

Week 10

Text Mining — Automatic Understanding of Text

Applied Data Science

Columbia University - Columbia Engineering

Course Agenda



- Week 1: Python Basics: How to Translate Procedures into Codes
- ❖ Week 2: Intermediate Python Data structures for Your Analysis
- Week 3: Relational Databases Where Big Data is Typically Stored
- Week 4: SQL Ubiquitous Database Format/Language
- Week 5: Statistical Distributions The Shape of Data
- Week 6: Sampling When You Can't or Won't Have ALL the Data

- Week 7:Hypothesis Testing Answering Questions about Your Data
- ❖ Week 8: Data Analysis and Visualization Using Python's NumPy for Analysis
- Week 9: Data analysis and visualization Using Python's Pandas for Data Wrangling
- *Week 10: Text Mining Automatic Understanding of Text
- Week 11: Machine learning Basic Regression and Classification
- Week 12: Machine learning Decision Trees and Clustering

Nlkt setup



Packages will start downloading in the background



Output should show "true" once it's downloaded

Out[1]: True

Data that needs to be imported



So, these two files you need to move from this location into the location where we downloaded all the 'nltk' download stuff.



Showing info

https://raw.githubusercontent.com/nltk/nltk_data/gh-pages/index.html

Data Cleaning with Pandas



In [*]: import nltk
 nltk.download()

showing info https://raw.githubusercontent.com/nltk/nltk_data/gh-pages/index.xml

Working with text!

Sentiment Analysis

Identify entities and emotions in a sentence and use these to determine if the entity is being viewed positively or negatively

Easy examples

- . I had an excellent souffle at the restaurant Cavity Maker
- Excellent is a positive word for both the souffle as well as for the restaurant

Not so easy examples

Often, looking at words alone is not enough to figure out the sentiment

The Girl on the Train is an excellent book for a 'stuck at home' snow day
 This one is easy since it includes an explicit positive opinion using a positive word

Showing info

https://raw.githubusercontent.com/nltk/nltk_data/gh-pages/index.html

Data Cleaning with Pandas



Not so easy examples

Often, looking at words alone is not enough to figure out the sentiment

- The Girl on the Train is an excellent book for a 'stuck at home' snow day I This one is easy since it includes an explicit positive opinion using a positive word
- The Girl on the Train is an excellent book for using as a liner for your cat's litter box Not so simple! The positive word "excellent" is used with a negative connotation.
- The Girl on the Train is better than Gone Girl
 The positive word is used as a comparator. Whether the writer likes The Girl on the Train or not depends on what he or she thinks of Gone Girl

Sentiment analysis is generally a starting point in analyzing a text and is then coupled with other techniques (e.g., topic analysis)

Sentiment analysis is usually done using a corpus of positive and negative words

- · Some sources compile lists of positive and negative words
- . Others include the polarity the degree of positivity or negativity of each word

Sources for Sentiment-Coded Words



NRG Emotion Lexicon: words coded into emotional categories (many languages)

http://ptrckprry.com/course/ssd/data/positive-words.txt http://ptrckprry.com/course/ssd/data/negative-words.txt

NRG Emotion Lexicon: words coded into emotional categories (many languages)

http://saifmohammad.com/WebPages/NRC-Emotion-Lexicon.htm

SentiWordNet: List of words weighted by positive or negative sentiments. Includes guidance on how to use the words

http://sentiwordnet.isti.cnr.it/

Vadar Sentiment tool: 7800 words with positive and negative polarity included with python nltk

Our examples

- · Compiled set of 15 reviews each of four neighborhood restaurants
- · Presidential inaugural addresses (from Washington to Trump)
- Some data from yelp (very limited!)

Simple Sentiment Analysis



```
]: def get words(url):
       import requests
       words = requests.get(url).content.decode('latin-1')
       word list = words.split('\n')
       index = 0
       while index < len(word list):
           word = word list[index]
           if ';' in word or not word:
               word list.pop(index)
           else:
               index+=1
       return word list
   #Get lists of positive and negative words
   p url = 'http://ptrckprry.com/course/ssd/data/positive-words.txt'
   n url = 'http://ptrckprry.com/course/ssd/data/negative-words.txt'
   positive words = get words(p url)
   negative words = get words(n url)
```

```
In [3]: len(positive words)
Out[3]: 2006
```

```
; Opinion Lexicon: Positive
; This file contains a list of POSITIVE opinion words (or sentiment words
; This file and the papers can all be downloaded from
    http://www.cs.uic.edu/~liub/FBS/sentiment-analysis.html
; If you use this list, please cite one of the following two papers:
   Minqing Hu and Bing Liu. "Mining and Summarizing Customer Reviews."
       Proceedings of the ACM SIGKDD International Conference on Knowled
       Discovery and Data Mining (KDD-2004), Aug 22-25, 2004, Seattle,
       Washington, USA,
   Bing Liu, Minqing Hu and Junsheng Cheng. "Opinion Observer: Analyzing
       and Comparing Opinions on the Web." Proceedings of the 14th
       International World Wide Web conference (WWW-2005), May 10-14,
       2005, Chiba, Japan.
 Notes:
```

- 1. The appearance of an opinion word in a sentence does not necessar mean that the sentence expresses a positive or negative opinion. See the paper below:
 - Bing Liu. "Sentiment Analysis and Subjectivity." An chapter in Handbook of Natural Language Processing, Second Edition, (editors: N. Indurkhya and F. J. Damerau), 2010.
- 2. You will notice many misspelled words in the list. They are not mistakes. They are included as these misspelled words appear frequently in social media content.

Simple Sentiment Analysis

le monde 5.33% 1.49%

3.85%



Compute sentiment by looking at the proportion of positive and negative words in the text

```
In [4]: from nltk import word tokenize
        cpos = cneg = lpos = lneg = 0
        for word in word tokenize(community):
            if word in positive words:
                cpos+=1
            if word in negative words:
                cneg+=1
        for word in word tokenize(le_monde):
            if word in positive words:
                lpos+=1
            if word in negative words:
                lneg+=1
        print("community {0:1.2f}%\t {1:1.2f}%\t {2:1.2f}%".format(cpos/len(word_tokenize(community))*100,
                                                                cneg/len(word_tokenize(community))*100,
                                                                (cpos-cneg)/len(word tokenize(community))*100))
        print("le monde {0:1.2f}%\t {1:1.2f}%\t {2:1.2f}%".format(lpos/len(word tokenize(le monde))*100,
                                                                lneg/len(word tokenize(le monde))*100,
                                                                (lpos-lneg)/len(word_tokenize(le_monde))*100))
        community 5.09% 1.12%
                                 3.97%
```

Simple Sentiment Analysis Using NRC Data



Simple sentiment analysis using NRC data

- NRC data codifies words with emotions
- 14,182 words are coded into 2 sentiments and 8 emotions

For example, the word abandonment is associated with anger, fear, sadness and has a negative sentiment

- · abandoned anger 1
- · abandoned anticipation 0
- · abandoned disgust 0
- · abandoned fear 1
- abandoned joy 0
- abandoned negative 1
- · abandoned positive 0
- abandoned sadness 1
- · abandoned surprise 0
- · abandoned trust 0

Read the NRC sentiment data

In [7]: emotion_dict



Caveat: We're only looking at one review snippet for each restaurant

- Download yelp python "pip install yelp"
- 2. Register with yelp https://www.yelp.com/developers/manage_apikeys (use anything for the host)
- 3. Copy the various keys into variables as below

We'll see what we can figure out from reviews of restaurants close to Columbia

```
First let's read in the yelp keys

In [8]: CONSUMER_KEY = ""

CONSUMER_SECRET = ""

TOKEN = ""

TOKEN_SECRET = ""
```

Steps to Follow

We need to do a few things

- . Get the latitude and longitude of our location
- . Set up the parameters for what data we want from yelp
- . Query yelp by passing authentication info as well as our parameters
- Extract review snippets

```
In [11]: #We'll use the get lat lng function we wrote way back in week 3
         def get lat lng(address):
             url = 'https://maps.googleapis.com/maps/api/geocode/json?address='
             url += address
             import requests
             response = requests.get(url)
             if not (response.status code == 200):
                 return None
             data = response.json()
             if not( data['status'] == 'OK'):
                 return None
             main result = data['results'][0]
             geometry = main result['geometry']
             latitude = geometry['location']['lat']
             longitude = geometry['location']['lng']
             return latitude, longitude
In [11]: lat,long = get lat lng("Columbia University")
```

```
In [11]: lat,long = get_lat_lng("Columbia University")
In [12]: #Now set up our search parameters
def set_search_parameters(lat,long,radius):
    #See the Yelp API for more details
    params = {}
    params["term"] = "restaurant"
    params["ll"] = "{},{}".format(str(lat),str(long))
    params["radius_filter"] = str(radius) #The distance around our point in metres
    params["limit"] = "10" #Limit ourselves to 10 results
    return params
```







```
In [19]: def get_results(params):
             import rauth
             consumer key = CONSUMER KEY
             consumer_secret = CONSUMER_SECRET
             token = TOKEN
             token secret = TOKEN SECRET
             session = rauth.OAuthlSession(
             consumer key = consumer key
             ,consumer_secret = consumer_secret
             ,access token = token
             ,access_token_secret = token_secret)
             request = session.get("http://api.yelp.com/v2/search",params=params)
             #print(request.url)
             #Transforms the JSON API response into a Python dictionary
             data = request.json()
             session.close()
             return data
```

Rauth Library to Handle Authentications



Extracting snippets

```
In [26]: response['businesses'][0]['id']
Out[26]: 'shaking-crab-new-york-4'
In [27]: all snippets = list()
         for business in response['businesses']:
             name = business['name']
             snippet = business['snippet text']
             id_no = business['id']
             all snippets.append((id no,name,snippet))
         all snippets
Out[27]: [('shaking-crab-new-york-4',
           'Shaking Crab',
           'Easily the best and freshest Cajun-style crawdads ever!
         s langostinos.. \n\nI asked for the hottest...'),
          ('community-food-and-juice-new-york',
           'Community Food & Juice',
           'Was brought here by some friends/regulars who live in t
         or! It was crispy and fluffy and light and...'),
          ('dig-inn-new-york-14',
           'Dig Inn',
           "Delicious! It's a great alternative to fast food; grea
         ions are pretty large, and the seasonings..."),
          / 'la manda nou work 2!
```

Functionalize this

```
In [28]: def get_snippets(response):
    all_snippets = list()
    for business in response['businesses']:
        name = business['name']
        snippet = business['snippet_text']
        id = business['id']
        all_snippets.append((id,name,snippet))
    return all_snippets
```

Functions Analyzing Emotions



We can analyze the emotional content of the review snippet

```
In [33]: print("%-12s %1s\t%1s %1s %1s %1s %1s %1s %1s %1s %1s"%(
                                                  "restaurant", "fear", "trust", "negative", "positive", "joy", "disgust", "anticip",
                                                  "sadness", "surprise"))
                           for snippet in all snippets:
                                      text = snippet[2]
                                      result = emotion analyzer(text)
                                      print("%-12s %1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.
                                                  snippet[1][0:10],result['fear'],result['trust'],
                                                       result['negative'],result['positive'],result['joy'],result['disgust'],
                                                        result['anticipation'], result['sadness'], result['surprise']))
                           restaurant
                                                                fear
                                                                                                trust negative positive joy
                                                                                                                                                                                      disgust anticip sadness
                                                                                                                                                                                                                                                          surprise
                           Shaking Cr
                                                                0.00
                                                                                                0.00
                                                                                                                       0.00
                                                                                                                                              0.00
                                                                                                                                                                                            0.00
                                                                                                                                                                                                                   0.00
                                                                                                                                                                                                                                         0.00
                                                                                                                                                                    0.00
                                                                                                                                                                                                                                                                 0.00
                           Community
                                                                0.04
                                                                                               0.00
                                                                                                                       0.04
                                                                                                                                              0.00
                                                                                                                                                                    0.00
                                                                                                                                                                                            0.00
                                                                                                                                                                                                                  0.00
                                                                                                                                                                                                                                         0.04
                                                                                                                                                                                                                                                                0.00
                                                                0.00
                                                                                                0.04
                                                                                                                       0.00
                                                                                                                                              0.04
                                                                                                                                                                    0.04
                                                                                                                                                                                            0.00
                                                                                                                                                                                                                   0.04
                                                                                                                                                                                                                                         0.00
                                                                                                                                                                                                                                                                 0.00
                           Dig Inn
                           Le Monde
                                                                0.00
                                                                                                0.07
                                                                                                                       0.04
                                                                                                                                              0.04
                                                                                                                                                                    0.04
                                                                                                                                                                                            0.00
                                                                                                                                                                                                                  0.00
                                                                                                                                                                                                                                         0.00
                                                                                                                                                                                                                                                                0.00
                                                                                                                                              0.03
                                                                                                                                                                                                                  0.00
                           Mill Korea
                                                               0.00
                                                                                                0.00
                                                                                                                      0.00
                                                                                                                                                                    0.00
                                                                                                                                                                                            0.00
                                                                                                                                                                                                                                         0.00
                                                                                                                                                                                                                                                                 0.00
                           Symposium
                                                                0.00
                                                                                               0.00
                                                                                                                      0.00
                                                                                                                                              0.00
                                                                                                                                                                    0.00
                                                                                                                                                                                            0.00
                                                                                                                                                                                                                  0.00
                                                                                                                                                                                                                                         0.00
                                                                                                                                                                                                                                                                0.00
                           Tom's Rest
                                                               0.00
                                                                                               0.14
                                                                                                                      0.03
                                                                                                                                              0.24
                                                                                                                                                                    0.14
                                                                                                                                                                                            0.00
                                                                                                                                                                                                                  0.10
                                                                                                                                                                                                                                         0.00
                                                                                                                                                                                                                                                                 0.00
                                                               0.00
                                                                                               0.00
                                                                                                                      0.00
                                                                                                                                              0.03
                                                                                                                                                                    0.00
                                                                                                                                                                                            0.00
                                                                                                                                                                                                                  0.00
                                                                                                                                                                                                                                         0.00
                                                                                                                                                                                                                                                                0.00
                           Mel's Burg
                           Amigos
                                                                0.00
                                                                                               0.00
                                                                                                                      0.00
                                                                                                                                              0.00
                                                                                                                                                                    0.00
                                                                                                                                                                                            0.00
                                                                                                                                                                                                                  0.03
                                                                                                                                                                                                                                         0.00
                                                                                                                                                                                                                                                                 0.00
                                                                                               0.11
                                                                                                                       0.00
                                                                                                                                              0.07
                                                                                                                                                                    0.04
                                                                                                                                                                                            0.00
                                                                                                                                                                                                                  0.07
                                                                                                                                                                                                                                         0.00
                           The Height
                                                               0.00
                                                                                                                                                                                                                                                                 0.00
```

Functions Analyzing Emotions



Functionalize with the yelp

```
In [34]: def comparative emotion analyzer(text tuples):
                                  print("%-20s %1s\t%1s %1s %1s %1s %1s %1s %1s %1s %1s"%(
                                                       "restaurant", "fear", "trust", "negative", "positive", "joy", "disgust", "anticip",
                                                       "sadness", "surprise"))
                                  for text tuple in text tuples:
                                            text = text tuple[2]
                                            result = emotion analyzer(text)
                                           print("%-20s %1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.
                                                       text tuple[1][0:20], result['fear'], result['trust'],
                                                            result['negative'], result['positive'], result['joy'], result['disgust'],
                                                            result['anticipation'], result['sadness'], result['surprise']))
                       #And test it
                       comparative emotion_analyzer(all_snippets)
                                                                                                          trust negative positive joy
                                                                                                                                                                                        disgust anticip sadness surprise
                       restaurant
                                                                             fear
                       Shaking Crab
                                                                              0.00
                                                                                                          0.00
                                                                                                                              0.00
                                                                                                                                                   0.00
                                                                                                                                                                        0.00
                                                                                                                                                                                            0.00
                                                                                                                                                                                                                  0.00
                                                                                                                                                                                                                                       0.00
                                                                                                                                                                                                                                                           0.00
                                                                                                                                                   0.00
                                                                                                                                                                        0.00
                                                                                                                                                                                                                  0.00
                       Community Food & Jui 0.04
                                                                                                          0.00
                                                                                                                              0.04
                                                                                                                                                                                            0.00
                                                                                                                                                                                                                                       0.04
                                                                                                                                                                                                                                                           0.00
                                                                                                                                                                        0.04
                                                                                                                                                                                                                  0.04
                       Dig Inn
                                                                             0.00
                                                                                                          0.04
                                                                                                                              0.00
                                                                                                                                                   0.04
                                                                                                                                                                                            0.00
                                                                                                                                                                                                                                      0.00
                                                                                                                                                                                                                                                           0.00
                                                                             0.00
                                                                                                                                                                                                                 0.00
                       Le Monde
                                                                                                          0.07
                                                                                                                              0.04
                                                                                                                                                   0.04
                                                                                                                                                                        0.04
                                                                                                                                                                                            0.00
                                                                                                                                                                                                                                      0.00
                                                                                                                                                                                                                                                           0.00
                                                                             0.00
                                                                                                                                                   0.03
                                                                                                                                                                                                                 0.00
                                                                                                                                                                                                                                      0.00
                       Mill Korean
                                                                                                          0.00
                                                                                                                              0.00
                                                                                                                                                                        0.00
                                                                                                                                                                                            0.00
                                                                                                                                                                                                                                                           0.00
                       Symposium Greek Rest 0.00
                                                                                                          0.00
                                                                                                                              0.00
                                                                                                                                                   0.00
                                                                                                                                                                        0.00
                                                                                                                                                                                            0.00
                                                                                                                                                                                                                 0.00
                                                                                                                                                                                                                                      0.00
                                                                                                                                                                                                                                                           0.00
                                                                                                                                                                                                                 0.10
                                                                                                                                                                                                                                      0.00
                                                                                                                                                                                                                                                           0.00
                       Tom's Restaurant
                                                                             0.00
                                                                                                          0.14
                                                                                                                              0.03
                                                                                                                                                   0.24
                                                                                                                                                                        0.14
                                                                                                                                                                                            0.00
                       Mel's Burger Bar
                                                                                                                                                                                                                 0.00
                                                                                                                                                                                                                                      0.00
                                                                                                                                                                                                                                                           0.00
                                                                             0.00
                                                                                                          0.00
                                                                                                                              0.00
                                                                                                                                                   0.03
                                                                                                                                                                        0.00
                                                                                                                                                                                            0.00
                       Amigos
                                                                             0.00
                                                                                                                              0.00
                                                                                                                                                   0.00
                                                                                                                                                                        0.00
                                                                                                                                                                                                                 0.03
                                                                                                                                                                                                                                      0.00
                                                                                                                                                                                                                                                           0.00
                                                                                                          0.00
                                                                                                                                                                                            0.00
                                                                                                                                                                                                                 0.07
                                                                                                                                                                                                                                      0.00
                       The Heights Bar & Gr 0.00
                                                                                                          0.11
                                                                                                                              0.00
                                                                                                                                                   0.07
                                                                                                                                                                        0.04
                                                                                                                                                                                            0.00
                                                                                                                                                                                                                                                           0.00
```

```
In [35]: def analyze nearby restaurants(address, radius):
             lat, long = get lat lng(address)
              params = set search parameters(lat,long,radius)
             response = get results(params)
              snippets = get snippets(response)
              comparative emotion analyzer(snippets)
         #And test it
         analyze nearby restaurants ("Community Food and Juice", 200)
         restaurant
                               fear
                                           trust negative positive joy
                                                                          disgust anticip sadness
         Shaking Crab
                               0.00
                                           0.00
                                                   0.00
                                                           0.00
                                                                    0.00
                                                                            0.00
                                                                                    0.00
                                                                                             0.00
                                                                   0.00
                                                                                    0.00
         Community Food & Jui 0.04
                                           0.00
                                                   0.04
                                                           0.00
                                                                            0.00
                                                                                             0.04
                                                                   0.04
         Dig Inn
                               0.00
                                           0.04
                                                   0.00
                                                           0.04
                                                                            0.00
                                                                                    0.04
                                                                                             0.00
         Le Monde
                               0.00
                                           0.07
                                                   0.04
                                                           0.04
                                                                    0.04
                                                                            0.00
                                                                                    0.00
                                                                                             0.00
         Mill Korean
                                                           0.03
                                                                    0.00
                                                                                    0.00
                               0.00
                                           0.00
                                                   0.00
                                                                            0.00
                                                                                             0.00
         Symposium Greek Rest 0.00
                                          0.00
                                                   0.00
                                                           0.00
                                                                   0.00
                                                                            0.00
                                                                                    0.00
                                                                                             0.00
         Tom's Restaurant
                                                   0.03
                                                           0.24
                                                                            0.00
                               0.00
                                           0.14
                                                                    0.14
                                                                                    0.10
                                                                                             0.00
         Mel's Burger Bar
                               0.00
                                           0.00
                                                   0.00
                                                           0.03
                                                                    0.00
                                                                            0.00
                                                                                    0.00
                                                                                             0.00
         Amigos
                               0.00
                                           0.00
                                                   0.00
                                                           0.00
                                                                   0.00
                                                                            0.00
                                                                                    0.03
                                                                                             0.00
         The Heights Bar & Gr 0.00
                                           0.11
                                                   0.00
                                                           0.07
                                                                   0.04
                                                                            0.00
                                                                                    0.07
                                                                                             0.00
In [36]: #Test it on some other place
         analyze nearby restaurants ("221 Baker Street", 200)
                               fear
                                           trust negative positive joy
                                                                          disgust anticip sadness
         restaurant
                                                                                    0.03
                               0.00
                                           0.03
                                                   0.00
                                                           0.03
                                                                   0.03
                                                                            0.00
                                                                                             0.00
         Opa
                                           0.15
                                                   0.10
                                                           0.25
                                                                    0.15
                                                                            0.05
                                                                                    0.25
         Beast + Barrel
                               0.05
                                                                                             0.05
         Brewhaus
                               0.00
                                           0.04
                                                   0.00
                                                           0.04
                                                                    0.04
                                                                            0.00
                                                                                    0.04
                                                                                             0.00
         Easy Bistro & Bar
                               0.00
                                           0.06
                                                   0.00
                                                           0.06
                                                                    0.06
                                                                            0.00
                                                                                    0.06
                                                                                             0.03
         Aretha Frankensteins 0.00
                                          0.00
                                                   0.00
                                                           0.03
                                                                   0.00
                                                                            0.00
                                                                                    0.00
                                                                                             0.00
                               0.00
                                                   0.00
                                                           0.00
                                                                   0.00
                                                                            0.00
                                                                                    0.00
                                                                                             0.00
         Milk & Honey
                                           0.00
                                           0.00
                                                   0.00
                                                           0.03
                                                                   0.00
                                                                            0.00
                                                                                    0.00
                                                                                             0.00
         Tony's Pasta Shop an 0.00
```

Simple Analysis: Word Clouds



Let's see what sort of words the snippets use

- · First we'll combine all snippets into one string
- . Then we'll generate a word cloud using the words in the string
- You may need to install wordcloud using pip
- · pip install wordcloud

```
In [ ]: all_snippets
In [ ]: text=''
    for snippet in all_snippets:
        text+=snippet[2]
    text
```





```
In [40]: from wordcloud import WordCloud, STOPWORDS
import matplotlib.pyplot as plt
%matplotlib inline

wordcloud = WordCloud(stopwords=STOPWORDS,background_color='white',width=3000,height=3000).generate(text)

plt.imshow(wordcloud)
plt.axis('off')
plt.show()
```



Detailed Analysis



Let's do a detailed comparison of local restaurants

I've saved a few reviews for each restaurant in four directories

We'll use the PlainTextCorpusReader to read these directories

PlainTextCorpusReader reads all matching files in a directory and saves them by file-ids

```
In []: import nltk
    from nltk.corpus import PlaintextCorpusReader
    community_root = "data/community"
    le_monde_root = "data/le_monde"
    community_files = "community.*"
    le_monde_files = "le_monde.*"
    heights_root = "data/heights"
    heights_files = "heights.*"
    amigos_root = "data/amigos"
    amigos_files = "amigos.*"
    community_data = PlaintextCorpusReader(community_root,community_files)
    le_monde_data = PlaintextCorpusReader(le_monde_root,le_monde_files)
    heights_data = PlaintextCorpusReader(heights_root,heights_files)
    amigos_data = PlaintextCorpusReader(amigos_root,amigos_files)
```

```
In [42]: amigos data.fileids()
Out[42]: ['amigos.1',
           'amigos.10',
           'amigos.11',
           'amigos.12',
           'amigos.13'.
           'amigos.14',
           'amigos.15'.
           'amigos.16',
           'amigos.17',
           'amigos.18',
           'amigos.19',
           'amigos.2',
           'amigos.20',
           'amigos.21',
           'amigos.3',
           'amigos.4',
           'amigos.5',
           'amigos.6',
           'amigos.7',
           1 --- 1 --- 01
   In [43]: amigos data.raw()
```

Out[43]: 'I see all these bad reviews, but speaking for recein the least—and I\'m here about once a weet 've never had a bad dish, although some are mere, and the margaritas are among the best I\'ve exat specials during the week—I believe one day it sounds so impossible even as I type it but I\'inco de Mayo but we wanted something close to he end telling me about this place a while back. So that it was packed BUT the receptionist explains \n\nWe found some seats at the bar and got some reat, the food was great, the drinks were AMA; FINITELY be going back with even more people! I



We need to modify comparitive_emotion_analyzer to tell it where the restaurant name and the text is in the tuple

```
In [ ]: def comparative_emotion_analyzer(text_tuples,name_location=1,text_location=2):
                          "restaurant", "fear", "trust", "negative", "positive", "joy", "disgust", "anticip",
                                            "sadness", "surprise"))
                          for text tuple in text tuples:
                                   text = text tuple[text location]
                                   result = emotion analyzer(text)
                                   print("%-20s %1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.2f\t%1.
                                            text tuple[name location][0:20],result['fear'],result['trust'],
                                                result['negative'], result['positive'], result['joy'], result['disgust'],
                                                result['anticipation'], result['sadness'], result['surprise']))
                  #And test it
                  comparative_emotion_analyzer(all_snippets)
                                                                                                                                               disgust anticip sadness surprise
                  restaurant
                                                            fear
                                                                                  trust negative positive joy
                  Shaking Crab
                                                            0.00
                                                                                  0.00
                                                                                                  0.00
                                                                                                                  0.00
                                                                                                                                   0.00
                                                                                                                                                   0.00
                                                                                                                                                                   0.00
                                                                                                                                                                                    0.00
                                                                                                                                                                                                    0.00
                 Community Food & Jui 0.04
                                                                                  0.00
                                                                                                  0.04
                                                                                                                  0.00
                                                                                                                                   0.00
                                                                                                                                                   0.00
                                                                                                                                                                   0.00
                                                                                                                                                                                   0.04
                                                                                                                                                                                                    0.00
                  Dig Inn
                                                                                  0.04
                                                                                                  0.00
                                                                                                                  0.04
                                                                                                                                   0.04
                                                                                                                                                   0.00
                                                                                                                                                                   0.04
                                                                                                                                                                                   0.00
                                                                                                                                                                                                    0.00
                                                            0.00
                 Le Monde
                                                            0.00
                                                                                  0.07
                                                                                                  0.04
                                                                                                                  0.04
                                                                                                                                   0.04
                                                                                                                                                   0.00
                                                                                                                                                                   0.00
                                                                                                                                                                                   0.00
                                                                                                                                                                                                    0.00
                                                                                                                                                                                                    0.00
                 Mill Korean
                                                            0.00
                                                                                  0.00
                                                                                                  0.00
                                                                                                                  0.03
                                                                                                                                   0.00
                                                                                                                                                   0.00
                                                                                                                                                                   0.00
                                                                                                                                                                                   0.00
                 Symposium Greek Rest 0.00
                                                                                  0.00
                                                                                                  0.00
                                                                                                                  0.00
                                                                                                                                   0.00
                                                                                                                                                   0.00
                                                                                                                                                                   0.00
                                                                                                                                                                                   0.00
                                                                                                                                                                                                    0.00
                 Tom's Restaurant
                                                            0.00
                                                                                  0.14
                                                                                                  0.03
                                                                                                                  0.24
                                                                                                                                   0.14
                                                                                                                                                   0.00
                                                                                                                                                                   0.10
                                                                                                                                                                                   0.00
                                                                                                                                                                                                    0.00
                 Mel's Burger Bar
                                                            0.00
                                                                                                  0.00
                                                                                                                  0.03
                                                                                                                                  0.00
                                                                                                                                                   0.00
                                                                                                                                                                   0.00
                                                                                                                                                                                   0.00
                                                                                                                                                                                                    0.00
                                                                                  0.00
                 Amigos
                                                            0.00
                                                                                  0.00
                                                                                                  0.00
                                                                                                                  0.00
                                                                                                                                   0.00
                                                                                                                                                   0.00
                                                                                                                                                                   0.03
                                                                                                                                                                                   0.00
                                                                                                                                                                                                    0.00
                 The Heights Bar & Gr 0.00
                                                                                  0.11
                                                                                                  0.00
                                                                                                                  0.07
                                                                                                                                   0.04
                                                                                                                                                   0.00
                                                                                                                                                                   0.07
                                                                                                                                                                                   0.00
                                                                                                                                                                                                    0.00
    In [45]: restaurant data = [('community',community data.raw()),('le monde',le monde data.raw())
                                                               ,('heights',heights data.raw()), ('amigos',amigos data.raw())]
                         comparative emotion analyzer(restaurant data,0,1)
                                                                                              trust negative positive joy
                         restaurant
                                                                       fear
                                                                                                                                                              disgust anticip sadness surprise
                         community
                                                                      0.00
                                                                                              0.03
                                                                                                                                0.05
                                                                                                                                                 0.03
                                                                                                                                                                                    0.02
                                                                                                               0.01
                                                                                                                                                                  0.01
                                                                                                                                                                                                     0.01
                                                                                                                                                                                                                      0.01
                         le monde
                                                                                                               0.01
                                                                                                                                                 0.02
                                                                                                                                                                  0.00
                                                                                                                                                                                    0.02
                                                                                                                                                                                                     0.00
                                                                                                                                                                                                                      0.01
                                                                      0.00
                                                                                              0.03
                                                                                                                                0.04
                         heights
                                                                                                                                                 0.03
                                                                      0.00
                                                                                              0.03
                                                                                                               0.01
                                                                                                                                0.04
                                                                                                                                                                  0.00
                                                                                                                                                                                    0.03
                                                                                                                                                                                                     0.01
                                                                                                                                                                                                                      0.01
                                                                                                                                                 0.03
                                                                                                                                                                  0.01
                         amigos
                                                                      0.01
                                                                                              0.03
                                                                                                               0.01
                                                                                                                                0.04
                                                                                                                                                                                    0.02
                                                                                                                                                                                                     0.01
                                                                                                                                                                                                                      0.01
```

Simple Analysis: Complexity



We'll look at four complexity factors

- · average word length: longer words adds to complexity
- · average sentence length: longer sentences are more complex (unless the text is rambling!)
- · vocabulary: the ratio of unique words used to the total number of words (more variety, more complexity)

token: A sequence (or group) of characters of interest. For e.g., in the below analysis, a token = a word

- . Generally: A token is the base unit of analysis
- . So, the first step is to convert text into tokens and nltk text object

Simple Analysis: Complexity



```
In [46]: #Construct tokens (words/sentences) from the text
         text = le monde data.raw()
         import nltk
         from nltk import sent tokenize, word tokenize
         sentences = nltk.Text(sent tokenize(text))
         print(len(sentences))
         words = nltk.Text(word_tokenize(text))
         print(len(words))
         188
         2595
In [47]: num_chars=len(text)
         num_words=len(word_tokenize(text))
         num sentences=len(sent tokenize(text))
         vocab = {x.lower() for x in word tokenize(text)}
         print(num chars,int(num chars/num words),int(num words/num sentences),(len(vocab)/num words))
         12332 4 13 0.29132947976878615
```

Let's functionalize this

```
In [48]: def get_complexity(text):
             num chars=len(text)
             num_words=len(word_tokenize(text))
             num sentences=len(sent tokenize(text))
             vocab = {x.lower() for x in word tokenize(text)}
             return len(vocab),int(num chars/num words),int(num words/num sentences),len(vocab)/num words
In [49]: get complexity(le monde data.raw())
Out[49]: (756, 4, 13, 0.29132947976878615)
In [51]: for text in restaurant data:
             (vocab, word size, sent size, vocab to text) = get complexity(text[1])
             print("{0:15s}\t{1:1.2f}\t{3:1.2f}\t{4:1.2f}\".format(text[0],vocab,word_size,sent_size,vocab_to_text))
          community
                         1029.00 4.00
                                        16.00
                                               0.28
                                               0.29
          le monde
                         756.00 4.00
                                        13.00
          heights
                         720.00 4.00
                                        16.00
                                               0.28
          amigos
                         792.00 4.00
                                        15.00 0.24
```





```
In [52]: texts = restaurant data
         from wordcloud import WordCloud, STOPWORDS
         import matplotlib.pyplot as plt
         *matplotlib inline
         #Remove unwanted words
         #As we look at the cloud, we can get rid of words that don't make sense by adding them to this variable
         DELETE WORDS = []
         def remove words(text string, DELETE_WORDS=DELETE_WORDS):
             for word in DELETE WORDS:
                  text string = text string.replace(word, ' ')
             return text string
         #Remove short words
         MIN LENGTH = 0
         def remove short words(text string,min length = MIN LENGTH):
             word list = text string.split()
             for word in word list:
                 if len(word) < min length:
                      text string = text string.replace(' '+word+' ',' ',1)
             return text string
         #Set up side by side clouds
         COL_NUM = 2
         ROW NUM = 2
         fig, axes = plt.subplots(ROW NUM, COL NUM, figsize=(12,12))
         for i in range(0,len(texts)):
             text string = remove words(texts[i][1])
             text string = remove short words(text string)
             ax = axes[i%2]
              ax = axes[i//2, i%2] #Use this if ROW NUM >=2
             ax.set title(texts[i][0])
             wordcloud = WordCloud(stopwords=STOPWORDS, background_color='white', width=1200, height=1000, max_words=20).generate(text)
             ax.imshow(wordcloud)
             ax.axis('off')
         plt.show()
```

output

Community



heights



Le monde



Amigos



Each restaurant can be analyzed based on the word clouds

Ntlk: Python's natural language toolkit



Comparing complexity of restaurant reviews won't get us anything useful. Let's look at something more useful

Ntlk documentation link:

http://www.nlkt.org/api/nlkt.html

Commands cheat sheet:

https://blogs.princeton.edu/etc/files/2014/03/Text-Analysis-with NLTK-Cheatsheet.pdf

Ntlk book

http://www.nltk.org/book

nltk contains a large corpora of pre-tokenized text

Load it using the command:

nltk.download()

Ntlk: Python's natural language toolkit



Often, a comparitive analysis helps us understand text better

Let's look at US President inaugural speeches

Copy the files 2013-Obama.txt and 2017-Trump.txt to the nltk_data/corpora/inaugural directory. nltk_data should be under your home directory

20

```
In [55]: inaugural.fileids()
Out[55]: ['1789-Washington.txt',
           '1793-Washington.txt',
           '1797-Adams.txt',
           '1801-Jefferson.txt',
           '1805-Jefferson.txt',
           '1809-Madison.txt',
           '1813-Madison.txt',
           '1817-Monroe.txt',
           '1821-Monroe.txt',
           '1825-Adams.txt'.
In [56]: inaugural.raw('1861-Lincoln.txt')
Out[56]: 'Fellow-Citizens of the United Stat
         ou to address you briefly and to ta
         o be taken by the President "before
         present for me to discuss those mat
         pprehension seems to exist among th
```

```
In [57]: texts = [('trump',inaugural.raw('2017-Trump.txt')),
                  ('obama',inaugural.raw('2009-Obama.txt')+inaugural.raw('2013-Obama.txt')),
                  ('jackson',inaugural.raw('1829-Jackson.txt')+inaugural.raw('1833-Jackson.txt')),
                  ('washington',inaugural.raw('1789-Washington.txt')+inaugural.raw('1793-Washington.txt'))]
             (vocab, word size, sent size, vocab to text) = get complexity(text[1])
             print("{0:15s}\t{1:1.2f}\t{2:1.2f}\t{4:1.2f}\t{4:1.2f}\".format(text[0],vocab,word size,sent size,vocab to text))
         trump
         obama
                        1349.00 5.00
                                       25.00
                                              0.27
                        813.00 5.00
         jackson
                                       45.00
                                              0.33
         washington
                        636.00 5.00
                                       62.00
                                              0.38
In [58]: from nltk.corpus import inaugural
           sentence lengths = list()
           for fileid in inaugural.fileids():
                sentence lengths.append(get complexity(' '.join(inaugural.words(fileid)))[2])
           plt.plot(sentence_lengths)
Out[58]: [<matplotlib.lines.Line2D at 0x1286352b0>]
            60
            50
            30
```

Digression Plot





https://tartarus.org/martin/PorterStemmer/

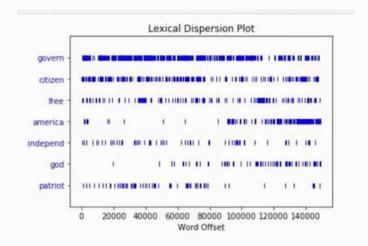
Digression Plot



We may want to use word stems rather than the part of speect form

- · For example: patriot, patriotic, patriotism all express roughly the same idea
- nltk has a stemmer that implements the "Porter Stemming Algorithm" https://tartarus.org/martin/PorterStemmer/
- · We'll push everything to lowercase as well

```
In []: from nltk.stem.porter import PorterStemmer
    p_stemmer = PorterStemmer()
    text = inaugural.raw()
    striptext = text.replace('\n\n', '')
    striptext = striptext.replace('\n', '')
    sentences = sent_tokenize(striptext)
    words = word_tokenize(striptext)
    text = nltk.Text([p_stemmer.stem(i).lower() for i in words])
    text.dispersion_plot(["govern", "citizen", "free", "america", 'independ', 'god', 'patriot'])
```



Weighted Word Analysis Using Vader



Vader contains a list of 7500 features weighted by how positive or negative they are

It uses these features to calculate stats on how positive, negative and neutral a passage is

And combines these results to give a compound sentiment (higher = more positive) for the passage

Human trained on twitter data and generally considered good for informal communication

10 humans rated each feature in each tweet in context from -4 to +4

Calculates the sentiment in a sentence using word order analysis

"marginally good" will get a lower positive score than "extremely good"

Computes a "compound" score based on heuristics (between -1 and +1)

Includes sentiment of emoticons, punctuation, and other 'social media' lexicon elements

```
In [62]: !pip install vaderSentiment

Collecting vaderSentiment

Downloading vaderSentiment-2.5-py2.py3-none-any.whl (102kB)

100% | 100k | 112kB 3.7MB/s ta 0:00:01

Installing collected packages: vaderSentiment

Successfully installed vaderSentiment-2.5
```



In [65]: restaurant data



```
In [63]: from vaderSentiment.vaderSentiment import SentimentIntensityAnalyzer
In []: headers = ['pos','neg','neu','compound']
    texts = restaurant_data
    analyzer = SentimentIntensityAnalyzer()
    for iI in range(len(texts)):
        name = texts[i][0]
        sentences = sent_tokenize(texts[i][1])
        pos=compound=neu=neg=0
        for sentence in sentences:
            vs = analyzer.polarity_scores(sentence)
            pos+=vs['pos']/(len(sentences))
            compound+=vs['compound']/(len(sentences))
            neu+=vs['neu']/(len(sentences))
            neg+=vs['neg']/(len(sentences))
            print(name,pos,neg,neu,compound)
```

```
Out[65]: [('community',
           'I ate here Monday night for a working dinner with a coworker. I immediately saw its reason for na
         ner\'s fandom of the show of the same name: the tablet setup are nearly all communal long tables. The
         w non-shared tables which is fine as not only is NYC filled with such setup but there\'s a livelines
         area vibe which prefers this setup. There\'s a shortage of such nice sit-down restaurants in the Col
         s place is not only capitalizing on its presence but doing so with great style and flare.\n\nI order
         r which was very juicy, perfectly lean, well sauced, cooked as desired (medium-well, for those who w
         d with a dill pickle, a handmade cole slaw that lacked mayo (which is good as mayo\'s pure trans-fat
         of fries. My coworker got the salmon sandwich, same review applies.\n\nAwesome bar alongside, but we
         e it. Vast amounts of seating and an open kitchen where tons of food is prepared live makes for a fu
         ccidentally left my headphones, called to ask if it was found, and fetched them the day later. The s
         ful, kind, and soft-spoken about recovering and returning it that I\'m duly humbled. Thanks, all!\nJ
         uite good, and I have no complaints in that area. The reason for only two stars is the cost of that
         y overpriced. My wife and I went here for breakfast. We each ordered blueberry pancakes and a cup of
         fruit, no juice, no bacon. The pancakes were $13 and the coffee was $4, so the initial bill was $34.
         and I left a tip of $6--resulting in a final total of $43 for a breakfast of flapjacks and coffee. N
         NY and costs here are higher--but $43! Ugh! There are a lot of nice places in the city where you c
         d breakfast for 1/3rd to 1/2 less than Community. I recommend you seek them out. \nThis is such a goo
         e food, service & ambiance was exactly what we were looking for. We sat at an adorable table on the
```

Weighted Word Analysis Using Vader

```
In [67]: headers = ['pos','neg','neu','compound']
    texts = restaurant_data
    analyzer = SentimentIntensityAnalyzer()
    for i in range(len(texts)):
        name = texts[i][0]
        sentences = sent_tokenize(texts[i][1])
        pos=compound=neu=neg=0
        for sentence in sentences:
            vs = analyzer.polarity_scores(sentence)
            pos+=vs['pos']/(len(sentences))
            compound+=vs['compound']/(len(sentences))
            neu+=vs['neu']/(len(sentences))
            neg+=vs['neu']/(len(sentences))
            print(name,pos,neg,neu,compound)
```

community 0.20305855855855845 0.02542342342342342343 0.771527027027026 0.3362022522522523 le monde 0.1749627659574468 0.04054255319148937 0.7845053191489363 0.2215510638297871 heights 0.18679354838709677 0.03675483870967742 0.7764387096774197 0.28074129032258055 amigos 0.19447887323943652 0.0497042253521127 0.755830985915493 0.25482488262910796

And functionalize this as well

```
In [68]: def vader comparison(texts):
             from vaderSentiment.vaderSentiment import SentimentIntensityAnalyzer
             headers = ['pos', 'neg', 'neu', 'compound']
             print("Name\t",' pos\t', 'neg\t', 'neu\t', 'compound')
             analyzer = SentimentIntensityAnalyzer()
             for i in range(len(texts)):
                 name = texts[i][0]
                 sentences = sent tokenize(texts[i][1])
                 pos=compound=neu=neg=0
                 for sentence in sentences:
                     vs = analyzer.polarity_scores(sentence)
                     pos+=vs['pos']/(len(sentences))
                     compound+=vs['compound']/(len(sentences))
                     neu+=vs['neu']/(len(sentences))
                     neg+=vs['neg']/(len(sentences))
                 print('%-10s'%name,'%1.2f\t'%pos,'%1.2f\t'%neg,'%1.2f\t'%neu,'%1.2f\t'%compound)
In [69]: vader comparison(restaurant data)
                                          compound
                                          0.34
         community 0.20 0.03
                                 0.77
         le monde 0.17 0.04
                                 0.78
                                          0.22
         heights
                   0.19 0.04
                                 0.78
                                          0.28
         amigos
                   0.19 0.05
                                 0.76
                                          0.25
```



Named Entities

COLUMBIA ENGINEERING EXECUTIVE EDUCATION

People, places, organizations

Named entities are often the subject of sentiments so identifying them can be very useful

Named entity detection is based on Part-of-speech tagging of words and chunks (groups of words)

- · Start with sentences (using a sentence tokenizer)
- · tokenize words in each sentence
- · chunk them. ne_chunk identifies likely chunked candidates (ne = named entity)
- · Finally build chunks using nltk's guess on what members of chunk represent (people, place, organization

```
In [71]: meaningful sents = list()
         i=0
         for sentence in sentences:
             if 'service' in sentence:
                  i+=1
                 meaningful sents.append((i,sentence))
         vader comparison(meaningful sents)
         Name
                          neg
                                   neu
                                           compound
                    pos
                    0.00
                          0.00
                                   1.00
                                           0.00
                    0.11 0.15
                                   0.73
                                           -0.17
                    0.53 0.00
                                   0.47
                                           0.84
                                   0.71
                    0.28 0.00
                                           0.49
                                   0.78
                                           0.49
                    0.23 0.00
                    0.36 0.00
                                   0.64
                                           0.74
```

```
In [70]: en={}
         try:
             sent_detector = nltk.data.load('tokenizers/punkt/english.pickle')
             sentences = sent_detector.tokenize(community_data.raw().strip())
             for sentence in sentences:
                     tokenized = nltk.word tokenize(sentence)
                     tagged = nltk.pos tag(tokenized)
                     chunked = nltk.ne chunk(tagged)
                     for tree in chunked:
                          if hasattr(tree, 'label'):
                             ne = ' '.join(c[0] for c in tree.leaves())
                             en[ne] = [tree.label(), '.join(c[1] for c in tree.leaves())]
          except Exception as e:
             print(str(e))
          import pprint
          pp = pprint.PrettyPrinter(indent=4)
          pp.pprint(en)
             'America': ['GPE', 'NNP'],
              'Awesome': ['GPE', 'NNP'],
             'BEST': ['ORGANIZATION', 'NNP'],
             'Bill': ['PERSON', 'NN'],
             'Boston': ['GPE', 'NNP'],
             'Bottomless': ['GPE', 'NNP'],
             'Broadway': ['GPE', 'NNP'],
             'Brooklyn': ['GPE', 'NNP'],
             'Brunch': ['PERSON', 'NNP'],
             'CU': ['ORGANIZATION', 'NNP'],
             'Came': ['GPE', 'NN'],
```

Named Entities



```
In [72]: def get affect(text,word,lower=False):
             import nltk
             from vaderSentiment.vaderSentiment import SentimentIntensityAnalyzer
             analyzer = SentimentIntensityAnalyzer()
             sent detector = nltk.data.load('tokenizers/punkt/english.pickle')
             sentences = sent detector.tokenize(text.strip())
             sentence count = 0
             running total = 0
             for sentence in sentences:
                 if lower: sentence = sentence.lower()
                 if word in sentence:
                      vs = analyzer.polarity scores(sentence)
                     running total += vs['compound']
                      sentence count += 1
             if sentence count == 0: return 0
             return running total/sentence count
In [73]: get affect(community data.raw(), 'service', True)
Out[73]: 0.4321428571428571
In [74]: get affect(le monde data.raw(), 'service', True)
Out[74]: 0.22273125
```

The nltk function concordance returns text fragments around a word

```
In [81]: nltk.Text(community_data.words()).concordance('service',100)

Displaying 7 of 7 matches:

This is such a good brunch spot! The food, service & ambiance was exactly what we were looking f I dug right in due to brunch hunger and slow service. I'd come back again, but be sure to have. Will be exploring more items on the menu. Service is great - our water was also filled, our co All in all, a pleasant experience with nice service and nice company. It's not every day I get place was packed. Even though it was packed service was very good. The beat salad and rice bowl e with a friend - we were seated promptly and service was attentive. The decor was beautiful and the juice non site coffee. Thanks to the kind service as well, i will remember this restaurant as
```

Text Summarization



Text summarization is useful because you can generate a short summary of a large piece of text automatically. Then, these summaries can serve as an input into a topic analyzer to figure out what the main topic of the text is

A naïve form of summary is to identify the most frequent words in a piece of text and use the occurrence of these words in sentence to rate the importance of a sentence.

('i', 73), ('the', 46),

('food', 20),

('good', 19),

('place', 15),

('ordered', 13), ('pancakes', 13), ('fresh', 13), ('came', 13), ('my', 12), ('coffee', 12), ('delicious', 12),

('nice', 13),

('we', 10),

('community', 10),

('brunch', 20),

```
In []: from nltk.tokenize import word_tokenize
    from nltk.tokenize import sent_tokenize
    from nltk.probability import FreqDist
    from nltk.corpus import stopwords
    from collections import OrderedDict
    import pprint
```

Then prep the text. Get did of end of line chars

```
In []: text = community_data.raw()
    summary_sentences = []
    candidate_sentences = {}
    candidate_sentence_counts = {}
    striptext = text.replace('\n\n', '')
    striptext = striptext.replace('\n', '')
```

Construct a list of words after getting rid of unimportant ones and numbers

Construct word frequencies and choose the most common n (20)

```
word_frequencies = FreqDist(lowercase_words)
most_frequent_words = FreqDist(lowercase_words).most_common(20)
pp = pprint.PrettyPrinter(indent=4)
pp.pprint(most_frequent_words)
```

Text Summarization



```
In [93]: sentences = sent tokenize(striptext)
          for sentence in sentences:
              candidate sentences[sentence] = sentence.lower()
          candidate sentences
Out[93]: {'"I have a degree from Columbia, and now I have to get one from America."': '"i have
          i have to get one from america."',
           '- Jeff Winger, Community.': '- jeff winger, community.',
           '- The truffle egg dish: had sausage, potatoes, some leafy greens, sautéed mushroor
          he truffle egg dish: had sausage, potatoes, some leafy greens, sautéed mushrooms, as
           '- Troy Barnes Man I love that show.': '- troy barnes man i love that show.',
           '-7 grains waffle: had vanilla cream, raspberry, cherry, and toasted nuts.': '-7 gr
          aspberry, cherry, and toasted nuts.',
           '15 minuets wait for the wait staff, 45 minuets before we got the wrong dish!': '1!
          45 minuets before we got the wrong dish!',
           '2 HOURS AND NO FOOD!': '2 hours and no food!',
In [94]: most frequent words
                                          In [96]: for long, short in candidate sentences.items():
Out[94]: [('i', 73),
                                                       count = 0
          ('the', 46),
                                                       for freq word, frequency score in most frequent words:
          ('food', 20),
                                                           if freq word in short:
          ('brunch', 20),
                                                               count += frequency score
          ('good', 19),
                                                       candidate sentence counts[long] = count
          ('place', 15),
                                          In [97]: candidate sentence counts
                                         Out[97]: {'"I have a degree from Columbia, and now I have to get one
                                                    '- Jeff Winger, Community.': 83,
                                                    '- The truffle egg dish: had sausage, potatoes, some leafy
                                                    '- Troy Barnes Man I love that show.': 73,
                                                    '-7 grains waffle: had vanilla cream, raspberry, cherry, a
```

Text Summarization



Packaging all this into a function

```
In [ ]: def build naive summary(text):
            from nltk.tokenize import word tokenize
            from nltk.tokenize import sent tokenize
            from nltk.probability import FreqDist
            from nltk.corpus import stopwords
            from collections import OrderedDict
            summary sentences = []
            candidate sentences = {}
            candidate sentence counts = {}
            striptext = text.replace('\n\n', '')
            striptext = striptext.replace('\n', '')
            words = word tokenize(striptext)
            lowercase words = [word.lower() for word in words
                              if word not in stopwords.words() and word.isalpha()]
            word frequencies = FreqDist(lowercase words)
            most frequent words = FreqDist(lowercase words).most common(20)
            sentences = sent tokenize(striptext)
            for sentence in sentences:
                candidate sentences[sentence] = sentence.lower()
            for long, short in candidate sentences.items():
                count = 0
                for freq word, frequency score in most frequent words:
                    if freq word in short:
                        count += frequency score
                        candidate sentence counts[long] = count
            sorted sentences = OrderedDict(sorted(
                                candidate sentence counts.items(),
                                key = lambda x: x[1],
                                reverse = True)[:4])
            return sorted sentences
```

```
In [101]: summary = '\n'.join(build_naive_summary(community_data.raw()))
print(summary)
```

I've come here several times with a friend for brunch and once for dinner -- we reakfast foods available; my favorite is the brioche French toast with blackber fles.

There are a lot of nice places in the city where you can get a very good breakf ty.

Came here for brunch with my wife after she found the good review, and it did r The beans were seasoned perfectly, as was the rest of the tomato sauce, The tor with delicious guacamole and sour cream.

```
In [102]: summary = '\n'.join(build_naive_summary(le_monde_data.raw()))
    print(summary)
```

i would recommend going for brunch over any other meal, as everything i've orde shakshuka is great, and you can't anything else like it in the area.

I would love to try them out for dinner, as I said, the service was great and v I would give them a better review because the food was good but the waiting real The food was really good and they are known for their brunch menu.

Summarizing President's speech

In [104]: build naive summary(inaugural.raw('2013-Obama.txt'))

Gensim - Text Summarizer



```
In [107]: from wordcloud import WordCloud, STOPWORDS
          import matplotlib.pyplot as plt
          *matplotlib inline
          import nltk
          from nltk.corpus import PlaintextCorpusReader
          from nltk import sent tokenize, word tokenize
          from nltk.book import *
In [108]: import nltk
          from nltk.corpus import PlaintextCorpusReader
          community root = "data/community"
          le monde root = "data/le monde"
          community files = "community.*"
          le monde files = "le monde.*"
          heights root = "data/heights"
          heights files = "heights.*"
          amigos root = "data/amigos"
          amigos files = "amigos.*"
          community data = PlaintextCorpusReader(community root,community files)
          le monde data = PlaintextCorpusReader(le monde root,le monde files)
          heights data = PlaintextCorpusReader(heights root,heights files)
          amigos data = PlaintextCorpusReader(amigos root, amigos files)
In [109]: type(community data)
Out[109]: nltk.corpus.reader.plaintext.PlaintextCorpusReader
```

```
In []: text = community_data.raw()
    summary_sentences = []
    candidate_sentence = {}
    candidate_sentence_counts = {}
    striptext = text.replace('\n\n', '')
    striptext = striptext.replace('\n', '')

In [*]: import gensim.summarization
In []: import gensim.summarization

In []: import gensim.summarization

In []: import gensim.summarization

In []: import gensim.summarization

In []: import gensim.summarization
In []: import gensim.summarization
```

```
In [113]: summary = gensim.summarization.summarize(striptext, word_count=100)
    print(summary)
```

I've come here several times with a friend for brunch and once for dinner -- we've both really enjoyed a lot of the b reakfast foods available; my favorite is the brioche French toast with blackberry and lemon curd δ she loves heir waf fles.

I ordered the Country Breakfast and the eggs were delicious and fluffy, the biscuit was moist and flavorful, the carr ot hash browns were warm and comforting, and the ham was sweet and juicy.

This place is really good for breakfast and has great pancakes, sausages, eggs.

Topic Modeling



The goal of topic modeling is to identify the major concepts underlying a piece of text Topic modelling uses Unsupervised Learning. No prior knowledge is necessay, though it is helpful in cleaning up results!

LDA: Latent Dirichlet Allocation Model

- Identifies potential topics using pruning techniques like "upward closure"
- Computes conditional probablities for topic word sets
- Identifies the most likely topics
- Does this over multiple passes probabilities for topic
- Good intuitive explanation: http://blog.echen.me/2011/08/22/introduction-to-latent-dirichlet-allocation/

```
In [116]: from gensim import corpora
from gensim.models.ldamodel import LdaModel
from gensim.parsing.preprocessing import STOPWORDS
import pprint
```

Topic Modeling



Prepare the text

```
In [117]: text = PlaintextCorpusReader("data/", "Nikon coolpix 4300.txt").raw()
          striptext = text.replace('\n\n', '')
          striptext = striptext.replace('\n', ' ')
          sentences = sent tokenize(striptext)
          #words = word tokenize(striptext)
          #tokenize each sentence into word tokens
          texts = [[word for word in sentence.lower().split()
                  if word not in STOPWORDS and word.isalnum()]
                  for sentence in sentences]
          len(texts)
Out[117]: 360
 In [118]: texts
 Out[118]:
            [['annotated', 'minging', 'hu', 'bing'],
             ['department',
               'sicence',
              'university',
              'illinois'.
              'chicago',
              'product',
              'nikon',
              'coolpix',
```

Create a (word, frequency) dictionary for each word in the text

Do the LDA



Parameters:

- Number of topics: The number of topics we have generated. The larger the document, the more desirable topics
- Passes: The LDA model makes through the document. More passes, slower analysis

See results

```
In [134]: pp = pprint.PrettyPrinter(indent=4)
    pp.pprint(lda.print_topics(num_words=3))

[       (0, '0.039*"camera" + 0.020*"pictures" + 0.019*"card"'),
            (1, '0.035*"camera" + 0.025*"nikon" + 0.017*"4300"'),
            (2, '0.021*"pics" + 0.010*"use" + 0.010*"camera"'),
            (3, '0.049*"camera" + 0.022*"battery" + 0.012*"little"'),
            (4, '0.038*"camera" + 0.019*"picture" + 0.016*"use"')]
```

See results

Do the LDA



We're using sentences as documents here, so this is less than ideal

Making Sense of the Topic



Draw wordclouds

```
In [ ]: def draw wordcloud(lda,topicnum,min size=0,STOPWORDS=[]):
            word list=[]
            prob total = 0
            for word, prob in lda.show topic(topicnum, topn=50):
                prob total +=prob
            for word, prob in lda.show topic(topicnum, topn=50):
                if word in STOPWORDS or len(word) < min size:
                     continue
                freq = int(prob/prob total*1000)
                alist=[word]
                word list.extend(alist*freg)
            from wordcloud import WordCloud, STOPWORDS
            import matplotlib.pyplot as plt
            %matplotlib inline
            text = ' '.join(word list)
            wordcloud = WordCloud(stopwords=STOPWORDS, background color='white', width=3000, height=3000).generate('
            plt.imshow(wordcloud)
            plt.axis('off')
            plt.show()
```

```
In [143]: draw_wordcloud(lda,2)
```



Roughly,

- · Ida looks for candidate topics assuming that there are many such candidates
- · looks for words related to the candidate topics
- · assign probablilites to those words

Presidential Addresses



Let's look at Presidential addresses to see what sorts of topics emerge from there

- · Each document will be analyzed for topic
- · The corpus will consist of 58 documents, one per presidential address

```
In [ ]: REMOVE_WORDS = {'shall', 'generally', 'spirit', 'country', 'people', 'nation', 'nations', 'great', 'better'}
        #Create a word dictionary (id, word)
        texts = [[word for word in sentence.lower().split()
                if word not in STOPWORDS and word not in REMOVE WORDS and word.isalnum()]
                for sentence in sentences |
        dictionary = corpora.Dictionary(texts)
        #Create a corpus of documents
        text list = list()
        for fileid in inaugural.fileids():
            text = inaugural.words(fileid)
            doc=list()
            for word in text:
                if word in STOPWORDS or word in REMOVE WORDS or not word.isalpha() or len(word) <5:
                     continue
                doc.append(word)
            text list.append(doc)
        by address corpus = [dictionary.doc2bow(text) for text in text list]
```

```
In [146]: len(by_address_corpus)
Out[146]: 58
```





```
In [147]: lda = LdaModel(by address corpus,
                        id2word=dictionary,
                        num topics=20,
                        passes=10)
In [148]: pp = pprint.PrettyPrinter(indent=4)
          pp.pprint(lda.print_topics(num words=10))
                  '0.032*"course" + 0.030*"things" + 0.028*"knowledge" + 0.023*"process" '
                  '+ 0.020*"years" + 0.019*"order" + 0.019*"little" + 0.017*"thought" + '
                  '0.017*"support" + 0.017*"familiar"'),
                  '0.003*"world" + 0.002*"power" + 0.002*"place" + 0.001*"future" + '
                  '0.001*"children" + 0.001*"course" + 0.001*"years" + 0.001*"service" + '
                  '0.001*"office" + 0.001*"order"'),
              ( 2,
                  '0.105*"power" + 0.026*"control" + 0.024*"given" + 0.019*"hands" + '
                  '0.016*"state" + 0.014*"subject" + 0.013*"produce" + 0.013*"important" '
                  '+ 0.013*"department" + 0.013*"intended"'),
                  '0.043*"years" + 0.031*"power" + 0.028*"action" + 0.019*"future" + '
                  '0.018*"office" + 0.017*"service" + 0.017*"period" + 0.016*"opinion" + '
                  '0.015*"trust" + 0.015*"ability"'),
              ( 4,
                  '0.033*"office" + 0.022*"state" + 0.022*"highest" + 0.022*"recommend" '
                  '+ 0.022*"performance" + 0.012*"believe" + 0.011*"right" + '
                  '0.011*"power" + 0.011*"control" + 0.011*"action"'),
                  '0.035*"world" + 0.028*"support" + 0.022*"service" + 0.021*"power" + '
                  '0.017*"years" + 0.017*"given" + 0.017*"important" + 0.016*"order" + '
```





```
In [ ]: len(by address corpus)
In [152]: from operator import itemgetter
          sorted(lda.get document topics(by address corpus[56], minimum probability=0, per word topics=False), key=itemgetter(1), reve
Out[152]: [(8, 0.51051380624923337),
            (11, 0.48152159191111937),
            (5, 0.00044247788373801006),
            (17, 0.00044247788355170585),
            (19, 0.00044247788285562616),
            (12, 0.00044247788249148379),
            (13, 0.00044247788194282128),
            (3, 0.00044247788192305253),
            (0, 0.00044247788185107641),
            (2, 0.0004424778815170437),
            (9, 0.00044247788129546736),
            (6, 0.00044247788065216576),
            (10, 0.00044247787985220237),
            (15, 0.00044247787900034022),
                                                      In [153]: draw wordcloud(lda,11)
            (4, 0.00044247787791450085),
            (14, 0.00044247787663700809),
                                                                      sliple sliplefriends friends *... week
            (1, 0.00044247787610634767),
                                                                       place place secrifice secrifice
            (7, 0.00044247787610627199),
            (16, 0.00044247787610620217),
            (18, 0.00044247787610619501)1
                                                         In [ ]: print(lda.show topic(12,topn=5))
                                                                  print(lda.show topic(18,topn=5))
```

Presidential Addresses: Create the Model



Given a corpus of documents, when a new document arrives, find the document that is the most similar

```
In [ ]: doc list = [community data,le monde data,amigos data,heights data]
        all text = community data.raw() + le monde data.raw() + amigos data.raw() + heights data.raw()
         documents = [doc.raw() for doc in doc list]
        texts = [[word for word in document.lower().split()
                if word not in STOPWORDS and word.isalnum()]
                 for document in documents]
        dictionary = corpora.Dictionary(texts)
        corpus = [dictionary.doc2bow(text) for text in texts]
In [ ]: from gensim.similarities.docsim import Similarity
        from gensim import corpora, models, similarities
        lsi = models.LsiModel(corpus, id2word=dictionary, num topics=2)
        doc = """
        Many, many years ago, I used to frequent this place for their amazing french toast.
        It's been a while since then and I've been hesitant to review a place I haven't been to in 7-8 years...
        but I passed by French Roast and, feeling nostalgic, decided to go back.
        It was a great decision.
        Their Bloody Mary is fantastic and includes bacon (which was perfectly cooked!!), olives,
        cucumber, and celery. The Irish coffee is also excellent, even without the cream which is what I ordered.
        Great food, great drinks, a great ambiance that is casual yet familiar like a tiny little French cafe.
        I highly recommend coming here, and will be back whenever I'm in the area next.
        Juan, the bartender, is great!! One of the best in any brunch spot in the city, by far.
        vec bow = dictionary.doc2bow(doc.lower().split())
        vec lsi = lsi[vec bow]
        index = similarities.MatrixSimilarity(lsi[corpus])
        sims = index[vec lsi]
        sims = sorted(enumerate(sims), key=lambda item: -item[1])
```





```
In [156]: sims
Out[156]: [(1, 0.98765731), (0, 0.95548683), (3, 0.79290682), (2, 0.7698077)]
 In [ ]: doc="""
          I went to Mexican Festival Restaurant for Cinco De Mayo because I had been there years
          prior and had such a good experience. This time wasn't so good. The food was just
          mediocre and it wasn't hot when it was brought to our table. They brought my friends food out
          10 minutes before everyone else and it took forever to get drinks. We let it slide because the place was
          packed with people and it was Cinco De Mayo. Also, the margaritas we had were slamming! Pure tequila.
          But then things took a turn for the worst. As I went to get something out of my purse which was on
          the back of my chair, I looked down and saw a huge water bug. I had to warn the lady next to me because
          it was so close to her chair. We called the waitress over and someone came with a broom and a dustpan and
          swept it away like it was an everyday experience. No one seemed phased.
          Even though our waitress was very nice, I do not think we will be returning to Mexican Festival again.
          It seems the restaurant is a shadow of its former self.
          vec bow = dictionary.doc2bow(doc.lower().split())
          vec lsi = lsi[vec bow]
          index = similarities.MatrixSimilarity(lsi[corpus])
          sims = index[vec lsi]
          sims = sorted(enumerate(sims), key=lambda item: -item[1])
```



www.emeritus.org