase
def TWEET_lowerCase(TWEET):
    return TWEET.lower()
# remove numbers
def remove_numbers(TWEET):
    result = re.sub(r'(\d+)', '', TWEET)
    return result
# remove punctuation
def remove_punctuation(TWEET):
    translator = str.maketrans('', '', string.punctuation)
    return TWEET.translate(translator)
# tokenize
def tokenize(TWEET):
    TWEET = word_tokenize(TWEET)
    return TWEET
# remove stopwords
stop_words = set(stopwords.words('english'))
def remove_stopwords(TWEET):
    TWEET = [i for i in TWEET if not i in stop_words]
    return TWEET
# lemmatize
lemmatizer = WordNetLemmatizer()
def lemmatize(TWEET):
    TWEET = [lemmatizer.lemmatize(token) for token in TWEET]
    return TWEET

def preprocessing(TWEET):
    TWEET = TWEET.lowerCase(TWEET)
    TWEET = remove_urls(TWEET)
    TWEET = remove_numbers(TWEET)
    TWEET = remove_punctuation(TWEET)
    TWEET = tokenize(TWEET)
    TWEET = remove_stopwords(TWEET)
    TWEET = lemmatize(TWEET)
    TWEET = ' '.join(TWEET)
    return TWEET
```

In [82]:

```
df.head()
```

Out [82]:

|   | TWEET   | pp_TWEET  |
|---|---|---|
| 0 | seriously? racist mcdonalds™s sign is obvious...  | seriously racist mcdonalds™s sign obviously h...  |
| 1 | hoax: mcdonald's issues official statement on...  | hoax mcdonalds issue official statement racis...  |
| 2 | spsa please do not drink any pepsi soda, a wor... | psa please drink pepsi soda worker company put... |
| 3 | deep-fried left wings demo-crab cakes barack-a... | deepfried left wing democrab cake barackamol...   |
| 4 | 42 million dead in bloodiest black friday week... | million dead bloodiest black friday weekend re... |

In [83]:

```
final_TWEET_data = list(df['pp_TWEET'])
```

In [86]:

```
from sklearn.feature_extraction.text import TfidfVectorizer

tf=TfidfVectorizer()

# the vectorizer must be fit onto the entire corpus
fitted_vectorizer = tf.fit(final_TWEET_data)

transform_all = fitted_vectorizer.transform(df['pp_TWEET'])
```

In [87]:

```
print(transform_all)
```

Out [87]:

```
((0, 5676) 0.06521763769677442
(0, 4386) 0.3848889279674288
(0, 4271) 0.39177991493899984
(0, 3849) 0.3848899279674288
(0, 3361) 0.46397160887798154
(0, 3832) 0.44666634547925484
(0, 2263) 0.3712263400816987
(1, 5679) 0.054772357129968916
(1, 5678) 0.35611582655709963
(1, 4586) 0.328632393181262
(1, 4358) 0.32256983623391985
(1, 4272) 0.38966188228819826
(1, 3849) 0.32256983623391985
(1, 3375) 0.24231715876571722
(1, 3851) 0.35611582655709963
(1, 2529) 0.34531642786961864
(1, 2283) 0.31177643154652886
(2, 5296) 0.24231653638144783
(2, 4453) 0.2873154504280958
(2, 4162) 0.1938863175354777
(2, 3854) 0.23384238746402559
(2, 3782) 0.2873154504280958
(2, 3823) 0.4529628457632189
(2, 3450) 0.24811143169880563
(2, 2453) 0.22337641720438547

(2305, 2344) 0.36781422126680346
(2305, 1137) 0.36781422126680346
(2305, 1084) 0.36781422126680346
(2305, 108) 0.27891954658947497
(2306, 5372) 0.320977690923629
(2306, 5228) 0.2691627288007266
(2306, 5079) 0.048864263597869135
(2306, 4876) 0.320777690923629
(2306, 4761) 0.23935770865440845
(2306, 3402) 0.19271967759691673
(2306, 3269) 0.1844776570798465
(2306, 3131) 0.3287277696933629
(2307, 2985) 0.2815329290821064
(2306, 2677) 0.24399915913390298
(2306, 2412) 0.29896774854764475
(2306, 947) 0.281168877644837
(2306, 870) 0.3688074947120958
(2306, 153) 0.320777690923629
(2307, 6079) 0.0708263152545684
(2307, 5088) 0.47712810795962923
(2307, 2985) 0.40857177590373817
(2307, 2677) 0.3527951604970869
(2307, 2412) 0.43387387826984286
(2307, 947) 0.2901163111426487
(2307, 870) 0.4452512288760943
```