**Feature Engineering**

```
1    from sklearn.feature_extraction import text
2    from os import path
3    ObjRead = open("/content/AST_A.txt", encoding='utf-8')
4    corpus = ObjRead.read();
5    ObjRead.close()
6    input=[corpus]
7    vectorizer = text.CountVectorizer(binary=True).fit(input)
8    vectorized_text = vectorizer.transform(input)
9    print(vectorized_text.todense())
```

[[1 1 1 ... 1 1 1]]

The TF-IDF transformation is a technique that, after counting how many times a token appears in a phrase, divides the value by the number of documents in which the token appears. Using this technique, the vectorizer deems a word less important, even if it appears many times in a text, when it also finds that word in other texts.

```
1    from sklearn.feature_extraction.text import TfidfTransformer
2    TfidF = text.TfidfTransformer(norm='l1')
3    tfidf = TfidF.fit_transform(vectorized_text)
4    phrase=0
5    total=2
6    for word in vectorizer.vocabulary_:
7      pos = vectorizer.vocabulary_[word]
8      value = list(tfidf.toarray()[phrase])[pos]
9      if value !=0:
10       print ("%10s: %0.3f" % (word, value))
11       total += value
12   print ('\nSummed values of a phrase: %0.1f' % total)
```

```
         after: 0.001
        nearly: 0.001
          five: 0.001
         month: 0.001
         delay: 0.001
        caused: 0.001
            by: 0.001
       missing: 0.001
   information: 0.001
     officials: 0.001
            at: 0.001
           the: 0.001
       federal: 0.001
      aviation: 0.001
administration: 0.001
          have: 0.001
         begun: 0.001
     reviewing: 0.001
       license: 0.001
   application: 0.001
           for: 0.001
      proposed: 0.001
     spaceport: 0.001
        camden: 0.001
       kenneth: 0.001
          wong: 0.001
       manager: 0.001
            of: 0.001
           faa: 0.001
     licensing: 0.001
           and: 0.001
    evaluation: 0.001
      division: 0.001
      informed: 0.001
        county: 0.001
    commission: 0.001
      chairman: 0.001
            in: 0.001
        letter: 0.001
         dated: 0.001
          june: 0.001
            28: 0.001
          that: 0.001
       stalled: 0.001
        review: 0.001
           had: 0.001
          been: 0.001
      restarted: 0.001
            on: 0.001
```

```
       feb: 0.001
        12: 0.001
      2019: 0.001
        we: 0.001
 determined: 0.001
        to: 0.001
        be: 0.001
       not: 0.001
  complete: 0.001
    enough: 0.001
```

Using this new TF-IDF model rescales the values of important words and makes them comparable between each text in the corpus. To recover part of the ordering of the text before the BoW transformation, adding n-grams

**Bi-Grams**. The following example uses CountVectorizer to model n-grams in the range of (2, 2), that is, bigrams.

```
1   bigrams = text.CountVectorizer(ngram_range=(2,2))
2   print (bigrams.fit(input).vocabulary_)
```

{'after nearly': 159, 'nearly five': 2772, 'five month': 1564, 'month delay': 2704, 'delay caused': 1190, 'caused by': 881, 'by missing': 784, 'missing information': 26

**Stemming and Wordstop Removal**. The following example demonstrates how to perform stemming and remove stop words from a sentence. It begins by training an algorithm to perform the required analysis using a test sentence. Afterward, the example checks a second sentence for words that appear in the first. The first output shows the stemmed words. All the stop words are missing as well. For example, you don't see the words so, he, all, or the. The second output shows how many times each stemmed word appears in the test sentence.

```
1    import nltk
2    nltk.download('punkt')
3    nltk.download('stopwords')
4    # load data
5    filename = '/content/AST_A.txt'
6    file = open(filename, 'rt')
7    text = file.read()
8    file.close()
9    # split into words
10   from nltk.tokenize import word_tokenize
11   tokens = word_tokenize(text)
12   # convert to lower case
13   tokens = [w.lower() for w in tokens]
14   # remove punctuation from each word
15   import string
16   table = str.maketrans('', '', string.punctuation)
17   stripped = [w.translate(table) for w in tokens]
18   # remove remaining tokens that are not alphabetic
19   words = [word for word in stripped if word.isalpha()]
20   # filter out stop words
21   from nltk.corpus import stopwords
22   stop_words = set(stopwords.words('english'))
23   words = [w for w in words if not w in stop_words]
24   ' '.join(words)
```

[nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data]   Package punkt is already up-to-date!
[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data]   Package stopwords is already up-to-date!
'nearly fivemonth delay caused missing information officials federal aviation administration begun reviewing license application proposed spaceport camden kenneth wong manager faa licensing evaluation division informed camden county commission chairman letter dated june stalled review restarted feb determined application complete enough begin review due four outstanding issues wong wrote letter requested information concerning environmental review mitigation potential risk fire analysis individual risk ability account manage population might exposed risk overflight launch vehicle received additional information june ca
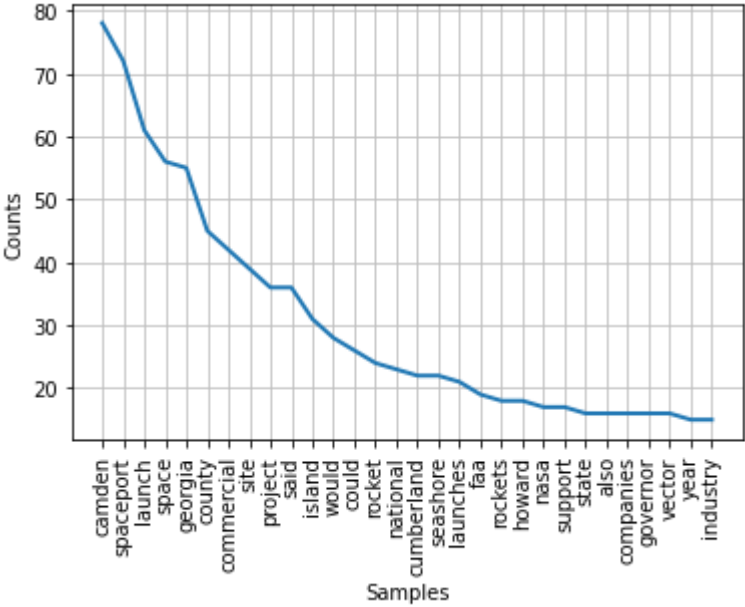
**Frequency Distribution**

```
1    import nltk
2    # Tokenize text
3    from nltk.tokenize import sent_tokenize
4    from os import path
5    ObjRead = open("/content/AST_Clean.txt", encoding='utf-8')
6    text = ObjRead.read();
7    ObjRead.close()
8    tokenized_text=sent_tokenize(text)
9    print(tokenized_text)
10
11   # Tokenize words
12   from nltk.tokenize import word_tokenize
13   tokenized_word=word_tokenize(text)
14   print(tokenized_word)
```

```
15
16    # Frequency Distribution
17    from nltk.probability import FreqDist
18    fdist = FreqDist(tokenized_word)
19    print(fdist)
20
21    # Ten most common words
22    fdist.most_common(10)
23
24    # Frequency Distribution Plots
25    import matplotlib.pyplot as plt
26    fdist.plot(30,cumulative=False)
27    plt.show()
28
```

['nearly fivemonth delay caused missing information officials federal aviation administration begun reviewing license application proposed spaceport camden ke
['nearly', 'fivemonth', 'delay', 'caused', 'missing', 'information', 'officials', 'federal', 'aviation', 'administration', 'begun', 'reviewing', 'license', 'application', 'proposed'
<FreqDist with 1622 samples and 3949 outcomes>



## Text Summarization

```
1    from os import path
2    ObjRead = open("/content/AST_A.txt")
3    txtContent = ObjRead.read();
4    ObjRead.close()
```

```
1    from gensim.summarization import summarize
2    print (summarize(txtContent))
```

After a nearly five-month delay caused by missing information, officials at the Federal Aviation Administration have begun reviewing a license application for the
Camden County is planning to build and operate a commercial spaceport where operators would launch liquid-fueled , small to medium-large vehicles .
We are now at T-minus 1 , a final decision by the FAA is the only outstanding item , Steve Howard , Spaceport Camden Project lead and Camden County admi
I have proudly supported Spaceport Camden from the first time I heard about this amazing project and all it can do for Georgia , said Gov.
The state of Georgia is firmly behind Camden County and we encourage the FAA to swiftly approve its launch site operator's license application .
It's too dangerous , Weinkle said , because the rockets would have to launch over Little Cumberland Island where the private owners and guests can't be force
Camden's rocket launching fantasy will only degrade these precious resources: polluting the salt marsh and tidal creeks , closing waters to commercial and rec
The Camden (Georgia) County Board of Commissioners sees the spaceport as an economic boon , one that could both attract commercial space flight compan
Unlike other communities along Georgia's 100-mile-long coast that have developed thriving tourist communities on their barrier islands , Camden County can't
With the Federal Aviation Administration poised to release a final environmental impact statement on the proposed launch facility by mid-December , with a go-
It's a scenario that paints a picture of as many as a dozen commercial launches a year that would arc across the northern end of Cumberland Island National S
While NASA long has handled space launches for the United States , recent years have seen the rise of private rocket companies that see a profit to be made
We have nothing against rocket ships unless they're being launched over one of the largest maritime wildernesses on the East Coast , said Emily Jones , who
Spaceport Camden is committed to learning from NASA's best practices for nighttime lighting and is already working with environmental groups and federal and
Kevin Lang , an Athens , Georgia , attorney working with a group opposed to the spaceport , pointed out during a phone call that the national seashore is visite
One of our biggest concerns is the fact that a lot of rocket failures do occur early , in the first stage of launch , and a rocket failure over the Cumberland Island N
Analysis of a proposed spaceport in Camden County shows the project has potential to create new jobs , increase tourism , boost the local economy and positi
McKay said the university study looked at the impacts of both construction of a spaceport and operations once launches are conducted at the site .
What makes Camden County such an attractive site when other commercial spaceports have struggled is the central location near other facilities including site
Bob Scaringe , a consultant with AVG Communications , said the commercial space industry continues to grow because private launch sites can launch rockets
Governor Casey each presented benefits Spaceport Camden will offer the State of Georgia and commercial space operators .
The Spaceport Camden project utilizes Georgia's high-technology industry and cutting-edge research capabilities , bringing significant economic opportunities
Governor Casey Cagle echoed Governor Deals sentiments , noting Camden is the ideal region for space vehicle manufacturing due to its location … and also
Spaceport Camden , coupled with NASA's presence in Alabama and Florida along with Georgia's expertise in aerospace engineering and manufacturing has th
Other candidates for Georgia Governor , such as Secretary of State Brian Kemp and former Minority Leader Stacey Abrams , have also given support for Spac
Today , the Camden County Board of Commissioners formally submitted its application for a Launch Site Operator License LSOL to the Federal Aviation Admir
Commission Vice Chairman Gary Blount added , This launch site operator license application has been a strategic priority of the Board of County Commissione
Camden will retain that title again , said Spaceport Camden project leader and County Administrator Steve Howard .
Spaceport Camden is a visionary project that will bring high-paying aerospace jobs to Southeast Georgia while supporting a new wave of STEM education and
Former Speaker of the House and member of the National Space Council Users Advisory Group , Newt Gingrich , joined Camden officials in celebrating the su
Spaceport Camden is an ideal location for the safe launch of the small satellites and rockets that are rapidly becoming the most important segment of this new
The Spaceport , local colleges and the nearby opportunities for the development , manufacture and test of satellites , ground hardware and rockets are valuabl
In June , both US Senators and the entire House Delegation from Georgia sent letters to the FAA supporting the Spaceport Camden draft environmental impac
After a two-year open comment process , a potential commercial space site on Georgia's coast is moving forward with bipartisan political and business leaders
Should the site be approved , the project could almost immediately begin bringing some of the soon-to-be trillion-dollar commercial space market to the state .

But Camden County administrator and project lead Steve Howard says that's a waste .
Additional support notes were submitted by former Speaker of the House and National Space Council Advisory Group member Newt Gingrich ; the private Con
Now , Camden County will submit an official launch site operator license application , after which the FAA has six months to issue a final decision on a launch s
Howard says the FAA will examine all facets of the project — economic impact , environmental factors and more — to determine whether the site will be appro
The support from elected officials , commercial space companies , and space policy advocates all reaches the same conclusion : The United States needs mor
Should Spaceport Camden become a reality , Howard says they have already had multiple commercial space companies reach out expressing interest in using
This launch represents both the first customer-funded launch operation for the new space commercial launch industry , as well as the first launch out of the his
Our historic launch today is a testament to the hard work of the Vector team , as well as support from NASA and Spaceport Camden .
A primary objective of the launch was to demonstrate and evaluate a next-generation engine injector developed through collaborative research with the NASA
Georgia looks forward to working with commercial space companies , like Vector , as we begin the next chapter of space exploration and innovation , he said .
An ongoing environmental assessment conducted by the Federal Aviation Administration due to be completed later this year will likely determine if rockets are
Detractors say the county has already spent millions of dollars on a project that will never be approved because most of the launch trajectories are over Cumbe
Built through partnerships with private companies , the area could become the first exclusively commercial spaceport on the East Coast ; the others in Florida a
Spaceport Camden is strategically located on the coast , which would allow rockets to launch east over mostly open ocean , posing little risk to populated areas
So far , many commercial space companies have shown interest in the place , according to county officials .
Camden County , which the spaceport is named after , is currently working to get the site licensed for launches by the Federal Aviation Administration .
I think we have the opportunity to build the first exclusive non-federal range on the East Coast , Steve Howardm , the Camden County administrator and head c
Then , a little less than five years ago , a space company reached out to the state of Georgia , looking for a new location to launch its rockets .
That initial meeting prompted Camden County officials to meet with other private spaceflight companies to gauge interest in a Georgia spaceport where the boc
For a while , the Spaceport Camden initiative was just a promising idea , but now the state of Georgia has shown it's serious about the project .
In August , a spaceflight startup named Vector launched one of its test rockets from the site — though the vehicle didn't reach orbit .
Currently , 10 sites throughout the US hold FAA licenses to operate as commercial spaceports .
A report done by space consulting firm Astralytical found that Spaceport Camden could be a great home for businesses that not only launch from the site , but

```
1   print ('Summary:')
2   print (summarize(txtContent, word_count=600))
```

Summary:
We are now at T-minus 1 , a final decision by the FAA is the only outstanding item , Steve Howard , Spaceport Camden Project lead and Camden County admi
Camden's rocket launching fantasy will only degrade these precious resources: polluting the salt marsh and tidal creeks , closing waters to commercial and rec
The Camden (Georgia) County Board of Commissioners sees the spaceport as an economic boon , one that could both attract commercial space flight compa
With the Federal Aviation Administration poised to release a final environmental impact statement on the proposed launch facility by mid-December , with a go-
While NASA long has handled space launches for the United States , recent years have seen the rise of private rocket companies that see a profit to be made
Analysis of a proposed spaceport in Camden County shows the project has potential to create new jobs , increase tourism , boost the local economy and positi
The Spaceport Camden project utilizes Georgia's high-technology industry and cutting-edge research capabilities , bringing significant economic opportunities t
Today , the Camden County Board of Commissioners formally submitted its application for a Launch Site Operator License LSOL to the Federal Aviation Admir
The Spaceport , local colleges and the nearby opportunities for the development , manufacture and test of satellites , ground hardware and rockets are valuabl
Additional support notes were submitted by former Speaker of the House and National Space Council Advisory Group member Newt Gingrich ; the private Con
Should Spaceport Camden become a reality , Howard says they have already had multiple commercial space companies reach out expressing interest in using
This launch represents both the first customer-funded launch operation for the new space commercial launch industry , as well as the first launch out of the hist
Camden County , which the spaceport is named after , is currently working to get the site licensed for launches by the Federal Aviation Administration .
Then , a little less than five years ago , a space company reached out to the state of Georgia , looking for a new location to launch its rockets .

## Topic Modeling using NMF

```
1    from sklearn.feature_extraction.text import TfidfVectorizer
2    vectorizer = TfidfVectorizer(stop_words='english')
3    tfidf = vectorizer.fit_transform(input)
4    from sklearn.decomposition import NMF
5    n_topics = 5
6    nmf = NMF(n_components=n_topics, random_state=101).fit(tfidf)
7    feature_names = vectorizer.get_feature_names()
8    n_top_words = 15
9    for topic_idx, topic in enumerate(nmf.components_):
10       print ("Topic #%d:" % (topic_idx+1),)
11       print (" ".join([feature_names[i] for i in
12                 topic.argsort()[:-n_top_words - 1:-1]]))
```

Topic #1:
camden spaceport launch space georgia county commercial site said project island rocket national cumberland launches
Topic #2:
casey support economic small near flight park smart mid location georgia long approved brian lead
Topic #3:
12 running canaveral knowledge ocean parking governor advance owned seashore bring residents 000 officials press
Topic #4:
rural operators central lie followers visited steven story mitigate asked discussed assessment examine offer betting
Topic #5:
vertical homes structures involves officials party sector loses 15 hold december lead attract wallops south

By using the argsort method, you can get the indexes of the top associations, whose high values indicate that they are the most representative words.

```
1    print (nmf.components_[0,:].argsort()[:-n_top_words-1:-1])
2    # Gets top words for topic 0
```

[ 222 1327  779 1323  606  340  289 1298 1235 1102  732 1218  926  353
   781]

## Identify the representative words

```
1  print (vectorizer.get_feature_names()[222])
2  print (vectorizer.get_feature_names()[1327])
3  print (vectorizer.get_feature_names()[779])
4  print (vectorizer.get_feature_names()[1323])
5  print (vectorizer.get_feature_names()[606])
6  print (vectorizer.get_feature_names()[340])
```

```
camden
spaceport
launch
space
georgia
county
```

## Text Classification with SVC

```
1  import numpy as np
2  import pandas as pd
3  from pandas import DataFrame
4  dataset = pd.read_csv('/content/MIA.csv')
5  data=pd.DataFrame(data=dataset, columns=['text','sentiment'])
6  data.head()
```

|   | text | sentiment |
|---|------|-----------|
| 0 | Few parts of the core of Central Florida are s... | 0 |
| 1 | "We are a busy airport," Brown said. | 0 |
| 2 | "There are not too many places in metropolita... | 0 |
| 3 | Brown said his expectation is that flight path... | 0 |
| 4 | "It will narrow the tracks that aircraft will ... | 0 |

```
1  X=data['text']
2  y=data['sentiment']
3  from sklearn.model_selection import train_test_split
4  X_train, X_test, y_train, y_test = train_test_split(X,y, test_size=.5, random_state=125)
```

```
1  from sklearn.feature_extraction import text
2  vectorizer = text.CountVectorizer(ngram_range=(1,2),
3           stop_words='english').fit(X)
4  TfidF = text.TfidfTransformer()
5  X = TfidF.fit_transform(vectorizer.transform(X))
6  Xt = TfidF.transform(vectorizer.transform(X_test))
```

```
1  from sklearn.svm import LinearSVC
2  from sklearn.model_selection import GridSearchCV
3  param_grid = {'C': [0.01, 0.1, 1.0, 10.0, 100.0]}
4  clf = GridSearchCV(LinearSVC(loss='hinge',
5           random_state=101), param_grid)
6  clf = clf.fit(X, y)
7  y_pred=clf.predict(Xt)
8  print ("Best parameters: %s" % clf.best_params_)
```

```
Best parameters: {'C': 0.01}
```

```
1  from sklearn.metrics import accuracy_score
2  solution = clf.predict(Xt)
3  print("Achieved accuracy: %0.3f" %
4       accuracy_score(y_test, solution))
```

```
Achieved accuracy: 0.600
```

```
1  print(X_test[y_test!=solution])
```

```
49   I don't give a damn about good windows, I live...
24   Gilderman said the current distribution of the...
10   The FAA plans to conduct its environmental stu...
14   Members of the Miami International Airport Noi...
45   NextGen's promises to save on fuel and travel ...
33   Public input will be allowed during the enviro...
16   Noise complaints have spiked in communities ne...
19   The federal agency said NextGen will save $15....
28   Satellite-based routes have resulted in flight...
```

```
26    It's going to be nonstop "He predicts concerne...
Name: text, dtype: object
```

**Classification Report**

```
1    from sklearn.metrics import classification_report, confusion_matrix
2    print(classification_report(y_test, y_pred))
```

```
              precision    recall  f1-score   support

           0       0.60      1.00      0.75        15
           1       0.00      0.00      0.00         5
           2       0.00      0.00      0.00         5

    accuracy                           0.60        25
   macro avg       0.20      0.33      0.25        25
weighted avg       0.36      0.60      0.45        25
```

**Confusion Matrix**

```
1    from sklearn.metrics import classification_report, confusion_matrix
2    cm=confusion_matrix(y_test, y_pred)
3    cm[::-1, ::-1]
```

```
array([[ 0,  0,  5],
       [ 0,  0,  5],
       [ 0,  0, 15]])
```

**Text Classification**

How to create a classifier

First create a training and test sample

```
1    train = [
2    ...    ('I love this sandwich.', 'pos'),
3    ...    ('this is an amazing place!', 'pos'),
4    ...    ('I feel very good about these beers.', 'pos'),
5    ...    ('this is my best work.', 'pos'),
6    ...    ("what an awesome view", 'pos'),
7    ...    ('I do not like this restaurant', 'neg'),
8    ...    ('I am tired of this stuff.', 'neg'),
9    ...    ("I can't deal with this", 'neg'),
10   ...    ('he is my sworn enemy!', 'neg'),
11   ...    ('my boss is horrible.', 'neg')
12   ... ]
13   test = [
14   ...    ('the beer was good.', 'pos'),
15   ...    ('I do not enjoy my job', 'neg'),
16   ...    ("I ain't feeling dandy today.", 'neg'),
17   ...    ("I feel amazing!", 'pos'),
18   ...    ('Gary is a friend of mine.', 'pos'),
19   ...    ("I can't believe I'm doing this.", 'neg')
20   ... ]
```

```
1    import nltk
2    nltk.download('punkt')
```

```
[nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data]   Package punkt is already up-to-date!
True
```

```
1    from textblob.classifiers import NaiveBayesClassifier
2    cl = NaiveBayesClassifier(train)
```

Classify a text

```
1    cl.classify("This is an amazing library!")
```

```
'pos'
```

```
1    prob_dist = cl.prob_classify("This one's a doozy.")
```

```
2    print('Overall Sentiment: ',prob_dist.max())
```

     Overall Sentiment:  pos

## Provide the probability distributions

```
1    prob_dist = cl.prob_classify("This one's a doozy.")
2    prob_dist.max() # classify probability of statement
```

     'pos'

```
1    round(prob_dist.prob("pos"), 2) # probability of text being positive
```

     0.63

```
1    round(prob_dist.prob("neg"), 2) # probability of text being negative
```

     0.37

```
1    cl.classify("I am really upset!") # test on a sentence
```

     'neg'

## **Sentiment Analysis** with TextBlob

```
1    from textblob import TextBlob
2    from os import path
3    ObjRead = open("/content/AST_A.txt")
4    txtContent = ObjRead.read();
5    ObjRead.close()
6    testimonial = TextBlob(txtContent)
7    testimonial.sentiment
```

     Sentiment(polarity=0.08435739381391548, subjectivity=0.4159659832920702)

```
1    import nltk
2    nltk.download('movie_reviews')
3    nltk.download('punkt')
4    from textblob import TextBlob
5    from textblob.sentiments import NaiveBayesAnalyzer
6    from os import path
7    ObjRead = open("/content/AST_A.txt")
8    txtContent = ObjRead.read();
9    ObjRead.close()
10   blob = TextBlob(txtContent, analyzer=NaiveBayesAnalyzer())
11   blob.sentiment
```

     [nltk_data] Downloading package movie_reviews to /root/nltk_data...
     [nltk_data]   Unzipping corpora/movie_reviews.zip.
     [nltk_data] Downloading package punkt to /root/nltk_data...
     [nltk_data]   Package punkt is already up-to-date!
     Sentiment(classification='pos', p_pos=1.0, p_neg=6.99671245902884e-109)

## **Tokenize Text**

```
1    from nltk.tokenize import word_tokenize
2    tokenizer = word_tokenize(txtContent)
3    print(tokenizer)
```

     ['After', 'a', 'nearly', 'five-month', 'delay', 'caused', 'by', 'missing', 'information', ',', 'officials', 'at', 'the', 'Federal', 'Aviation', 'Administration', 'have', 'begun', 'reviewin

## **Chunker**

```
1    import nltk
2    nltk.download('conll2000')
3    from textblob import TextBlob
4    from textblob.np_extractors import ConllExtractor
5    extractor = ConllExtractor()
6    blob = TextBlob(txtContent, np_extractor=extractor)
7    blob.noun_phrases
```

     [nltk_data] Downloading package conll2000 to /root/nltk_data...

```
[nltk_data]   Unzipping corpora/conll2000.zip.
WordList(['five-month delay', 'federal aviation administration', 'license application', 'proposed spaceport camden', 'kenneth wong', 'faa licensing', 'evaluation divis
```

## Entity Recognition

```
1    from nltk import word_tokenize, pos_tag, ne_chunk
2    nltk.download('maxent_ne_chunker')
3    nltk.download('words')
4    my_string= 'Antonio Guttierez is the Secretary General of the United Nations in New York in the USA'
5    print(ne_chunk(pos_tag(word_tokenize(my_string))))

     [nltk_data] Downloading package maxent_ne_chunker to
     [nltk_data]     /root/nltk_data...
     [nltk_data]   Unzipping chunkers/maxent_ne_chunker.zip.
     [nltk_data] Downloading package words to /root/nltk_data...
     [nltk_data]   Unzipping corpora/words.zip.
     (S
       (PERSON Antonio/NNP)
       (PERSON Guttierez/NNP)
       is/VBZ
       the/DT
       Secretary/NNP
       General/NNP
       of/IN
       the/DT
       (ORGANIZATION United/NNP Nations/NNPS)
       in/IN
       (GPE New/NNP York/NNP)
       in/IN
       the/DT
       (ORGANIZATION USA/NNP))
```

**POS Tagger.** List of taggers: https://www.learntek.org/blog/categorizing-pos-tagging-nltk-python/

```
1    import nltk
2    nltk.download('averaged_perceptron_tagger')
3    from textblob import TextBlob
4    from textblob.taggers import NLTKTagger
5    nltk_tagger = NLTKTagger()
6    blob = TextBlob(txtContent, pos_tagger=nltk_tagger)
7    blob.pos_tags

     ('Georgians', 'NNPS'),
     ('have', 'VBP'),
     ('recognized', 'VBN'),
     ('that', 'IN'),
     ('our', 'PRP$'),
     ('thriving', 'VBG'),
     ('economy', 'NN'),
     ('and', 'CC'),
     ('fantastic', 'JJ'),
     ('quality', 'NN'),
     ('of', 'IN'),
     ('life', 'NN'),
     ('depend', 'NN'),
     ('on', 'IN'),
     ('the', 'DT'),
     ('clean', 'JJ'),
     ('water', 'NN'),
     ('productive', 'JJ'),
     ('fisheries', 'NNS'),
     ('and', 'CC'),
     ('spectacular', 'JJ'),
     ('scenery', 'NN'),
     ('of', 'IN'),
     ('the', 'DT'),
     ('salt', 'NN'),
     ('marshes', 'NNS'),
     ('barrier', 'JJR'),
     ('islands', 'NNS'),
     ('and', 'CC'),
     ('coastal', 'JJ'),
     ('waterways', 'NNS'),
     ('Camden', 'NNP'),
     (''', 'NNP'),
     ('s', 'NN'),
     ('rocket', 'NN'),
     ('launching', 'VBG'),
     ('fantasy', 'NN'),
     ('will', 'MD'),
     ('only', 'RB'),
     ('degrade', 'VB'),
     ('these', 'DT'),
     ('precious', 'JJ'),
```

('resources', 'NNS'),
      ('polluting', 'VBG'),
      ('the', 'DT'),
      ('salt', 'NN'),
      ('marsh', 'NN'),
      ('and', 'CC'),
      ('tidal', 'JJ'),
      ('creeks', 'NN'),
      ('closing', 'VBG'),
      ('waters', 'NNS'),
      ('to', 'TO'),
      ('commercial', 'JJ'),
      ('and', 'CC'),
      ('recreational', 'JJ'),
      ('fishing', 'NN'),
      ('disrupting', 'VBG'),
      ('traffic', 'NN'),
      ('to', 'TO'),

## Parser

```
1   from textblob import TextBlob
2   from textblob.parsers import PatternParser
3   blob = TextBlob(txtContent, parser=PatternParser())
4   blob.parse()
```

'After/IN/B-PP/B-PNP a/DT/B-NP/I-PNP nearly/RB/I-NP/I-PNP five-month/JJ/I-NP/I-PNP delay/NN/I-NP/I-PNP caused/VBN/B-VP/I-PNP by/IN/B-PP/B-PNP mi
ssing/VBG/B-VP/I-PNP information/NN/B-NP/I-PNP ,/,/O/O officials/NNS/B-NP/O at/IN/B-PP/B-PNP the/DT/B-NP/I-PNP Federal/NNP/I-NP/I-PNP Aviation/NN
P/I-NP/I-PNP Administration/NNP/I-NP/I-PNP have/VBP/B-VP/O begun/VBN/I-VP/O reviewing/VBG/I-VP/O a/DT/B-NP/O license/NN/I-NP/O application/NN/I-
NP/O for/IN/B-PP/O the/DT/O/O proposed/VBN/B-VP/O Spaceport/NNP/B-NP/O Camden/NNP/I-NP/O ./,/O/O\nKenneth/NNP/B-NP/O Wong/NNP/I-NP/O ,/,/
O/O manager/NN/B-NP/O of/IN/B-PP/B-PNP the/DT/B-NP/I-PNP FAA/NNP/I-NP/I-PNP Licensing/NNP/I-NP/I-PNP and/CC/O/O Evaluation/NN/B-NP/O Divisi

## Compare Taggers

```
1   from textblob import Blobber
2   from textblob.taggers import NLTKTagger
3   tb = Blobber(pos_tagger=NLTKTagger())
4   blob1 = tb("This is a blob.")
5   blob2 = tb("This is another blob.")
6   blob1.pos_tagger is blob2.pos_tagger
```

True

## Evaluate overall sentiment

```
1   from textblob import TextBlob
2   blob = TextBlob("The beer is good. But the hangover is horrible.", classifier=cl)
3   blob.classify()
```

'pos'

```
1   cl.show_informative_features(5)
```

```
Most Informative Features
        contains(my) = True          neg : pos   =      1.7 : 1.0
        contains(an) = False         neg : pos   =      1.6 : 1.0
         contains(I) = False         pos : neg   =      1.4 : 1.0
         contains(I) = True          neg : pos   =      1.4 : 1.0
        contains(my) = False         pos : neg   =      1.3 : 1.0
```

```
1   blob = TextBlob("The beer is good. But the hangover is horrible.", classifier=cl)
2   for s in blob.sentences:
3       print(s)
4       print(s.classify())
```

```
The beer is good.
pos
But the hangover is horrible.
neg
```

```
1   print('Classifier Accuracy: %.3f' % cl.accuracy(test))
```

Classifier Accuracy: 0.833

```
1   print('Extract features: ', cl.extract_features(txtContent))
```

Extract features:  {'contains(an)': True, 'contains(best)': True, 'contains(sworn)': False, 'contains(like)': True, 'contains(horrible)': False, 'contains(my)': True, 'con

```
1   from textblob import TextBlob
2   from textblob.classifiers import MaxEntClassifier
3   print('Maximum Entropy Classifier: ', MaxEntClassifier(test))
```

Maximum Entropy Classifier:  <MaxEntClassifier trained on 6 instances>

```
1   new_data = [('She is my best friend.', 'pos'),
2   ...            ("I'm happy to have a new friend.", 'pos'),
3   ...            ("Stay thirsty, my friend.", 'pos'),
4   ...            ("He ain't from around here.", 'neg')]
5   cl.update(new_data)
6   print('Accuracy with new data: ', cl.accuracy(test))
```

Accuracy with new data:  1.0

## Did we update the data?

```
1   cl.update(new_data)
```

True

## WordCloud

```
1   import matplotlib.pyplot as plt
2   %matplotlib inline
3   import nltk
4   nltk.download('punkt')
5   nltk.download('stopwords')
6   # load data
7   filename = '/content/AST_A.txt'
8   file = open(filename, 'rt')
9   text = file.read()
10  file.close()
11  # split into words
12  from nltk.tokenize import word_tokenize
13  word_tokens = word_tokenize(text)
14  # convert to lower case
15  word_tokens = [w.lower() for w in word_tokens]
16  # remove punctuation from each word
17  import string
18  table = str.maketrans('', '', string.punctuation)
19  stripped = [w.translate(table) for w in word_tokens]
20  # remove remaining tokens that are not alphabetic
21  words = [word for word in stripped if word.isalpha()]
22  # filter out stop words
23  from nltk.corpus import stopwords
24  stop_words = set(stopwords.words('english'))
25  words = [w for w in words if not w in stop_words]
26  ' '.join(words)
27  from wordcloud import WordCloud
28  document_file_path = '/content/AST_A.txt'
29  text_from_file = open ( '/content/AST_A.txt') . read ()
30  stop_words = set ( nltk . corpus . stopwords . words ( 'english' ))
31  word_tokenize = word_tokenize( text_from_file )
32  filtered_sentence = [ w for w in word_tokens if not w in stop_words ]
33  wl_space_split = " " . join ( filtered_sentence )
34  my_wordcloud = WordCloud () . generate ( wl_space_split )
35  plt . imshow ( my_wordcloud )
36  plt . axis ( "off" )
37  plt . show ()
```

[nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data]   Package punkt is already up-to-date!
[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data]   Unzipping corpora/stopwords.zip.

## Finding Synonyms

```
1   import nltk
2   nltk.download('wordnet')
3   from nltk.corpus import wordnet
4   syn=wordnet.synsets("automobile")
5   print(syn)
```

```
[nltk_data] Downloading package wordnet to /root/nltk_data...
[nltk_data]   Unzipping corpora/wordnet.zip.
[Synset('car.n.01'), Synset('automobile.v.01')]
```

```
1   import nltk
2   nltk.download('punkt')
3   nltk.download('stopwords')
```

```
[nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data]   Package punkt is already up-to-date!
[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data]   Package stopwords is already up-to-date!
True
```

## Cosine Similarity

```
1   # importing all necessary modules
2   from nltk.tokenize import sent_tokenize, word_tokenize
3   import warnings
4   import nltk
5
6   warnings.filterwarnings(action = 'ignore')
7
8   import gensim
9   from gensim.models import Word2Vec
10
11  #  Reads 'the text' file
12  sample = open("/content/AST_Clean.txt", "r")
13  s = sample.read()
14
15  # Replaces escape character with space
16  f = s.replace("\n", " ")
17
18  data = []
19
20  # iterate through each sentence in the file
21  for i in sent_tokenize(f):
22      temp = []
23
24      # tokenize the sentence into words
25      for j in word_tokenize(i):
26          temp.append(j.lower())
27
28      data.append(temp)
29
30  #Create CBOW model
31  model1 = gensim.models.Word2Vec(data, min_count = 1, size = 100, window = 5)
32
33  #Print results
34  print("Cosine similarity between 'spaceport' " +
35          "and 'environmental' - CBOW : ",
36      model1.similarity('spaceport', 'environmental'))
37  #Print results
38  print("Cosine similarity between 'spaceport' " +
39          "and 'booming' - CBOW : ",
40      model1.similarity('spaceport', 'booming'))
```

```
Cosine similarity between 'spaceport' and 'environmental' - CBOW :  0.09163346
Cosine similarity between 'spaceport' and 'booming' - CBOW :  -0.28629297
```

## Jaccard Index

```
1   import nltk
2   nltk.download('wordnet')
```

```
[nltk_data] Downloading package wordnet to /root/nltk_data...
[nltk_data]   Package wordnet is already up-to-date!
```

True

```
1    from nltk import word_tokenize
2    from nltk.stem import WordNetLemmatizer
3    from sklearn.feature_extraction.text import TfidfVectorizer
4    from sklearn.metrics.pairwise import cosine_similarity
5    lemmatizer = WordNetLemmatizer()
6    pair1 = ["the agency renewed launch certification ","launches will resume soon"]
7    pair2 = ["leaders come together ", "georgia has the support of local leaders"]
8    pair3 = ["the area is in the wilderness", "wild life is significant"]
9    def extract_text_similarity_jaccard (text1, text2):
10       words_text1 = [lemmatizer.lemmatize(word.lower()) for word in word_tokenize(text1)]
11       words_text2 = [lemmatizer.lemmatize(word.lower()) for word in word_tokenize(text2)]
12       nr = len(set(words_text1).intersection(set(words_text2)))
13       dr = len(set(words_text1).union(set(words_text2)))
14       jaccard_sim = nr/dr
15       return jaccard_sim
16   print("Text Similarity: \n", extract_text_similarity_jaccard(pair1[0],pair1[1]))
```

Text Similarity:
 0.125

## Similarity across pairs

```
1    tfidf_model = TfidfVectorizer()
2    corpus = [pair1[0], pair1[1], pair2[0], pair2[1], pair3[0], pair3[1]]
3    tfidf_results = tfidf_model.fit_transform(corpus).todense()
4    cosine_similarity(tfidf_results[0],tfidf_results[1])
```

array([[0.]])