

marta_hack

October 29, 2016

```
In [35]: !pip install mpld3
```

Requirement already satisfied (use --upgrade to upgrade): mpld3 in /Library/Frameworks/Python.framework/Versions/2.7/Resources/Python.framework/Versions/2.7/Headers/Python.h

```
In [36]: import pandas as pd
```

```
import numpy as np
import pandas as pd
import nltk
import re
import os
import codecs
from sklearn import feature_extraction
import mpld3
from bs4 import BeautifulSoup
from nltk import clean_html
data = pd.read_csv('https://raw.githubusercontent.com/johnymontana/harvard')
```

```
In [37]: data.head(5)
```

```
Out[37]:
```

	groupId	name	id	urlkey
0	18616327	Open Source	563	opensource
1	18616327	CSS	1973	css
2	18616327	JavaScript	7029	javascript
3	18616327	New Technology	9696	newtech
4	18616327	Web Technology	10209	web

```
In [38]: #members data
```

```
#members_data = pd.read_csv('https://raw.githubusercontent.com/johnymontana/harvard')
```

```
In [39]: #members_data.head(5)
```

```
In [40]: #events
```

```
events_data = pd.read_csv('https://raw.githubusercontent.com/johnymontana/harvard')
```

```
In [41]: events_data.head(5)
```

```

Out[41]:
      id                name      time  utc_offset
0  3051299  Boston Java September Meetup  1095116400000  -14400000
1  2962561  Boston Wikipedia September Meetup  1095116400000  -14400000
2  3301994  Boston Wikipedia October Meetup  1097535600000  -14400000
3  3301979  Boston Java October Meetup  1098140400000  -14400000
4  3466780  Boston Wikipedia November Meetup  1099958400000  -18000000

      group_id  venue_id  status  \
0      87071      NaN  past
1     165993      NaN  past
2     165993      NaN  past
3      87071      NaN  past
4     165993      NaN  past

      description
0  <p>This is an introductory Java meetup - I hop...
1      NaN
2      NaN
3  <p>How about we'll talk about the Web Framewor...
4      NaN

```

```
In [42]: #groups
```

```
      #groups_data = pd.read_csv('https://raw.githubusercontent.com/johnymontana
```

```
In [43]: #groups_data.head(5)
```

```
In [44]: #join data
```

```
      #data1 = pd.merge(data, members_data, left_on='groupId', right_on='groupI
```

```
In [45]: #data1.head(5)
```

```
In [46]: #data2 = pd.merge(data1, events_data, left_on='groupId', right_on='group_
```

```
In [47]: #data2.head()
```

```
In [48]: events_data.shape
```

```
Out[48]: (9311, 8)
```

```
In [49]: df = events_data[['name', 'description']] # Slice to remove redundant colu
```

```
In [50]: df.shape
```

```
Out[50]: (9311, 2)
```

```
In [51]: print (df['name'][:200])
```

0 Boston Java September Meetup
 1 Boston Wikipedia September Meetup
 2 Boston Wikipedia October Meetup
 3 Boston Java October Meetup
 4 Boston Wikipedia November Meetup
 5 Boston Java November Meetup
 6 Boston Wikipedia December Meetup
 7 Boston Java December Meetup
 8 Boston Wikipedia January Meetup
 9 Boston Java January Meetup
 10 Boston Wikipedia
 11 Boston Wikipedia Pre-V-Day Meetup
 12 Boston Java February Meetup
 13 The Boston Wikipedia March Meetup
 14 The Boston Java March Meetup
 15 Post-Ides party planning
 16 The Boston Wikipedia April Meetup
 17 The Boston Java April Meetup
 18 Evening meetup (with jwales?)
 19 The Boston Wikipedia May Meetup
 20 The Boston Java May Meetup
 21 Wikimagic planning meeting
 22 The Boston Wikipedia June Meetup
 23 The Boston Wikipedia July Meetup
 24 The Boston Java July Meetup
 25 The Boston Wikipedia August Meetup
 26 The Boston Wikipedia September Meetup
 27 The Boston Java September Meetup
 28 The Boston Wikipedia September Meetup
 29 The Boston Wikipedia October Meetup
 ...
 170 The Cambridge Python October Meetup
 171 The Boston Java October Meetup
 172 Recent High Tech University Spin-Outs - Who is...
 173 IT Job Seekers Greater Boston (Java, .NET, SQL...
 174 The Boston Wikipedia November Meetup
 175 The Cambridge Enterprise Web 2.0 November Meetup
 176 GLOBAL ENTREPRENEURSHIP WEEK Nov.17-23 - NORT...
 177 The Boston Java November Meetup: Wicket
 178 GLOBAL ENTREPRENEURSHIP WEEK Nov.17-23 - NORT...
 179 DIGITAL MEDIA - (GLOBAL ENTREPRENEURSHIP WEEK)
 180 IT Job Seekers Greater Boston (Java, .NET, SQL...
 181 The Boston Wikipedia December Meetup
 182 The Cambridge Enterprise Web 2.0 December Meetup
 183 The Boston Java December Meetup: OpenMRS
 184 Nanotechnology - Timing Is Everything
 185 The Web Content Management Systems Meetup
 186 IT Job Seekers Greater Boston (Java, .NET, SQL...

```

187             The Boston Wikipedia January Meetup
188     The Cambridge Enterprise Web 2.0 January Meetup
189             Boston WordPress Kick-off
190             Going to the GYM - Web Media & Search
191             Cambridge Python January Meetup
192             First Saturday meetup of the year
193             The Boston Java January Meetup
194     IT Job Seekers Greater Boston (Java, .NET, SQL...
195             The Boston Wikipedia February Meetup
196     The Cambridge Enterprise Web 2.0 February Meetup
197             Blogging @ The Brewery
198     GREENING OF AMERICA - Conversion of Bio-Agricu...
199             The Cambridge Python February Meetup
Name: name, dtype: object

```

```
In [52]: print (df['description'][:200])
```

```

0     <p>This is an introductory Java meetup - I hop...
1                                             NaN
2                                             NaN
3     <p>How about we'll talk about the Web Framewor...
4                                             NaN
5     <p>November Java Meetup - if nothing else, let...
6                                             NaN
7     <p>Our monthly gathering to discuss the busine...
8                                             NaN
9     <p>Happy New Year! We invite anybody to share ...
10    <p>Meet Jimbo after his conference for early d...
11    <p>Discuss V-Pedia, sweethearts, sweetmeats, f...
12                                             NaN
13                                             NaN
14                                             NaN
15    <p>It's time to have a real Spring party. The ...
16    <p>Yahoo!ligans unite... come have fun in Club...
17    <p>Dear Java Enthusiasts,<br/>I wanted to thro...
18    <p>Jimmy Wales will be in town next Tuesday, w...
19                                             NaN
20    <p>Hello all Java enthusiasts!</p> <p>As we pr...
21    <p>Planning for upcoming speakers and events.<...
22    <p>Meet inside Pho Pasteur. Bring a notebook a...
23    <p>There will be a speaker this evening; so co...
24    <p>Hello, Java Meetup members!</p> <p>At our n...
25    <p>There will likely be a speaker this evening...
26    <p>j here: SJ has a conflict at about 7 pm. De...
27    <p>Hello, Java Meetup members!<br/>this month,...
28    <p>Discussion of great local ice-cream joints...
29    <p>Discussion of great local ice-cream joints...

```

```

...
170 <p>Follow Snake Charmers,</p> <p>October's Mee...
171 <p>Hi,</p> <p><br/>We don't have a presenter f...
172 <p>Recent High Tech University Spin-Outs - Who...
173 <p>6:00pm to 6:30pm: Refreshments and networki...
174 <p>We're moving to dinner from ice cream, and ...
175 <p>We meet at <a href="http://maps.google.com/...
176 <p>Northeastern University is hosting events f...
177 <p>Dave Rafkind will lead a discussion of Wick...
178 <p>Northeastern University is hosting events f...
179 <p>[<b>The Youth Movement: Digital Media's Lea...
180 <p>6:00pm to 6:30pm: Refreshments and networki...
181 <p>We're moving to dinner from ice cream, and ...
182 <p>We meet at <a href="http://maps.google.com/...
183 <p>Darius Jazayeri, lead software developer at...
184 <p>Nanotechnology - Timing is Everything : Kno...
185 <p>This meeting will be postponed until furthe...
186 <p>6:00pm to 6:30pm: Refreshments and networki...
187 <p>We're moving to dinner, and this one will a...
188 <p>We meet at <a href="http://maps.google.com/...
189 NaN
190 <p>GOING TO THE GYM : WEB-MEDIA & SEARCH -...
191 <p>Alexander Fairley will be presenting <a hre...
192 <p>Lots of food and project discussions!</p> <...
193 <p>Meeting for dinner & drinks at the asga...
194 <p>6:00pm to 6:30pm: Refreshments and networki...
195 <p>We're moving to dinner, and this one will a...
196 <p>We meet at <a href="http://maps.google.com/...
197 <p>The Cambridge Brewing Co. is a great venue ...
198 <p><b>GREENING OF AMERICA - The Conversion Of ...
199 <p>Topic: "A Whirlwind Excursion through Pytho...
Name: description, dtype: object

```

```
In [53]: clean_df = df.dropna(subset = ['description'])
```

```
In [54]: import json
import nltk
```

```
# Download nltk packages used in this example
nltk.download('stopwords')
```

```
#blog_data = load_json('blog')
```

```
# Customize your list of stopwords as needed. Here, we add common
# punctuation and contraction artifacts.
```

```

stop_words = nltk.corpus.stopwords.words('english') + [
    '.',
    ',',
    '--',
    '\s',
    '?',
    ')',
    '(',
    ':',
    '\'',
    '\re',
    '"',
    '-',
    '}',
    '{',
    '<p>',
    '>',
    '<',
    '/p',
    'p',
    '\xa0',
    u'--',
]

```

```

[nltk_data] Downloading package stopwords to
[nltk_data]      /Users/sbuciuma/nltk_data...
[nltk_data]   Package stopwords is already up-to-date!

```

```

In [55]: def cleanHtml(html):
         soup = BeautifulSoup(html, "lxml")
         return soup.get_text()

```

```

In [56]: test = clean_df['description'][:25]

```

```

In [57]: blog_posts = []
         for post in test:
             #print(post)
             blog_posts.append({'content': cleanHtml(post)})

```

```

/Library/Frameworks/Python.framework/Versions/3.5/lib/python3.5/site-packages/bs4/k
self.parser.feed(markup)
/Library/Frameworks/Python.framework/Versions/3.5/lib/python3.5/site-packages/bs4/k
self.parser.feed(markup)
/Library/Frameworks/Python.framework/Versions/3.5/lib/python3.5/site-packages/bs4/k
self.parser.feed(markup)
/Library/Frameworks/Python.framework/Versions/3.5/lib/python3.5/site-packages/bs4/k
self.parser.feed(markup)
/Library/Frameworks/Python.framework/Versions/3.5/lib/python3.5/site-packages/bs4/k
self.parser.feed(markup)

```

```

        self.parser.feed(markup)
/Library/Frameworks/Python.framework/Versions/3.5/lib/python3.5/site-packages/bs4/k
        self.parser.feed(markup)
/Library/Frameworks/Python.framework/Versions/3.5/lib/python3.5/site-packages/bs4/k
        self.parser.feed(markup)
/Library/Frameworks/Python.framework/Versions/3.5/lib/python3.5/site-packages/bs4/k
        self.parser.feed(markup)
/Library/Frameworks/Python.framework/Versions/3.5/lib/python3.5/site-packages/bs4/k
        self.parser.feed(markup)
/Library/Frameworks/Python.framework/Versions/3.5/lib/python3.5/site-packages/bs4/k
        self.parser.feed(markup)
/Library/Frameworks/Python.framework/Versions/3.5/lib/python3.5/site-packages/bs4/k
        self.parser.feed(markup)
/Library/Frameworks/Python.framework/Versions/3.5/lib/python3.5/site-packages/bs4/k
        self.parser.feed(markup)
/Library/Frameworks/Python.framework/Versions/3.5/lib/python3.5/site-packages/bs4/k
        self.parser.feed(markup)
/Library/Frameworks/Python.framework/Versions/3.5/lib/python3.5/site-packages/bs4/k
        self.parser.feed(markup)
/Library/Frameworks/Python.framework/Versions/3.5/lib/python3.5/site-packages/bs4/k
        self.parser.feed(markup)
/Library/Frameworks/Python.framework/Versions/3.5/lib/python3.5/site-packages/bs4/k
        self.parser.feed(markup)
/Library/Frameworks/Python.framework/Versions/3.5/lib/python3.5/site-packages/bs4/k
        self.parser.feed(markup)
/Library/Frameworks/Python.framework/Versions/3.5/lib/python3.5/site-packages/bs4/k
        self.parser.feed(markup)
/Library/Frameworks/Python.framework/Versions/3.5/lib/python3.5/site-packages/bs4/k
        self.parser.feed(markup)
/Library/Frameworks/Python.framework/Versions/3.5/lib/python3.5/site-packages/bs4/k
        self.parser.feed(markup)
/Library/Frameworks/Python.framework/Versions/3.5/lib/python3.5/site-packages/bs4/k
        self.parser.feed(markup)
/Library/Frameworks/Python.framework/Versions/3.5/lib/python3.5/site-packages/bs4/k
        self.parser.feed(markup)
/Library/Frameworks/Python.framework/Versions/3.5/lib/python3.5/site-packages/bs4/k
        self.parser.feed(markup)
/Library/Frameworks/Python.framework/Versions/3.5/lib/python3.5/site-packages/bs4/k
        self.parser.feed(markup)

```

```
In [58]: len(blog_posts)
```

```
Out[58]: 25
```

```
In [59]: print (stop_words[:10])
```

```
['i', 'me', 'my', 'myself', 'we', 'our', 'ours', 'ourselves', 'you', 'your']
```

```
In [60]: # load nltk's SnowballStemmer as variable 'stemmer'
```

```
from nltk.stem.snowball import SnowballStemmer
stemmer = SnowballStemmer("english")
```

```
In [61]: # here I define a tokenizer and stemmer which returns the set of stems in
```

```
def tokenize_and_stem(text):
    # first tokenize by sentence, then by word to ensure that punctuation
    tokens = [word for sent in nltk.sent_tokenize(text) for word in nltk.w
    filtered_tokens = []
    # filter out any tokens not containing letters (e.g., numeric tokens,
    for token in tokens:
        if re.search('[a-zA-Z]', token):
            filtered_tokens.append(token)
    stems = [stemmer.stem(t) for t in filtered_tokens]
    return stems
```

```
def tokenize_only(text):
    # first tokenize by sentence, then by word to ensure that punctuation
    tokens = [word.lower() for sent in nltk.sent_tokenize(text) for word i
    filtered_tokens = []
    # filter out any tokens not containing letters (e.g., numeric tokens,
    for token in tokens:
        if re.search('[a-zA-Z]', token):
            filtered_tokens.append(token)
    return filtered_tokens
```

```
In [62]: type(df['description'])
```

```
Out[62]: pandas.core.series.Series
```

```
In [64]: # creating top words per each event and listing as a list, all are appende
```

```
counter = 0
master_list = []

for post in blog_posts:
    sentences = nltk.tokenize.sent_tokenize(post['content']) #will be used
    #print(sentences)
    words = [w.lower() for sentence in sentences for w in # this ones will
              nltk.tokenize.word_tokenize(sentence) if w[0] not in stop_wor
    #print(words)
    master_list.append(words)
    #print(master_list)

fdist = nltk.FreqDist(words) #will be used in k-means for creating the
```



```

# Basic stats

num_words = sum([i[1] for i in fdist.items()])
num_unique_words = len(fdist.keys())

# Hapaxes are words that appear only once

num_hapaxes = len(fdist.hapaxes())

top_10_words_sans_stop_words = [w for w in fdist.items() if w[0]
                                not in stop_words][:25]
total_corpus = [x[0] for x in top_10_words_sans_stop_words]
counter += 1
print ('====|=====|=====|=====|=====|=====|=====')
print("Blog number:", counter)
print ('_____')
print(sentences)
print ('=====')
print("len of each list of words", len(words))
print(words)
print ('.....')

====|=====|=====|=====|=====|=====|=====
Blog number: 1

-----
['This is an introductory Java meetup - I hope we can get together and talk about w
=====
len of each list of words 24
['this', 'java', 'i', 'hope', 'we', 'can', 'get', 'what', 'java', 'we', 'sorry', 'i
.....

====|=====|=====|=====|=====|=====|=====
Blog number: 2

-----
["How about we'll talk about the Web Frameworks people are using?It's just an idea
=====
len of each list of words 27
['how', 'we', 'web', 'frameworks', 'using', 'it', 'just', 'feel', 'free', '!', 'hi
.....

====|=====|=====|=====|=====|=====|=====
Blog number: 3

-----
['November Java Meetup - if nothing else, lets talk about Web frameworks again :)']
=====
len of each list of words 8
['november', 'java', 'meetup', 'nothing', 'else', 'lets', 'web', 'frameworks']
.....

====|=====|=====|=====|=====|=====|=====

```

Blog number: 4

```
['Our monthly gathering to discuss the business of being a Java developer.', "Decem
=====
len of each list of words 19
['our', 'gathering', 'business', 'being', 'java', 'december', 'java', 'meetup', 'le
.....
====|=====|=====|=====|=====|=====|=====
Blog number: 5
```

```
['Happy New Year!', 'We invite anybody to share an interesting Java application he
=====
len of each list of words 8
['happy', 'new', 'year', '!', 'we', 'java', 'he', 'worked']
.....
====|=====|=====|=====|=====|=====|=====
Blog number: 6
```

```
['Meet Jimbo after his conference for early dinner, 4-7pm.', "Jimmy Wales will be i
=====
len of each list of words 45
['meet', 'jimbo', 'his', 'conference', 'for', 'early', '4-7pm', 'jimmy', 'wales', '
.....
====|=====|=====|=====|=====|=====|=====
Blog number: 7
```

```
['Discuss V-Pedia, sweethearts, sweetmeats, fine bookshops, good mid-town ski-jumps
=====
len of each list of words 24
['discuss', 'v-pedia', 'fine', 'bookshops', 'good', 'free', 'conferences', 'essenti
.....
====|=====|=====|=====|=====|=====|=====
Blog number: 8
```

```
["It's time to have a real Spring party.", 'The Boston Cyberarts Festival is coming
=====
len of each list of words 31
['it', 'have', 'real', 'spring', 'the', 'boston', 'cyberarts', 'festival', 'coming
.....
====|=====|=====|=====|=====|=====|=====
Blog number: 9
```

```
['Yahoo!ligans unite... come have fun in Club Passim for an hour and talk about wor
=====
len of each list of words 35
['yahoo', '!', 'ligans', 'unite', 'come', 'have', 'fun', 'club', 'passim', 'for', '
.....
====|=====|=====|=====|=====|=====|=====
```

Blog number: 10

```
["Dear Java Enthusiasts,I wanted to throw out a little agenda for our April Meetup.
=====
len of each list of words 36
['dear', 'java', 'enthusiasts', 'i', 'wanted', 'little', 'for', 'april', 'meetup',
.....
====|=====|=====|=====|=====|=====|=====
Blog number: 11
```

```
['Jimmy Wales will be in town next Tuesday, which is fortuitously about the right t
=====
len of each list of words 36
['jimmy', 'wales', 'will', 'be', 'next', 'tuesday', 'which', 'fortuitously', 'right
.....
====|=====|=====|=====|=====|=====|=====
Blog number: 12
```

```
['Hello all Java enthusiasts!', 'As we promised, our next meetup will be different!
=====
len of each list of words 60
['hello', 'java', 'enthusiasts', '!', 'as', 'we', 'next', 'will', 'be', '!', 'we',
.....
====|=====|=====|=====|=====|=====|=====
Blog number: 13
```

```
['Planning for upcoming speakers and events.Library and community outreach; local m
=====
len of each list of words 46
['planning', 'for', 'upcoming', 'events.library', 'community', ';', 'local', 'educa
.....
====|=====|=====|=====|=====|=====|=====
Blog number: 14
```

```
['Meet inside Pho Pasteur.', 'Bring a notebook and a bib.', 'For more information,
=====
len of each list of words 12
['meet', 'pho', 'pasteur', 'bring', 'notebook', 'bib', 'for', 'http', '//boston.cit
.....
====|=====|=====|=====|=====|=====|=====
Blog number: 15
```

```
['There will be a speaker this evening; so come prepared for a 20-minute presentati
=====
len of each list of words 30
['there', 'will', 'be', 'evening', ';', 'come', 'for', '20-minute', 'followed', 'by
.....
====|=====|=====|=====|=====|=====|=====
```

Blog number: 16

```
['Hello, Java Meetup members!', 'At our next Meetup, David Thomson, creator of an o
=====
len of each list of words 60
['hello', 'java', 'meetup', '!', 'at', 'next', 'meetup', 'david', 'thomson', 'creat
.....
====|=====|=====|=====|=====|=====|=====
Blog number: 17
```

```
['There will likely be a speaker this evening, on usability; so come prepared for a
=====
len of each list of words 34
['there', 'will', 'likely', 'be', 'evening', 'usability', ';', 'come', 'for', '20-m
.....
====|=====|=====|=====|=====|=====|=====
Blog number: 18
```

```
['j here: SJ has a conflict at about 7 pm.', "Depending on when his meeting ends, w
=====
len of each list of words 40
['j', 'here', 'sj', 'has', 'conflict', '7', 'depending', 'when', 'his', 'ends', 'we
.....
====|=====|=====|=====|=====|=====|=====
Blog number: 19
```

```
['Hello, Java Meetup members!this month, Renaud Richardet from Wyona, agreed to int
=====
len of each list of words 60
['hello', 'java', 'meetup', '!', 'renaud', 'richardet', 'from', 'wyona', 'apache-co
.....
====|=====|=====|=====|=====|=====|=====
Blog number: 20
```

```
['Discussion of great local ice-cream joints.', 'Exploration and patronization of s
=====
len of each list of words 15
['discussion', 'great', 'local', 'joints', 'exploration', 'establishments', 'oh', '
.....
====|=====|=====|=====|=====|=====|=====
Blog number: 21
```

```
['Discussion of great local ice-cream joints.', 'Exploration and patronization of s
=====
len of each list of words 22
['discussion', 'great', 'local', 'joints', 'exploration', 'establishments', 'oh', '
.....
====|=====|=====|=====|=====|=====|=====
```

Blog number: 22

```
['Hello, Java Meetup members!', 'This month we want to gather in a social setting a
=====
len of each list of words 55
['hello', 'java', 'meetup', '!', 'this', 'we', 'want', 'gather', 'have', 'questions
.....
====|=====|=====|=====|=====|=====|=====
Blog number: 23
```

```
['Wikimania stuff and just general wiki and Wikipedia stuff and /fabulous!!', '!/ i
=====
len of each list of words 20
['wikimania', 'just', 'general', 'wiki', 'wikipedia', '/fabulous', '!', '!', '!', '
.....
====|=====|=====|=====|=====|=====|=====
Blog number: 24
```

```
['Hello, Java Meetup members!lets gather and talk Java.', 'There is so much news ar
=====
len of each list of words 32
['hello', 'java', 'meetup', '!', 'lets', 'gather', 'java', 'there', 'news', 'java',
.....
====|=====|=====|=====|=====|=====|=====
Blog number: 25
```

```
['Wikimania stuff and just general wiki and Wikipedia stuff and /fabulous!!', '!/ i
=====
len of each list of words 19
['wikimania', 'just', 'general', 'wiki', 'wikipedia', '/fabulous', '!', '!', '!', '
.....
```

```
In [66]: #writing the lists to local file "new_filename.txt" as each list is a line
         ##writelines in python3 to iterate thru lists
```

```
with open('new_filename.txt', 'w') as f:
    f.writelines("%s\n" % l for l in master_list)

#####words from each blog are selected,
#####has been scheduled 60 words what represent the content
```

```
In [67]: import numpy
         from nltk.cluster import KMeansClusterer, GAAClusterer, cosine_distance
         import nltk.corpus
         import nltk.stem
         stemmer_func = nltk.stem.snowball.SnowballStemmer("english").stem
         stopwords = set(nltk.corpus.stopwords.words('english'))
```



```

0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0,
0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0,
0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0,
0, 1, 0, 0, 0, 0, 0, 0], dtype=int16)]

```

```

In [74]: for title in job_titles:
         print(title)

```

```

['this', 'java', 'i', 'hope', 'we', 'can', 'get', 'what', 'java', 'we', 'sorry', 'i',
['how', 'we', 'web', 'frameworks', 'using', 'it', 'just', 'feel', 'free', '!', 'hi',
['november', 'java', 'meetup', 'nothing', 'else', 'lets', 'web', 'frameworks']
['our', 'gathering', 'business', 'being', 'java', 'december', 'java', 'meetup', 'le',
['happy', 'new', 'year', '!', 'we', 'java', 'he', 'worked']
['meet', 'jimbo', 'his', 'conference', 'for', 'early', '4-7pm', 'jimmy', 'wales', 'v',
['discuss', 'v-pedia', 'fine', 'bookshops', 'good', 'free', 'conferences', 'essenti',
['it', 'have', 'real', 'spring', 'the', 'boston', 'cyberarts', 'festival', 'coming',
['yahoo', '!', 'ligans', 'unite', 'come', 'have', 'fun', 'club', 'passim', 'for', 'v',
['dear', 'java', 'enthusiasts', 'i', 'wanted', 'little', 'for', 'april', 'meetup',
['jimmy', 'wales', 'will', 'be', 'next', 'tuesday', 'which', 'fortuitously', 'right',
['hello', 'java', 'enthusiasts', '!', 'as', 'we', 'next', 'will', 'be', '!', 'we',
['planning', 'for', 'upcoming', 'events.library', 'community', ';', 'local', 'educa',
['meet', 'pho', 'pasteur', 'bring', 'notebook', 'bib', 'for', 'http', '//boston.cit',
['there', 'will', 'be', 'evening', ';', 'come', 'for', '20-minute', 'followed', 'by',
['hello', 'java', 'meetup', '!', 'at', 'next', 'meetup', 'david', 'thomson', 'creat',
['there', 'will', 'likely', 'be', 'evening', 'usability', ';', 'come', 'for', '20-m',
['j', 'here', 'sj', 'has', 'conflict', '7', 'depending', 'when', 'his', 'ends', 'we',
['hello', 'java', 'meetup', '!', 'renaud', 'richardet', 'from', 'wyona', 'apache-co',
['discussion', 'great', 'local', 'joints', 'exploration', 'establishments', 'oh', 'v',
['discussion', 'great', 'local', 'joints', 'exploration', 'establishments', 'oh', 'v',
['hello', 'java', 'meetup', '!', 'this', 'we', 'want', 'gather', 'have', 'questions',
['wikimania', 'just', 'general', 'wiki', 'wikipedia', '/fabulous', '!', '!', '!', 'v',
['hello', 'java', 'meetup', '!', 'lets', 'gather', 'java', 'there', 'news', 'java', 'v',
['wikimania', 'just', 'general', 'wiki', 'wikipedia', '/fabulous', '!', '!', '!', 'v',

```

```

In [75]: cluster = KMeansClusterer(10, cosine_distance) ###cosine distance from nlt
         cluster.cluster([vectorspaced(title) for title in job_titles if title]) #v
         classified_examples = [cluster.classify(vectorspaced(title)) for title in

```

```

In [76]: for cluster_id, title in sorted(zip(classified_examples, job_titles)):
         print ("\n", "Cluster number:", cluster_id, "\n", "Words extracted from

```

Cluster number: 0

Words extracted from the event representing individual event by description:

```

['meet', 'pho', 'pasteur', 'bring', 'notebook', 'bib', 'for', 'http', '//boston.ci

```

Cluster number: 1

Words extracted from the event representing individual event by description:

['november', 'java', 'meetup', 'nothing', 'else', 'lets', 'web', 'frameworks']

Cluster number: 1

Words extracted from the event representing individual event by description:

['our', 'gathering', 'business', 'being', 'java', 'december', 'java', 'meetup', 'l

Cluster number: 2

Words extracted from the event representing individual event by description:

['how', 'we', 'web', 'frameworks', 'using', 'it', 'just', 'feel', 'free', '!', 'hi

Cluster number: 2

Words extracted from the event representing individual event by description:

['there', 'will', 'be', 'evening', ';', 'come', 'for', '20-minute', 'followed', 'k

Cluster number: 2

Words extracted from the event representing individual event by description:

['there', 'will', 'likely', 'be', 'evening', 'usability', ';', 'come', 'for', '20-

Cluster number: 2

Words extracted from the event representing individual event by description:

['this', 'java', 'i', 'hope', 'we', 'can', 'get', 'what', 'java', 'we', 'sorry', '

Cluster number: 3

Words extracted from the event representing individual event by description:

['dear', 'java', 'enthusiasts', 'i', 'wanted', 'little', 'for', 'april', 'meetup',

Cluster number: 3

Words extracted from the event representing individual event by description:

['hello', 'java', 'meetup', '!', 'lets', 'gather', 'java', 'there', 'news', 'java

Cluster number: 3

Words extracted from the event representing individual event by description:

['hello', 'java', 'meetup', '!', 'this', 'we', 'want', 'gather', 'have', 'question

Cluster number: 4

Words extracted from the event representing individual event by description:

['it', 'have', 'real', 'spring', 'the', 'boston', 'cyberarts', 'festival', 'coming

Cluster number: 4

Words extracted from the event representing individual event by description:

['j', 'here', 'sj', 'has', 'conflict', '7', 'depending', 'when', 'his', 'ends', 'w

Cluster number: 4

Words extracted from the event representing individual event by description:

['planning', 'for', 'upcoming', 'events.library', 'community', ';', 'local', 'educ

Cluster number: 4

Words extracted from the event representing individual event by description:

['wikimania', 'just', 'general', 'wiki', 'wikipedia', '/fabulous', '!', '!', '!',

Cluster number: 4

Words extracted from the event representing individual event by description:

['wikimania', 'just', 'general', 'wiki', 'wikipedia', '/fabulous', '!', '!', '!',

Cluster number: 5

Words extracted from the event representing individual event by description:

['hello', 'java', 'enthusiasts', '!', 'as', 'we', 'next', 'will', 'be', '!', 'we',

Cluster number: 6

Words extracted from the event representing individual event by description:

['discuss', 'v-pedia', 'fine', 'bookshops', 'good', 'free', 'conferences', 'essent

Cluster number: 6

Words extracted from the event representing individual event by description:

['jimmy', 'wales', 'will', 'be', 'next', 'tuesday', 'which', 'fortuitously', 'right

Cluster number: 6

Words extracted from the event representing individual event by description:

['meet', 'jimbo', 'his', 'conference', 'for', 'early', '4-7pm', 'jimmy', 'wales',

Cluster number: 6

Words extracted from the event representing individual event by description:

['yahoo', '!', 'ligans', 'unite', 'come', 'have', 'fun', 'club', 'passim', 'for',

Cluster number: 7

Words extracted from the event representing individual event by description:

```
['discussion', 'great', 'local', 'joints', 'exploration', 'establishments', 'oh',
```

Cluster number: 7

Words extracted from the event representing individual event by description:

```
['discussion', 'great', 'local', 'joints', 'exploration', 'establishments', 'oh',
```

Cluster number: 8

Words extracted from the event representing individual event by description:

```
['hello', 'java', 'meetup', '!', 'at', 'next', 'meetup', 'david', 'thomson', 'crea
```

Cluster number: 8

Words extracted from the event representing individual event by description:

```
['hello', 'java', 'meetup', '!', 'renaud', 'richardet', 'from', 'wyona', 'apache-c
```

Cluster number: 9

Words extracted from the event representing individual event by description:

```
['happy', 'new', 'year', '!', 'we', 'java', 'he', 'worked']
```

In []: