

In [1]:

```
#Loading NLTK
#TEXT MINING ANALYSIS
#1.NLTK IS A POWERFUL PACKAGE THAT PROVIDES A SET OF DIVERSE NATURAL LANGUAGES ALGORITHM.
#2.IT IS FREE,OPENSOURCE EASY TO USE AND WEL DOCUMENTED.
#3.NLTK CONSISTS OF THE MOST COMMON ALGORITHMS SUCH AS TOKENZING,PART OD SPEECH TAGGING,STEMMING,SENTIMENT ANALYSIS,
# TOPIC SEGMENTATION,AND NAMED ENTITY RECOGNITION NLTK HELPS THE COMPUTER TO ANALYSIS,PREPROCESS,AND UNDERSTAND THE WRITTEN TEXT.
import nltk
```

In [10]:

```
#Tokenization is the first step in Text Analytics.
#The Process of Breaking Down a Text Paragraph into Smaller Chunks Such as Words or Sentence is Called Tokenization.
#Token is Single Entity That is Building Blocks For Sentence or Paragraph.
#SENTENCE TOKENIZATION
from nltk.tokenize import sent_tokenize
text="""Hello Miss.Vanita,what are you doing today? the weather is great,and city is awesome. The Sky is Pinkish-Blue."""
tokenized_sent=sent_tokenize(text)
print(tokenized_sent)

['Hello Miss.Vanita,what are you doing today?', 'the weather is great,and city is awesome.', 'The Sky is Pinkish-Blue.']
```

In [11]:

```
# Word Tokenizer Breaks Text Paragraph into Words.
# WORD TOKENIZATION
from nltk.tokenize import word_tokenize
text="""Hello Miss.Vanita,what are you doing today? the weather is great,and city is awesome. The Sky is Pinkish-Blue."""
tokenized_word=word_tokenize(text)
print(tokenized_word)

['Hello', 'Miss.Vanita', ',', 'what', 'are', 'you', 'doing', 'today', '?', 'the', 'weather', 'is', 'great', ',', 'and', 'city', 'is', 'awesome', '.', 'The', 'Sky', 'is', 'Pinkish-Blue', '.']
```

In [8]:

```
#FREQUENCY DISTRIBUTION
from nltk.probability import FreqDist
fdist=FreqDist(tokenized_word)
print(fdist)

<FreqDist with 20 samples and 24 outcomes>
```

In [6]:

```
fdist.most_common(2)
```

Out[6]:

```
[('is', 3), ('', 2)]
```

In [9]:

```
#FREQUENCY DISTRIBUTION PLOT
import matplotlib.pyplot as plt
fdist.plot(30,cumulative=False)
plt.show()
```

In [10]:

```
import nltk
nltk.download('punkt')
nltk.download('wordnet')
```

Out[10]:

```
[nltk_data] Downloading package punkt to
[nltk_data] C:\Users\Owner\AppData\Roaming\nltk_data...
[nltk_data] Package punkt is already up-to-date!
[nltk_data] Downloading package wordnet to
[nltk_data] C:\Users\Owner\AppData\Roaming\nltk_data...
True
```

In [22]:

```
nltk.word_tokenize("hi How are you")
```

Out[22]:

```
['hi', 'How', 'are', 'you']
```

In [12]:

```
#STOPWORDS CONSIDERED AS NOISE IN THE TEXT.TEXT MAY CONTAIN STOP WORDS SUCH AS IS,AM,ARE,THIS,A,AN,THE,etc
#STOPWORDS
from nltk.corpus import stopwords
stop_words=stopwords.words("english")
print(stop_words)
```

Out[12]:

```
['i', 'me', 'my', 'myself', 'we', 'our', 'ours', 'ourselves', 'you', 'you're', 'you've', 'you'll', 'you'd', 'your', 'yours', 'yourself', 'yourselves', 'he', 'him', 'his', 'himself', 'she', 'she's', 'her', 'hers', 'herself', 'it', 'it's', 'its', 'itself', 'they', 'them', 'their', 'theirs', 'themselves', 'what', 'wh ich', 'who', 'whom', 'this', 'that', 'that'll', 'these', 'those', 'am', 'is', 'are', 'was', 'were', 'be', 'been', 'being', 'have', 'has', 'had', 'having', 'd o', 'does', 'did', 'doing', 'a', 'an', 'the', 'and', 'but', 'if', 'or', 'because', 'as', 'until', 'while', 'of', 'at', 'by', 'for', 'with', 'about', 'again s t', 'between', 'into', 'through', 'during', 'before', 'after', 'above', 'below', 'to', 'from', 'up', 'down', 'in', 'out', 'on', 'off', 'over', 'under', 'agai n', 'further', 'then', 'once', 'here', 'there', 'when', 'where', 'why', 'how', 'all', 'any', 'both', 'each', 'few', 'more', 'most', 'other', 'some', 'such', 'n', 'nor', 'not', 'only', 'own', 'same', 'so', 'than', 'too', 'very', 's', 't', 'can', 'will', 'just', 'don', 'don't', 'should', 'should've', 'now', 'd', 'l l', 'm', 'o', 're', 've', 'y', 'ain', 'aren', 'aren't', 'couldn', 'couldn't', 'didn', 'didn't', 'doesn', 'doesn't', 'hadn', 'hadn't', 'hasn', 'hasn't', 'have n', 'haven't', 'isn', 'isn't', 'ma', 'mightn', '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']
```

In [23]:

```
print(len(stopwords))
print(stopwords)
```

Out[23]:

```
179
['now', 'as', 'weren', 'of', 'between', 'aren't', 'didn't', 'once', 'won', 'about', 'that'll', 'yourself', 'i', 'some', 'needn', 'through', 'and', 'again', 'a in', 'own', 'wouldn't', 'were', 'further', 'mightn', 'these', 'both', 'll', 'your', 'with', 'you've', 'all', 'too', 'y', 'any', 'into', 'or', 'herself', 'at', 'down', 'shouldn', 'in', 'what', 'hadn't', 'the', 'shan', 'such', 'during', 'o', 'nor', 'being', 'you'd', 'was', 'above', 'who', 'after', 'there', 'for', 'di d', 'couldn't', 'm', 'it', 'just', 't', 'don't', 'which', 'him', 'doesn', 'you', 'my', 'more', 'on', 'shouldn't', 'so', 'will', 'no', 'this', 'by', 'we', 'hav ing', 'haven't', 'myself', 'are', 'where', 's', 'didn', 'yourselves', 'hadn', 'ma', 'to', 'its', 'himself', 'is', 'ours', 'mustn't', 'other', 'if', 'from', 'ou selves', 'than', 'not', 'aren', 'isn't', 'up', 'under', 'most', 'wasn', 'hasn', 'me', 'should', 'should've', 'be', 'are', 'hasn't', 'mightn't', 'hers', 'beca use', 'been', 'have', 'wouldn', 'each', 'they', 'weren't', 'she's', 'those', 'you'll', 'she', 'has', 'mustn', 'their', 'yours', 'our', 'itself', 'when', 'belo w', 'does', 'why', 'an', 'few', 'off', 'wasn't', 'whom', 'needn't', 'a', 'd', 'can', 'very', 'doesn't', 'doing', 'don', 'her', 'themselves', 'isn', 'had', 'ov er', 'shan't', 'while', 'against', 'them', 'am', 'be', 'how', 'same', 'haven', 'until', 'before', 'then', 'only', 'won't', 'that', 'you're', 'it's', 'theirs', 'here', 've', 'do', 'couldn', 'out', 'but', 'his']
```

In [18]:

```
stop_words.append('work')
print(stop_words)
```

Out[18]:

```
['i', 'me', 'my', 'myself', 'we', 'our', 'ours', 'ourselves', 'you', 'you're', 'you've', 'you'll', 'you'd', 'your', 'yours', 'yourself', 'yourselves', 'he', 'him', 'his', 'himself', 'she', 'she's', 'her', 'hers', 'herself', 'it', 'it's', 'its', 'itself', 'they', 'them', 'their', 'theirs', 'themselves', 'what', 'wh ich', 'who', 'whom', 'this', 'that', 'that'll', 'these', 'those', 'am', 'is', 'are', 'was', 'were', 'be', 'been', 'being', 'have', 'has', 'had', 'having', 'd o', 'does', 'did', 'doing', 'a', 'an', 'the', 'and', 'but', 'if', 'or', 'because', 'as', 'until', 'while', 'of', 'at', 'by', 'for', 'with', 'about', 'again s t', 'between', 'into', 'through', 'during', 'before', 'after', 'above', 'below', 'to', 'from', 'up', 'down', 'in', 'out', 'on', 'off', 'over', 'under', 'agai n', 'further', 'then', 'once', 'here', 'there', 'when', 'where', 'why', 'how', 'all', 'any', 'both', 'each', 'few', 'more', 'most', 'other', 'some', 'such', 'n', 'nor', 'not', 'only', 'own', 'same', 'so', 'than', 'too', 'very', 's', 't', 'can', 'will', 'just', 'don', 'don't', 'should', 'should've', 'now', 'd', 'l l', 'm', 'o', 're', 've', 'y', 'ain', 'aren', 'aren't', 'couldn', 'couldn't', 'didn', 'didn't', 'doesn', 'doesn't', 'hadn', 'hadn't', 'hasn', 'hasn't', 'have n', 'haven't', 'isn', 'isn't', 'ma', 'mightn', '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', 'work']
```

In [21]:

```
# Removing Stopwords
from nltk.tokenize import sent_tokenize,word_tokenize
from nltk.corpus import stopwords

data="AI was introduced in the year 1956 but it gained popularity recently."
stopwords=set(stopwords.words('english'))
words=word_tokenize(data)
wordsFiltered=[]

for w in words:
    if w not in stopwords:
        wordsFiltered.append(w)

print(wordsFiltered)
```

Out[21]:

```
['AI', 'introduced', 'year', '1956', 'gained', 'popularity', 'recently', '.']
```

In [15]:

```
#Stemming is The Process of Bringing Words Back to Their Root form This Way You End Up With Less Variance in the Data.
#For Example: Connection, Connected,Connected Word Reduce to a Common Word 'Connect'.
import nltk
from nltk.stem import PorterStemmer
#from nltk.tokenize import word_tokenize
stemmer=PorterStemmer()
Input_str="There are several types of stemming Algorithms."
Input_str=nltk.word_tokenize(Input_str)
for word in Input_str:
    print(stemmer.stem(word))
```

Out[15]:

```
there
are
sever
type
of
stem
algorithm
.
```

In [6]:

```
#Lemmatization is Same Like Stemming i.e it Hve Same Goals As Like Stemming But Does So in a More Gramatically Sensitive Way.
import nltk
wn=nltk.WordNetLemmatizer()
ps=nltk.PorterStemmer()
dir(wn)
```

Out[6]:

```
['_class__',
 '_delattr__',
 '_dict__',
 '_dir__',
 '_doc__',
 '_eq__',
 '_format__',
 '_ge__',
 '_getattr__',
 '_gt__',
 '_hash__',
 '_init__',
 '_init_subclass__',
 '_le__',
 '_lt__',
 '_module__',
 '_ne__',
 '_new__',
 '_reduce__',
 '_reduce_ex__',
 '_repr__',
 '_setattr__',
 '_sizeof__',
 '_str__',
 '_subclasshook__',
 '_weakref__',
 'lemmatize']
```

In [38]:

```
print(ps.stem('goose'))
print(ps.stem('geese'))
```

Out[38]:

```
goos
gees
```

In [39]:

```
print(wn.lemmatize('cactus'))
print(wn.lemmatize('cacti'))
```

Out[39]:

```
cactus
cactus
```

In [16]:

```
#Stemming Code
import nltk
from nltk.stem.porter import PorterStemmer
porter_stemmer=PorterStemmer()
text="studies studying cries cry"
tokenization=nltk.word_tokenize(text)
for w in tokenization:
    print("Stemming for {} is {}".format(w,porter_stemmer.stem(w)))
```

Out[16]:

```
Stemming for studies is studi
Stemming for studying is studi
Stemming for cries is cri
Stemming for cry is cri
```

In [17]:

```
# Lemmatization Code
import nltk
from nltk.stem import WordNetLemmatizer
wordnet_lemmatizer=WordNetLemmatizer()
text="studies studying cries cry"
tokenization=nltk.word_tokenize(text)
for w in tokenization:
    print("Lemma for {} is {}".format(w,wordnet_lemmatizer.lemmatize(w)))
```

Out[17]:

```
lemma for studies is study
lemma for studying is studying
lemma for cries is cry
lemma for cry is cry
```

In [24]:

```
# POS Tagging(Part of Speech Tagging) is The Process of attributing a Grammatical Label to Every Part Of Sentences.
import nltk
text=nltk.word_tokenize("It is a pleasant day today")
nltk.pos_tag(text)
```

Out[24]:

```
[('It', 'PRP'),
 ('is', 'VBZ'),
 ('a', 'DT'),
 ('pleasant', 'JJ'),
 ('day', 'NN'),
 ('today', 'NN')]
```

In [34]:

```
nltk.help.upenn_tagset('NNS')
```

Out[34]:

```
NNS: noun, common, plural
undergraduates scotches bric-a-brac products bodyguards facets coasts
divestitures storehouses designs clubs fragrances averages
subjectivists apprehensions muses factory-jobs ...
```

In [38]:

```
nltk.help.upenn_tagset('VB,*')
```

Out[38]:

```
VB: verb, base form
ask assemble assess assign assume atone attention avoid bake balkanize
bank begin behold believe bend benefit bevel beware bless boil bomb
boost brace break bring broil brush build ...

VBD: verb, past tense
dipped pleaded swiped regummed soaked tidied convened halted registered
cushioned exacted snubbed strode aimed adopted belied figgered
speculated wore appreciated contemplated ...

VBG: verb, present participle or gerund
telegraphing stirring focusing angering judging stalling lactating
hankerin' alleging veering capping approaching traveling besieging
encrypting interrupting erasing wincing ...

VBN: verb, past participle
multihulled floundered aerosolized chaired languished panelized sheared
experimented flourished imitated reunited factored condensed used
unsettled primed dubbed desired ...

VBP: verb, presnt tense, not 3rd person singular
predominate wrap resort sue twist spill cure lengthen brush terminate
appear tend stray glisten obtain comprise detest tease attract
emphasize mold postpone sever return wag ...

VBZ: verb, present tense, 3rd person singular
bases reconstructs marks mixes displeases seals carps weaves snatches
slumps stretches authorizes slanders pictures emerges stockpiles
seduces fizzes uses bolsters slaps speaks pleads ...
```

In [39]:

```
import nltk
text=nltk.word_tokenize("I cannot bear the pain of bear")
nltk.pos_tag(text)
```

Out[39]:

```
[('I', 'PRP'),
 ('can', 'MD'),
 ('not', 'RB'),
 ('bear', 'VB'),
 ('the', 'DT'),
 ('pain', 'NN'),
 ('of', 'IN'),
 ('bear', 'NN')]
```

In [51]:

```
# Bag of Words: Bag of Words is The Simplest Way of Structuring Textual data Every Document is Turned into a Word Vector.
import sklearn
from sklearn.feature_extraction.text import CountVectorizer
```

In [48]:

```
phrases=["the quick brown fox jumped over the lazy dog"]
```

In [46]:

```
vect = CountVectorizer()
vect.fit(phrases)
```

Out[46]:

```
CountVectorizer()
```

In [47]:

```
print("Vocabulary size: {}".format(len(vect.vocabulary_)))
print("Vocabulary content:\n {}".format(vect.vocabulary_))
```

Out[47]:

```
Vocabulary size: 8
{'the': 7, 'quick': 6, 'brown': 0, 'fox': 2, 'jumped': 3, 'over': 5, 'lazy': 4, 'dog': 1}
```

In [49]:

```
bag_of_words = vect.transform(phrases)
```

In [50]:

```
print(bag_of_words)
```

Out[50]:

```
(0, 0)      1
(0, 1)      1
(0, 2)      1
(0, 3)      1
(0, 4)      1
(0, 5)      1
(0, 6)      1
(0, 7)      2
```

In [52]:

```
print("bag_of_words as an array:\n{}".format(bag_of_words.toarray()))
```

Out[52]:

```
bag_of_words as an array:
[[1 1 1 1 1 1 2]]
```

In [53]:

```
vect.get_feature_names()
```

Out[53]:

```
['brown', 'dog', 'fox', 'jumped', 'lazy', 'over', 'quick', 'the']
```

In []: