

In [1]:

```
1 '''Imports'''
2 from IPython.display import Audio, display, clear_output
3 from preamble import *
4
5 cwd = os.getcwd()
6 dataPath = os.path.join(cwd, 'data')
7 savePath = os.path.join(cwd, 'saves')
8
9 print('Current Working Directory:\n', cwd)
10 print('\nData Directory:\n', dataPath)
11 print('\nSaves Direcory:\n', savePath)
```

Current Working Directory:

C:\Users\josh.galloway\Google Drive\Sync\Education_Northeastern University\03_DS 5230 Unsup ML\99_Project\05_Code

Data Directory:

C:\Users\josh.galloway\Google Drive\Sync\Education_Northeastern University\03_DS 5230 Unsup ML\99_Project\05_Code\data

Saves Direcory:

C:\Users\josh.galloway\Google Drive\Sync\Education_Northeastern University\03_DS 5230 Unsup ML\99_Project\05_Code\saves

Show Package Versions and Code to Download various Lexicons

Leave commented out unless needed.

In [2]:

```
1 # import sys
2 # import nltk
3 # import sklearn
4 # import matplotlib
5 # import scipy
6 # import wordcloud
7
8 # print('Python version {}'.format(sys.version))
9 # print('NumPy version {}'.format(np.__version__))
10 # print('Matplotlib version {}'.format(matplotlib.__version__))
11 # print('Seaborn version {}'.format(sns.__version__))
12 # print('Pandas version {}'.format(pd.__version__))
13 # print('Scipy version {}'.format(scipy.__version__))
14 # print('nltk version {}'.format(nltk.__version__))
15 # print('scikit-Learn version {}'.format(sklearn.__version__))
16 # print('wordcloud version {}'.format(wordcloud.__version__))
17
18
19 # # Get List of Stop Words
20 # nltk.download('stopwords')
21 # # Get Tokenizer
22 # nltk.download('punkt')
23 # # Get Wordnet Lemmatizer
24 # nltk.download('wordnet')
25 # nltk.download('vader_Lexicon')
26
27 versions = ['Python version 3.7.3\n\t(default, Mar 27 2019, 17:13:21)\n\t[MSC v.1915 64 bit (AMD64)]',
28     'NumPy version 1.16.2',
29     'Matplotlib version 3.0.3',
30     'Seaborn version 0.10.1',
31     'Pandas version 0.24.2',
32     'Scipy version 1.4.1',
33     'nltk version 3.5',
34     'scikit-learn version 0.23.1',
35     'wordcloud version 1.7.0']
36 print('*'*10 + ' Package Versions ' + '*'*10)
37 for v in versions:
38     print(v)
```

===== Package Versions =====

```
Python version 3.7.3
    (default, Mar 27 2019, 17:13:21)
    [MSC v.1915 64 bit (AMD64)]
NumPy version 1.16.2
Matplotlib version 3.0.3
Seaborn version 0.10.1
Pandas version 0.24.2
Scipy version 1.4.1
nltk version 3.5
scikit-learn version 0.23.1
wordcloud version 1.7.0
```

```

In [3]: 1  '''Save or recall objects'''
2  def pickle_object(obj, filename, save=True, addTime=False):
3      if save:
4          # save the results to disk
5          if addTime:
6              temp = filename.split('.')
7              filename = temp[0] + datetime.now().strftime("%Y%m%d-%H%M%S")+'.' + temp[-1]
8              pickle.dump(obj, open(filename, 'wb'))
9              print("Saved File: ",filename)
10             print('At Time: ',datetime.now().strftime("%m/%d/%Y, %H:%M:%S"))
11             return obj
12         else:
13             return pickle.load(open(filename, 'rb'))
14
15  '''Display Topics From LDA and NMF Models'''
16  def display_topics(model, feature_names, no_top_words):
17      for i in range(len(model.components_[:,0])):
18          print("Topic {:d} ======".format(i))
19          idx = (-model.components_[i,:]).argsort()[:no_top_words]
20          print("|".join([feature_names[k] for k in idx]))
21
22  def display_centroids(model, feature_names, no_top_words):
23      vs = vader.SentimentIntensityAnalyzer()
24      for i in range(len(model.cluster_centers_[:,0])):
25          print("Centroid {:d} ======".format(i))
26          idx = (-model.cluster_centers_[i,:]).argsort()[:no_top_words]
27          print("|".join([feature_names[k] for k in idx]))
28          print('Sentiment: ',
29                  vs.polarity_scores(" ".join([feature_names[k] for k in idx])))
30
31 beepFile = os.path.join(cwd,'data','Pager Beeps.mp3')
32 def done():
33     '''Beep to get attention at end of code execution'''
34     display(Video(beepFile, autoplay=True))
35
36
37 def plot_dendrogram(data, title, linkage_='ward', model = None):
38     '''Build Hierarchical Model and Plot Dendrogram'''
39     # Create Linkage matrix and then plot the dendrogram
40     if model is None:
41         model = AgglomerativeClustering(distance_threshold=0,
42                                         n_clusters=None, linkage=linkage_)
43     model = model.fit(data)
44
45     # create the counts of samples under each node
46     counts = np.zeros(model.children_.shape[0])
47     n_samples = len(model.labels_)
48     for i, merge in enumerate(model.children_):
49         current_count = 0
50         for child_idx in merge:
51             if child_idx < n_samples:
52                 current_count += 1 # Leaf node
53             else:
54                 current_count += counts[child_idx - n_samples]
55         counts[i] = current_count
56
57     linkage_matrix = np.column_stack([model.children_, model.distances_,
58                                     counts]).astype(float)
59
60     # Plot the corresponding dendrogram
61     plt.figure(figsize=(15,7),facecolor='#d3d3d3')
62     plt.title(title,fontweight='bold')
63     dendrogram(linkage_matrix)
64     plt.yticks([])
65     _=plt.xlabel("Number of points in node (or index of point if no parenthesis).")
66     return model
67
68 def plot_clustering(cluster_centers_, data_df, title, angles, fsize=(15,15)):
69     '''Plot Clustering with DataFrame columns neg, neu, pos, Label'''
70     fig = plt.figure(tight_layout = True, figsize=fsize, facecolor='#d3d3d3')
71     marks = ['o','^','s','*','x']

```

```
72     lgd = 'Centroid {:d}'
73     #angles = [(None,30),(10,90),(100,0),(10,0)]
74     for j,angle in enumerate(angles):
75         ax = fig.add_subplot(2,2,j+1, projection='3d')
76
77         # build dict for labels and colors
78         lbls = {}
79         for i,lbl in enumerate(set(data_df['Label'])):
80             lbls[lbl] = i
81
82         # Plot points
83         for i, topic in enumerate(data_df[['neg','neu','pos']].values):
84             cl = lbls[data_df['Label'].loc[i]]
85             ax.scatter(topic[0], topic[1], topic[2],
86                         c=COLORS[cl], marker=marks[cl], s=25)
87
88         # Plot Centers or Ad
89         if cluster_centers_ is None:
90             # build custom legend
91             le = []
92             for key in lbls.keys():
93                 cl = lbls[key]
94                 le.append(Line2D([0], [0], marker=marks[cl], color=COLORS[cl],
95                                 label=lgd.format(key), markerfacecolor=COLORS[cl],
96                                 markersize=10))
97             ax.legend(handles = le, loc='best')
98         else:
99             # Plot Centers as bigger markers
100             for i, center in enumerate(cluster_centers_):
101                 ax.scatter(center[0], center[1], center[2],
102                             c=COLORS[i], marker=marks[i], s=150,
103                             label = lgd.format(i))
104             ax.legend(loc='best')
105
106             ax.set_xlabel('Negative Sentiment',fontweight='bold')
107             ax.set_ylabel('Neutral Sentiment',fontweight='bold')
108             ax.set_zlabel('Positive Sentiment',fontweight='bold')
109             ax.set_title(title,fontweight='bold')
110             ax.view_init(*angle)
```

Load Datasets

```
In [4]: 1 abcFile = 'abcnews-date-text.csv'
2 redditFile = 'RedditNews.csv'
3 abcPath = os.path.join(dataPath,abcFile)
4 redditPath = os.path.join(dataPath,redditFile)
5
6 abcDF = pd.read_csv(str(abcPath))
7 redditDF = pd.read_csv(str(redditPath))
8 display(abcDF.head())
9 display(redditDF.head())
```

	publish_date	headline_text
0	20030219	aba decides against community broadcasting lic...
1	20030219	act fire witnesses must be aware of defamation
2	20030219	a g calls for infrastructure protection summit
3	20030219	air nz staff in aust strike for pay rise
4	20030219	air nz strike to affect australian travellers

	Date	News
0	2016-07-01	A 117-year-old woman in Mexico City finally re...
1	2016-07-01	IMF chief backs Athens as permanent Olympic host
2	2016-07-01	The president of France says if Brexit won, so...
3	2016-07-01	British Man Who Must Give Police 24 Hours' Not...
4	2016-07-01	100+ Nobel laureates urge Greenpeace to stop o...

Clean and Synthesize into one Dataframe

```
In [5]: 1 # Align Column Names
2 abcDF.columns = redditDF.columns
3 print('ABC Dataframe')
4 display(abcDF.info())
5 print('Reddit Dataframe')
6 display(redditDF.info())
```

ABC Dataframe
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1186018 entries, 0 to 1186017
Data columns (total 2 columns):
Date 1186018 non-null int64
News 1186018 non-null object
dtypes: int64(1), object(1)
memory usage: 18.1+ MB

None

Reddit Dataframe
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 73608 entries, 0 to 73607
Data columns (total 2 columns):
Date 73608 non-null object
News 73608 non-null object
dtypes: object(2)
memory usage: 1.1+ MB

None

Add Features to capture information that may be eliminated otherwise

Since we are about to alter the headlines from the original form and then combine them into one dataframe, add features to try to capture that data. Similarly, as part of the cleaning process digits will be removed, so mark the headlines that have numbers.

In [6]:

```
1 # Add new features
2 hasNum = lambda x: int(bool(re.search(r'\d', x)) == True)
3
4 abcDF['Source'] = 'ABC' # Set Source
5 abcDF['Length'] = abcDF.News.apply(len) # Save original headline length
6 abcDF['HasNumbers'] = abcDF.News.apply(hasNum)
7
8 redditDF['Source'] = 'Reddit' # Set Source
9 redditDF['Length'] = redditDF.News.apply(len) # Save original headline length
10 redditDF['HasNumbers'] = redditDF.News.apply(hasNum)
11
12 print('ABC Dataframe')
13 display(abcDF.info())
14 print('Reddit Dataframe')
15 display(redditDF.info())
```

```
ABC Dataframe
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1186018 entries, 0 to 1186017
Data columns (total 5 columns):
Date      1186018 non-null int64
News      1186018 non-null object
Source     1186018 non-null object
Length    1186018 non-null int64
HasNumbers 1186018 non-null int64
dtypes: int64(3), object(2)
memory usage: 45.2+ MB
```

```
None
```

```
Reddit Dataframe
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 73608 entries, 0 to 73607
Data columns (total 5 columns):
Date      73608 non-null object
News      73608 non-null object
Source     73608 non-null object
Length    73608 non-null int64
HasNumbers 73608 non-null int64
dtypes: int64(2), object(3)
memory usage: 2.8+ MB
```

```
None
```

Convert the Dates to Datetime objects

In [7]:

```
1 # Convert to Date Time 'Date' Column
2 abcDF.Date = pd.to_datetime(abcDF.Date,format='%Y%m%d')
3 redditDF.Date = pd.to_datetime(redditDF.Date,format='%Y-%m-%d')
4 print('ABC Dataframe')
5 display(abcDF.info())
6 print('Reddit Dataframe')
7 display(redditDF.info())
```

ABC Dataframe
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1186018 entries, 0 to 1186017
Data columns (total 5 columns):
Date 1186018 non-null datetime64[ns]
News 1186018 non-null object
Source 1186018 non-null object
Length 1186018 non-null int64
HasNumbers 1186018 non-null int64
dtypes: datetime64[ns](1), int64(2), object(2)
memory usage: 45.2+ MB

None

Reddit Dataframe
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 73608 entries, 0 to 73607
Data columns (total 5 columns):
Date 73608 non-null datetime64[ns]
News 73608 non-null object
Source 73608 non-null object
Length 73608 non-null int64
HasNumbers 73608 non-null int64
dtypes: datetime64[ns](1), int64(2), object(2)
memory usage: 2.8+ MB

None

Combine into one Dataframe

In [8]:

```

1 # Combine in to one dataframe
2 totalDF = pd.concat([abcDF,redditDF], ignore_index=True)
3 display(totalDF.head())
4 display(totalDF.tail())
5 display(totalDF.info())
6 display(totalDF.describe())

```

	Date		News	Source	Length	HasNumbers
0	2003-02-19	aba decides against community broadcasting lic...	ABC	50	0	
1	2003-02-19	act fire witnesses must be aware of defamation	ABC	46	0	
2	2003-02-19	a g calls for infrastructure protection summit	ABC	46	0	
3	2003-02-19	air nz staff in aust strike for pay rise	ABC	40	0	
4	2003-02-19	air nz strike to affect australian travellers	ABC	45	0	

	Date		News	Source	Length	HasNumbers
1259621	2008-06-08	b'Man goes berzerk in Akihabara and stabs ever...	Reddit	79	1	
1259622	2008-06-08	b'Threat of world AIDS pandemic among heterose...	Reddit	75	0	
1259623	2008-06-08	b'Angst in Ankara: Turkey Steers into a Danger...	Reddit	66	0	
1259624	2008-06-08	b"UK: Identity cards 'could be used to spy on ...	Reddit	171	0	
1259625	2008-06-08	b'Marriage, they said, was reduced to the stat...	Reddit	175	0	

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1259626 entries, 0 to 1259625
Data columns (total 5 columns):
Date          1259626 non-null datetime64[ns]
News          1259626 non-null object
Source        1259626 non-null object
Length        1259626 non-null int64
HasNumbers    1259626 non-null int64
dtypes: datetime64[ns](1), int64(2), object(2)
memory usage: 48.1+ MB

```

None

	Length	HasNumbers
count	1.259626e+06	1.259626e+06
mean	4.478764e+01	8.092402e-02
std	2.550213e+01	2.727185e-01
min	3.000000e+00	0.000000e+00
25%	3.400000e+01	0.000000e+00
50%	4.200000e+01	0.000000e+00
75%	4.900000e+01	0.000000e+00
max	3.180000e+02	1.000000e+00

Clean the Headlines and Create New Column for the Altered Data

In [9]:

```

1 # Clean Dataframe
2
3 # All to lower case
4 t0 = time()
5 print('\nSetting all to lowercase...')
6 totalDF['Cleaned'] = totalDF.News.apply(str.lower)
7 print("done in %0.3fs." % (time() - t0))
8
9 # Remove Punctuation
10 remPunct = lambda x: x.translate(str.maketrans('', '', punctuation))
11 t0 = time()
12 print('\nRemoving punctuation...')
13 totalDF.Cleaned = totalDF.Cleaned.apply(remPunct)
14 print("done in %0.3fs." % (time() - t0))
15
16 # Remove Numbers
17 remNum = lambda x: x.translate(str.maketrans('', '', digits))
18 t0 = time()
19 print('\nRemoving numbers...')
20 totalDF.Cleaned = totalDF.Cleaned.apply(remNum)
21 print("done in %0.3fs." % (time() - t0))
22
23 # Get List of Stop words with punctuation removed
24 nopunctStopwords = [w.translate(str.maketrans('', '', punctuation)) for w in stopwords.words('english')]
25 remStop = lambda x: ' '.join([w for w in x.split() if w not in nopunctStopwords])
26 t0 = time()
27 print('\nRemoving stopwords...')
28 totalDF.Cleaned = totalDF.Cleaned.apply(remStop)
29 print("done in %0.3fs." % (time() - t0))
30
31 # Lemmatize or Stem Words
32 wnl = WordNetLemmatizer()
33 ps = PorterStemmer()
34 lemma = lambda x: ' '.join([wnl.lemmatize(t) for t in word_tokenize(x)])
35 stemm = lambda x: ' '.join([ps.stem(t) for t in word_tokenize(x)])
36 t0 = time()
37 print('\nLemmatizing/Stemming words...')
38 totalDF.Cleaned = totalDF.Cleaned.apply(lemma)
39 print("Completed in %0.3fs." % (time() - t0))
40
41 display(totalDF.head())
42 done()

```

Setting all to lowercase...
done in 0.307s.

Removing punctuation...
done in 3.657s.

Removing numbers...
done in 2.491s.

Removing stopwords...
done in 15.565s.

Lemmatizing/Stemming words...
Completed in 129.570s.

	Date	News	Source	Length	HasNumbers	Cleaned
0	2003-02-19	aba decides against community broadcasting lic...	ABC	50	0	aba decides community broadcasting licence
1	2003-02-19	act fire witnesses must be aware of defamation	ABC	46	0	act fire witness must aware defamation
2	2003-02-19	a g calls for infrastructure protection summit	ABC	46	0	g call infrastructure protection summit

	Date	News	Source	Length	HasNumbers	Cleaned
3	2003-02-19	air nz staff in aust strike for pay rise	ABC	40	0	air nz staff aust strike pay rise
4	2003-02-19	air nz strike to affect australian travellers	ABC	45	0	air nz strike affect australian traveller

0:04 / 0:04

Save off Dataframe for Later Recall

In [10]:

```

1 # # Save off the Dataframe object
2 # pickleFile = os.path.join(savePath, 'cleanedData.sav')
3 # csvFile = os.path.join(savePath, 'cleanedData.csv')
4
5 # pickle_object(totalDF, str(pickleFile), save=True, addTime=False)
6 # totalDF.to_csv(str(csvFile), index=False)

```

Create term frequency-inverse document frequency and term frequency matrices for LSI, PLSI and LDA

Set tokenizer to use Unigrams and Bigrams

In [11]:

```

1 no_features = 2000 # Number of words to truncate vocabulary, chosen arbitrarily
2
3 # SVD and NMF are able to use tf-idf
4 t0 = time()
5 print('\nCreating TF-IDF Matrix...')
6 tfidf_vectorizer = TfidfVectorizer(max_features=no_features,
7                                     lowercase=False, ngram_range=(1,2))
8 tfidf = tfidf_vectorizer.fit_transform(totalDF.Cleaned)
9 tfidf_feature_names = tfidf_vectorizer.get_feature_names()
10 print("Completed in %0.3fs." % (time() - t0))
11
12
13 # LDA can only use raw term counts for LDA because it is a probabilistic graphical model
14 t0 = time()
15 print('\nCreating TF Matrix...')
16 tf_vectorizer = CountVectorizer(max_features=no_features,
17                                 lowercase=False, ngram_range=(1,2))
18 tf = tf_vectorizer.fit_transform(totalDF.Cleaned)
19 tf_feature_names = tf_vectorizer.get_feature_names()
20 print("Completed in %0.3fs." % (time() - t0))
21 done()

```

Creating TF-IDF Matrix...
Completed in 20.257s.

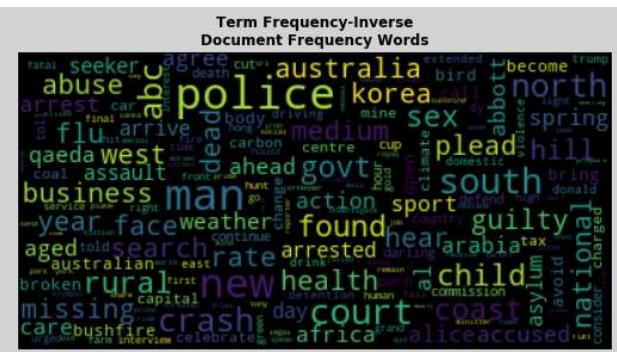
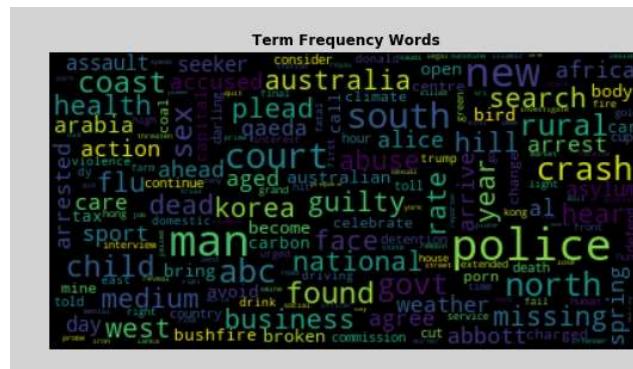
Creating TF Matrix...
Completed in 20.781s.

0:04 / 0:04

Visualize the terms and find average length of headline

```
In [87]: 1 # Find Average Length of word in feature set
2 charLenAvg = round(totalDF.Length.mean())
3 avg = 0
4 for w in tfidf_feature_names:
5     avg = len(w) + avg
6 print('TF-TDF Matrix Size: ',tfidf.shape)
7 print('TF-IDF Feature Set Average Word Length: ',round(avg/len(tfidf_feature_names)))
8 print('TF-IDF Feature Set Average Headline Words: ',charLenAvg/round(avg/len(tfidf_feature_names)))
9
10 avg = 0
11 for w in tf_feature_names:
12     avg = len(w) + avg
13
14 print('TF Matrix Size: ',tf.shape)
15 print('TF Feature Set Average Word Length: ',round(avg/len(tf_feature_names)))
16 print('TF Feature Set Average Headline Words: ',charLenAvg/round(avg/len(tf_feature_names)))
17
18
19 # plot wordcloud
20 plt.figure(tight_layout = True, figsize = (15, 8), facecolor = '#d3d3d3')
21 titles = ['Term Frequency Words', 'Term Frequency-Inverse\nDocument Frequency Words']
22 wc = WordCloud(background_color ='black', min_font_size = 5,max_font_size=30)
23 for i,terms in enumerate([tf_feature_names, tfidf_feature_names]):
24     # plot the WordCloud image
25     wc.generate(' '.join(terms))
26     plt.subplot(1,2,i+1)
27     plt.imshow(wc, interpolation='bilinear')
28     plt.axis("off")
29     plt.title(titles[i],fontweight='bold')
30
```

TF-TDF Matrix Size: (1259626, 2000)
TF-IDF Feature Set Average Word Length: 6
TF-IDF Feature Set Average Headline Words: 7.5
TF Matrix Size: (1259626, 2000)
TF Feature Set Average Word Length: 6
TF Feature Set Average Headline Words: 7.5



Save TF and TF-IDF Sparse Matrices for Later Recall

```
In [13]: 1 # # Save TF Matix
2 # pickle_object(tf_feature_names,str(os.path.join(savePath,'TF_Features.sav')),
3 #                 save=True,addTime=False)
4 # pickle_object(tf,str(os.path.join(savePath,'TFMatix.sav')),
5 #                 save=True,addTime=False)
6
7 # # Save TF-IDF Matrix
8 # pickle_object(tfidf_feature_names,str(os.path.join(savePath,'TFIDF_Features.sav')),
9 #                 save=True,addTime=False)
10 # pickle_object(tfidf,str(os.path.join(savePath,'TFIDFMatrix.sav')),
11 #                  save=True,addTime=False)
```

Create Test Train Split Stratified By Source

```
In [14]: 1 np.random.seed(42)
2 idx_train, idx_test = train_test_split(totalDF.index.values,
3                                         stratify=totalDF.Source,
4                                         test_size=0.25)
5 print('Original Dataset Sources: ')
6 cnts = totalDF.Source.value_counts()
7 print(cnts)
8 print('ABC: ',cnts['ABC']/len(totalDF.index)*100,'Reddit: ',cnts['Reddit']/len(totalDF.index)*100)
9 print('\nTraining Set Sources: ')
10 cnts = totalDF.loc[idx_train].Source.value_counts()
11 print(cnts)
12 print('ABC: ',cnts['ABC']/len(idx_train)*100,'Reddit: ',cnts['Reddit']/len(idx_train)*100)
13 print('\nTest Set Sources: ')
14 cnts = totalDF.loc[idx_test].Source.value_counts()
15 print(cnts)
16 print('ABC: ',cnts['ABC']/len(idx_test)*100,'Reddit: ',cnts['Reddit']/len(idx_test)*100)
17
```

Original Dataset Sources:

Source	Count
ABC	1186018
Reddit	73608

Name: Source, dtype: int64

Source	Count		
ABC	94.15636069754038	Reddit	5.843639302459619
Reddit	5.843639302459619		

Training Set Sources:

Source	Count
ABC	889513
Reddit	55206

Name: Source, dtype: int64

Source	Count		
ABC	94.15635760474808	Reddit	5.843642395251922
Reddit	5.843642395251922		

Test Set Sources:

Source	Count
ABC	296505
Reddit	18402

Name: Source, dtype: int64

Source	Count		
ABC	94.15636997589765	Reddit	5.843630024102354
Reddit	5.843630024102354		

Build Latent Semantic Analysis to find Optimal Number of Topics

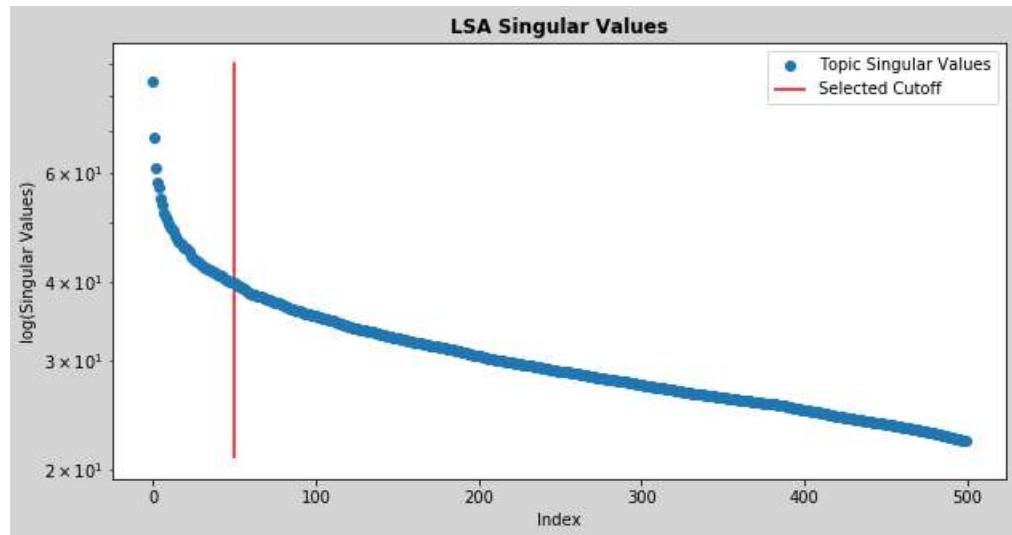
```
In [15]: 1 # Build a Latent Semantic Indexing Model using SVD
2 t0 = time()
3 print('\nCreating LSA/SVD Model...')
4 lsa_model = TruncatedSVD(n_components=500)
5 lsa = lsa_model.fit_transform(tfidf[idx_train,:])
6 print("Completed in %0.3fs." % (time() - t0))
7 done()
```

Creating LSA/SVD Model...
Completed in 335.923s.

0:04 / 0:04

```
In [16]: 1 # Plot Results to Find Knee
2 plt.figure(figsize=(10,5),facecolor='#d3d3d3')
3 plt.semilogy(lsa_model.singular_values_,'o',label='Topic Singular Values')
4 plt.vlines(50,*plt.ylim(),'r',label='Selected Cutoff')
5 plt.ylabel('log(Singular Values)')
6 plt.xlabel('Index')
7 _ = plt.title('LSA Singular Values',fontweight='bold')
8 plt.legend(loc='best')
```

Out[16]: <matplotlib.legend.Legend at 0x20e20100eb8>



Save LSA Model for Later Recall

```
In [17]: 1 # # Save LSA Model
2 # pickle_object(lsa_model, str(os.path.join(savePath, 'Lsa500TopicModel.sav')),
3 #                           save=True, addTime=False)
```

Build Latent Dirichlet Allocation Model

Build model using 50 topics as estimated roughly from the LSA model singular values.

In [18]:

```

1 # Build LDA Model
2 n_topics = 50 # number of topics in topic model
3
4 print("Fitting LDA models with tf features, "
5      "n_samples = %d and n_features = %d..." %
6      (tf[idx_train,:].shape[0],tf[idx_train,:].shape[1]))
7 lda = LatentDirichletAllocation(n_components=50, max_iter=5,
8                                 learning_method='online',
9                                 learning_offset=50.,
10                                random_state=0)
11 t0 = time()
12 lda.fit(tf[idx_train,:])
13 print("Completed in %.3fs." % (time() - t0))
14 done()

```

Fitting LDA models with tf features, n_samples = 944719 and n_features = 2000...
Completed in 634.803s.

0:04 / 0:04

Save LDA Model for Later Recall

In [19]:

```

1 # Save LDA Model
2 # pickle_object(Lda, str(os.path.join(savePath,'Lda50TopicModel.sav'))),
3 #           save=True, addTime=False)

```

Inspect Topics from LDA Model

In [20]:

```

1 print("\nTopics in LDA model:")
2 display_topics(lda, tf_feature_names, 9)

```

Topics in LDA model:

Topic 0 ======
wa|national|rise|labor|rate|boy|bus|nz|flu

Topic 1 ======
north|pay|force|industry|strike|safety|mining|blue|debate

Topic 2 ======
one|family|driver|arrested|japan|men|girl|break|sea

Topic 3 ======
man|child|charged|country|hope|hour|city|four|shooting

Topic 4 ======
abc|canberra|welcome|hill|cancer|expert|defence|staff|research

Topic 5 ======
murder|set|road|make|trial|brisbane|sign|give|nuclear

Topic 6 ======
drug|service|station|violence|target|stop|release|raid|impact

Topic 7 ======
face|go|public|community|rule|game|risk|black|truck

Topic 8 ======

Build Non-negative Matrix Factorization/ Probabilistic Latent Semantic Analysis Topic Model

NMF with beta loss set to KL divergence is equivalent to PLSA. Use 50 topics as before. Per for search over ratio of L1 to L2 norm used in distance calculation and find optimal value based on reconstruction error for a randomly sampled smaller dataset and then, train full on that value.

```
In [21]: 1  '''Setup Grid Search'''
2  def create_model(n_topics_ = 50, max_iter_ = 100, alpha_ = 0.1, l1_ratio_=0.5):
3      '''Build Tunable Model Function'''
4      model = NMF(n_components=n_topics_, random_state=42,
5                  beta_loss='kullback-leibler', solver='mu', init = 'random',
6                  max_iter=max_iter_, alpha=alpha_, l1_ratio=l1_ratio_)
7      return model
```

```
In [22]: 1  '''Setup Search Parameters'''
2  _alpha = np.linspace(0.8,1,10)
3  _l1_ratio = np.linspace(0,0.2,10)
4  AA, LL = np.meshgrid(_alpha,_l1_ratio)
5  aa, ll = AA.flatten(), LL.flatten()
6
7  # Build Small Test Set
8  np.random.seed(42)
9  _waste, idx_grid = train_test_split(idx_train, stratify=totalDF.loc[idx_train].Source,
10                                         test_size=0.05)
11 '''Start Grid Search'''
12 nt = 5
13 print('Grid Search Sample Set Size: ',len(idx_grid))
14 print('Topics for Grid Search: ',nt)
15 print('Beginning grid search on %d parameter combinations...' % (len(aa)))
16 results = []
17 t0 = time()
18 for i in range(len(aa)):
19     t1 = time()
20     mdl = create_model(n_topics_ = nt, max_iter_ = 100,
21                         alpha_ = aa[i], l1_ratio_ = ll[i])
22     print('\nLoop %d, alpha = %0.3f, l1_ratio = %0.3f' % (i, aa[i], ll[i]))
23     mdl.fit(tfidf[idx_grid,:])
24     merr = mdl.reconstruction_err_
25     results.append([aa[i],ll[i],merr])
26     print('Loop completed in %0.3fs. Error: %0.6f' %(time() - t1, merr))
27
28 clear_output(wait=True)
29 print("Completed in %0.3fs." % (time() - t0))
30
31 '''Display Results'''
32 df = pd.DataFrame(results,columns=['Alpha','L1 Ratio', 'Error'])
33 df.Error = (df.Error - df.Error.mean())/(df.Error.std())
34 display(df.loc[df.Error.idxmin()])
35 done()
```

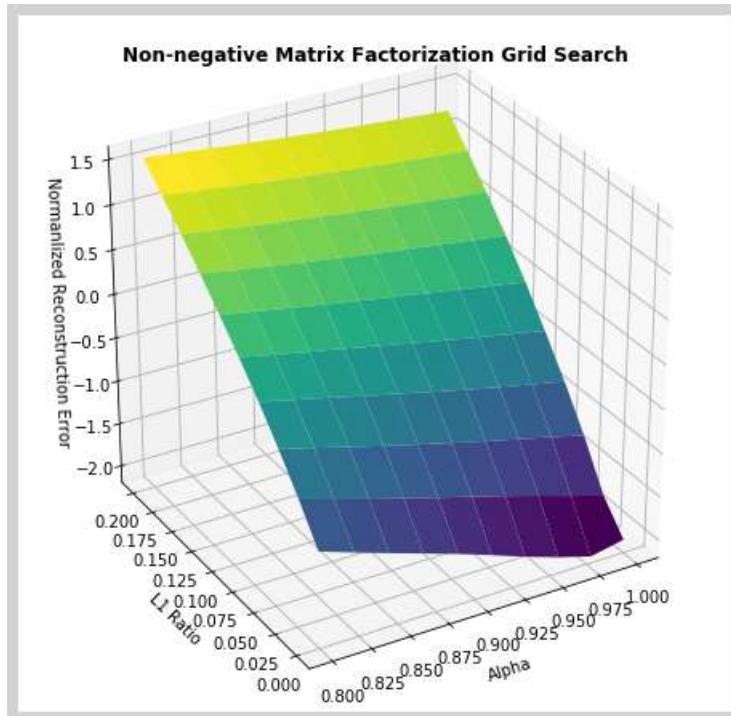
Completed in 793.911s.

Alpha	0.977778
L1 Ratio	0.000000
Error	-2.121364
Name:	8, dtype: float64

0.04 / 0:04

In [23]:

```
1  '''Plot Results of Grid Search'''
2  fig = plt.figure(figsize=(8,8), facecolor='#d3d3d3')
3  ax = fig.add_subplot(111, projection='3d')
4  Z = df.values[:,2].reshape(AA.shape)
5  #ax.plot_wireframe(AA, LL, Z, cmap='viridis', rstride=1, cstride=1)
6  ax.plot_surface(AA, LL, Z, cmap='viridis', edgecolor='none')
7  ax.set_xlabel('Alpha')
8  ax.set_ylabel('L1 Ratio')
9  ax.set_zlabel('Normalized Reconstruction Error')
10 ax.set_title('Non-negative Matrix Factorization Grid Search',
11              fontweight='bold')
12 ax.view_init(None, -120)
```



```
In [24]: 1 # Fit the NMF model
2 print("Fitting the NMF model (generalized Kullback-Leibler divergence) with "
3     "tf-idf features, n_samples=%d and n_features=%d..." %
4     (tfidf[idx_train,:].shape[0],tfidf[idx_train,:].shape[1]))
5 t0 = time()
6
7 # W is document membership weights (returned via transform), ...
8 # ... H is topic word dictionary (model.components_)
9 nmfKL = NMF(n_components=n_topics, random_state=42,
10             beta_loss='kullback-leibler', solver='mu', max_iter=1000, alpha=0.977,
11             l1_ratio=0)
12 nmf_doc_member_weights = nmfKL.fit_transform(tfidf[idx_train,:])
13 print("Completed in %0.3fs." % (time() - t0))
14 done()
```

Fitting the NMF model (generalized Kullback-Leibler divergence) with tf-idf features, n_samples=944719 and n_features=2000...
Completed in 475.833s.

0:04 / 0:04

Save NMF Model for Later Recall

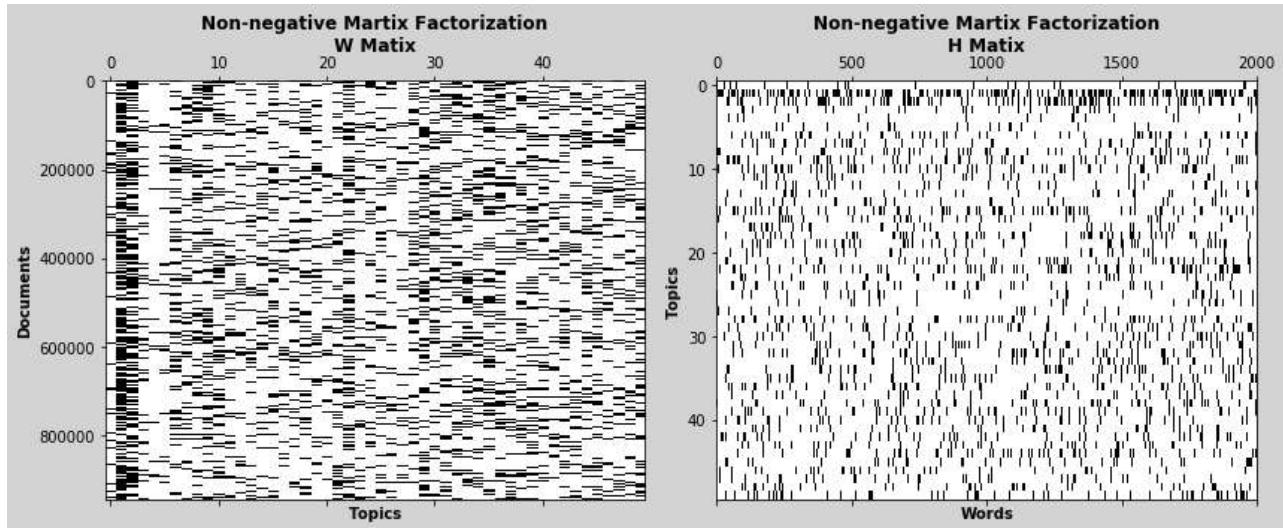
```
In [25]: 1 # # Save nmf Model
2 # pickle_object(nmfKL, str(os.path.join(savePath, 'nmfTopicModel.sav'))),
3 #           save=True, addTime=False)
```

Inspect Topics from PLSA Model

```
In [26]: 1 print("\nTopics in NMF model (generalized Kullback-Leibler divergence):")
2 display_topics(nmfKL, tfidf_feature_names, 9)
```

Topics in NMF model (generalized Kullback-Leibler divergence):
Topic 0 ======
interview|john|nr1|david|michael|smith|peter|scott|james
Topic 1 ======
security|stop|uk|rescue|fined|gun|bay|illegal|link
Topic 2 ======
port|give|crisis|demand|rudd|meeting|development|shire|see
Topic 3 ======
police|probe|chief|officer|arrest|station|investigate|hunt|target
Topic 4 ======
new|name|zealand|new zealand|new year|york|new york|technology|cabinet
Topic 5 ======
say|man|man charged|charged|man jailed|expert|need|man arrested|stabbing
Topic 6 ======
win|tour|england|award|star|title|race|king|stage
Topic 7 ======
fire|house|force|warning|threat|season|blaze|bushfire|risk
Topic 8 ======

```
In [27]: 1 #np.argsort(nmf_doc_member_weights[:,range(len(nmfKL.components_[:,0]))])[:,::-1].shape
2 nmf_doc_member_weights.shape
3 fig = plt.figure(tight_layout = True,figsize=(12,5),
4                   facecolor ='#d3d3d3')
5 ax = fig.add_subplot(121)
6 ax.spy(nmf_doc_member_weights)
7 ax.set_aspect('auto')
8 ax.set_title('Non-negative Martix Factorization\nW Matix', fontweight = 'bold')
9 ax.set_ylabel('Documents', fontweight = 'bold')
10 ax.set_xlabel('Topics', fontweight = 'bold')
11
12 ax = fig.add_subplot(122)
13 ax.spy(nmfKL.components_)
14 ax.set_aspect('auto')
15 ax.set_title('Non-negative Martix Factorization\nH Matix', fontweight = 'bold')
16 ax.set_ylabel('Topics', fontweight = 'bold')
17 _=ax.set_xlabel('Words', fontweight = 'bold')
```



Calculate Sentiment Scores for LDA and PLSA/NMF Topics

```
In [28]: 1 # Score each topic for sentiment
2 cutoff = 0.1
3 nmf_sent = []
4 vs = vader.SentimentIntensityAnalyzer()
5 for i in range(len(nmfKL.components_[:,0])):
6     x = np.argwhere(nmfKL.components_[i,:] > cutoff).ravel()
7     nmf_sent.append(vs.polarity_scores(" ".join([tfidf_feature_names[xi] for xi in x])))
8
9 nmf_sent_df = pd.DataFrame(nmf_sent)
```

In [29]:

```

1 # Score each topic for sentiment
2 cutoff = 0.1
3 lda_sent = []
4 vs = vader.SentimentIntensityAnalyzer()
5 for i in range(len(lda.components_[:,0])):
6     x = np.argwhere(lda.components_[i,:] > cutoff).ravel()
7     lda_sent.append(vs.polarity_scores(" ".join([tf_feature_names[xi] for xi in x])))
8
9 lda_sent_df = pd.DataFrame(lda_sent)
10 lda_sent_df.head()

```

Out[29]:

	compound	neg	neu	pos
0	-0.6369	0.215	0.670	0.115
1	0.3612	0.143	0.724	0.133
2	-0.5106	0.147	0.713	0.140
3	-0.7845	0.265	0.603	0.132
4	0.1280	0.190	0.599	0.212

Visualize Distribution for LDA and PLSA/NMF Sentiment Scores

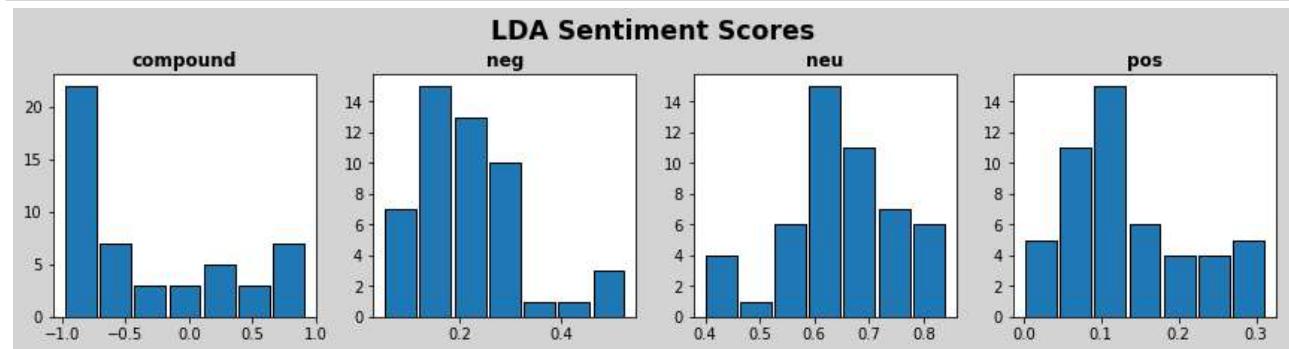
1. LDA Sentiment

In [30]:

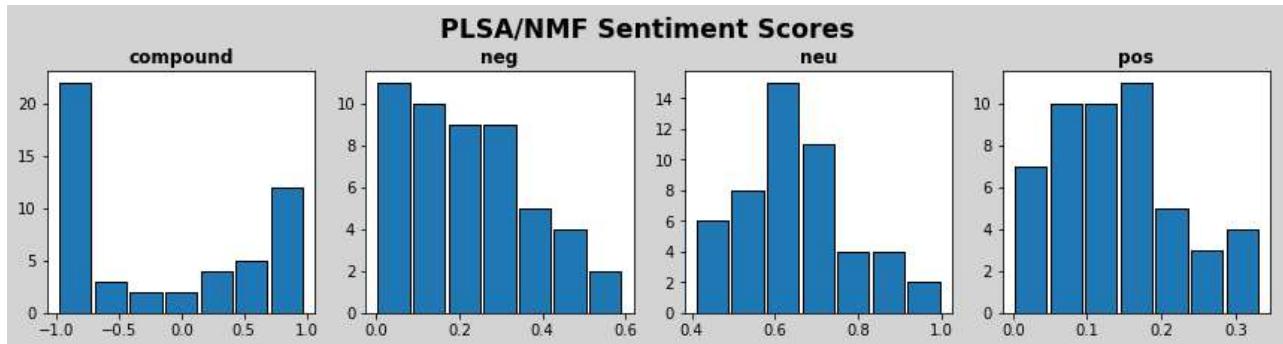
```

1 def plot_sent_hist(df,title,fsiz = (12,3)):
2     n_bins = int(np.ceil(1 + 3.22*np.log10(len(df))))
3     f,axes = plt.subplots(nrows=1, ncols=4, tight_layout = True, figsize = fsiz,
4                           facecolor='#d3d3d3')
5     df.hist(bins=n_bins, ax= axes, rwidth=0.9, align='mid', ec='k',grid=False)
6     for ax in axes.ravel():
7         ax.set_title(ax.get_title(),fontweight='bold')
8         _ = f.suptitle(title,fontweight='bold',
9                         x=0.5, y=1.05, ha='center', fontsize='xx-large')
10
11 plot_sent_hist(lda_sent_df,'LDA Sentiment Scores')

```



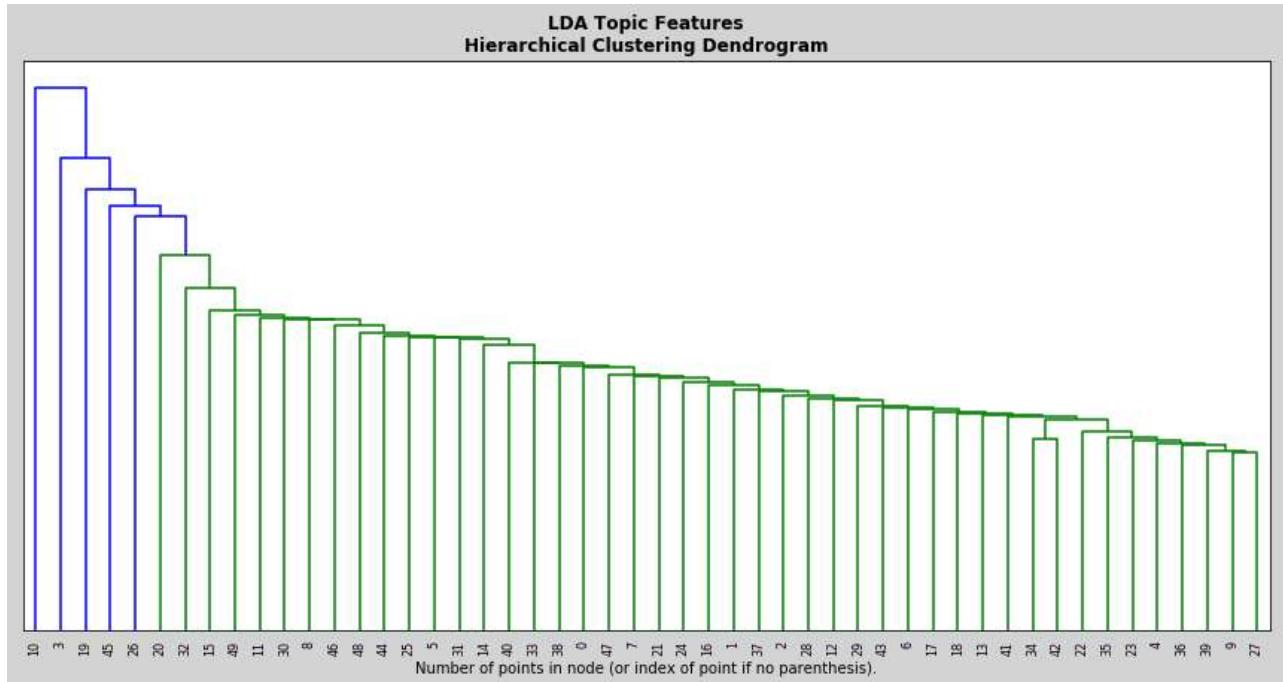
```
In [31]: 1 plot_sent_hist(nmf_sent_df, 'PLSA/NMF Sentiment Scores')
```



Hierarchical Clustering

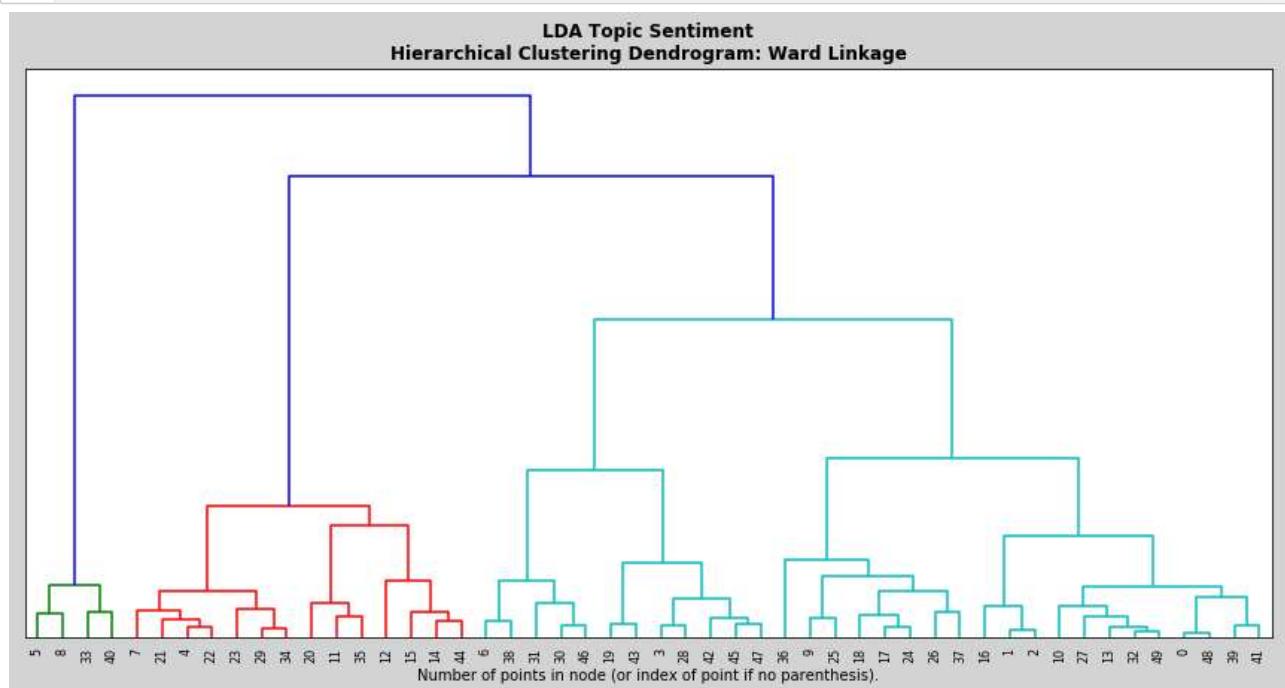
Hierarchical Clustering for LDA Topics

```
In [32]: 1 # copy dataframes
2 hc_lda_df = lda_sent_df.copy(deep=True)
3 hc_nmf_df = nmf_sent_df.copy(deep=True)
4
5 _ = plot_dendrogram(lda.components_,
6                      'LDA Topic Features\nHierarchical Clustering Dendrogram' )
```

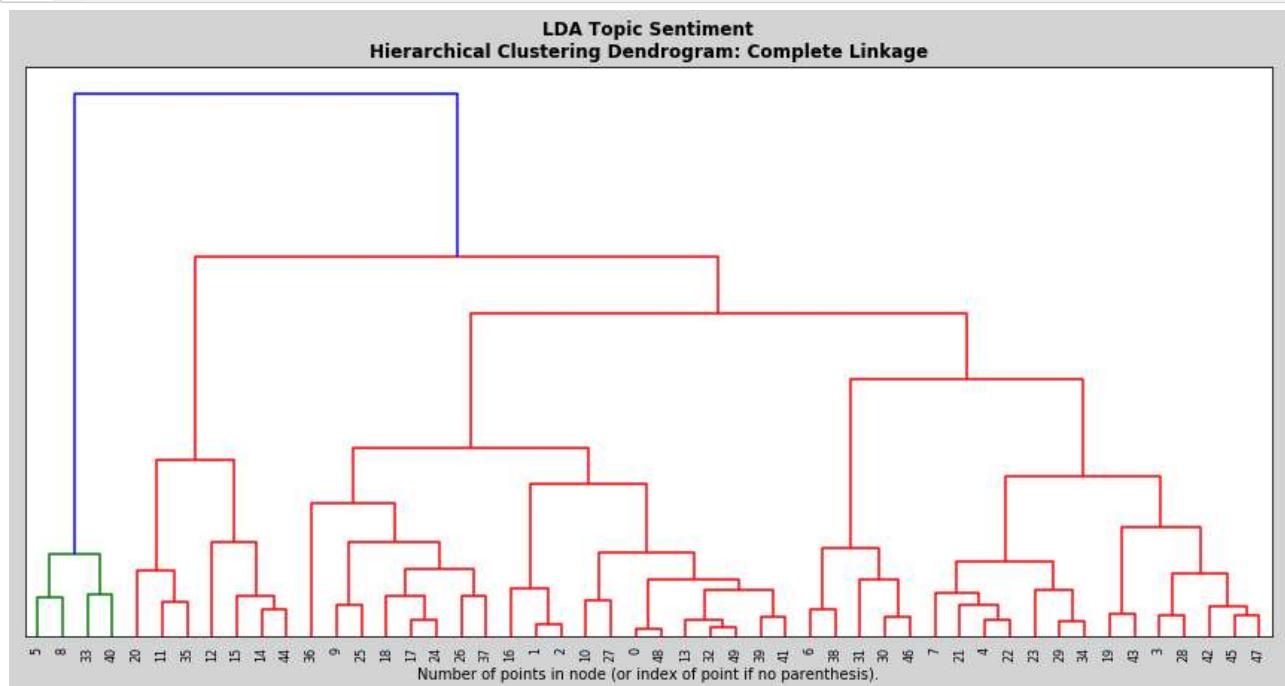


Hierarchical Clustering LDA Topic Sentiments

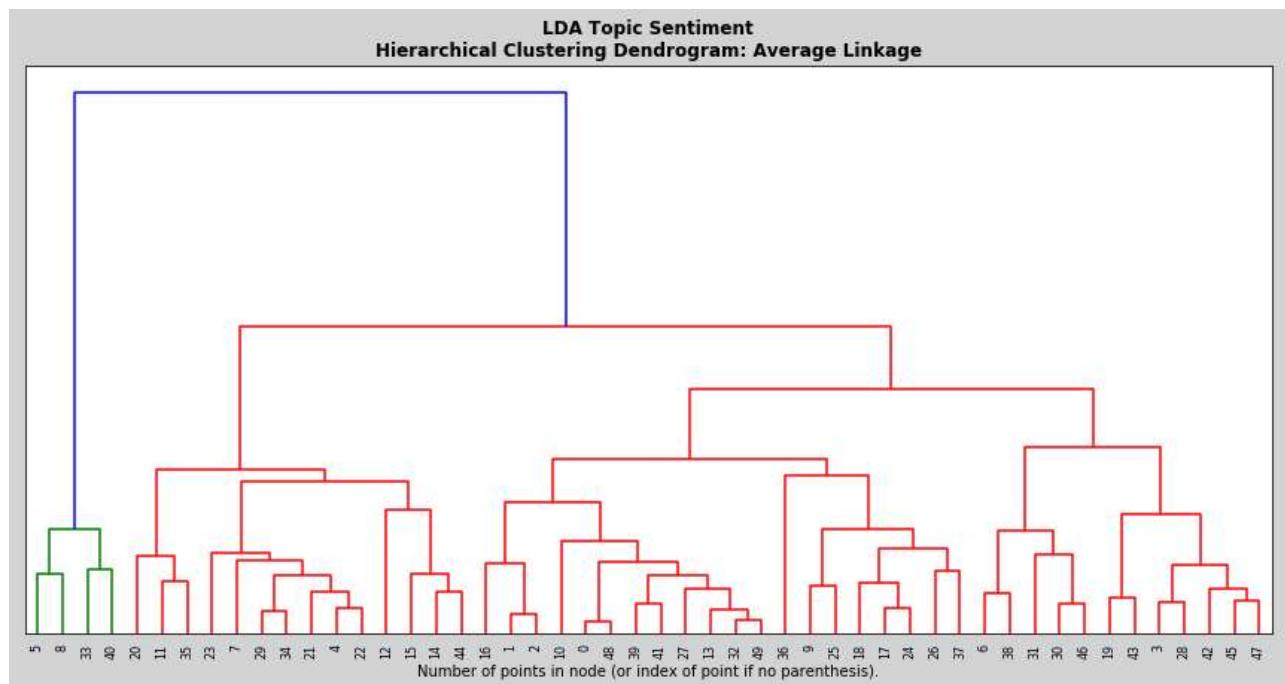
```
In [33]: 1 _ = plot_dendrogram(hc_lda_df[['neg', 'neu', 'pos']],
2                           'LDA Topic Sentiment\nHierarchical Clustering Dendrogram: Ward Linkage')
```



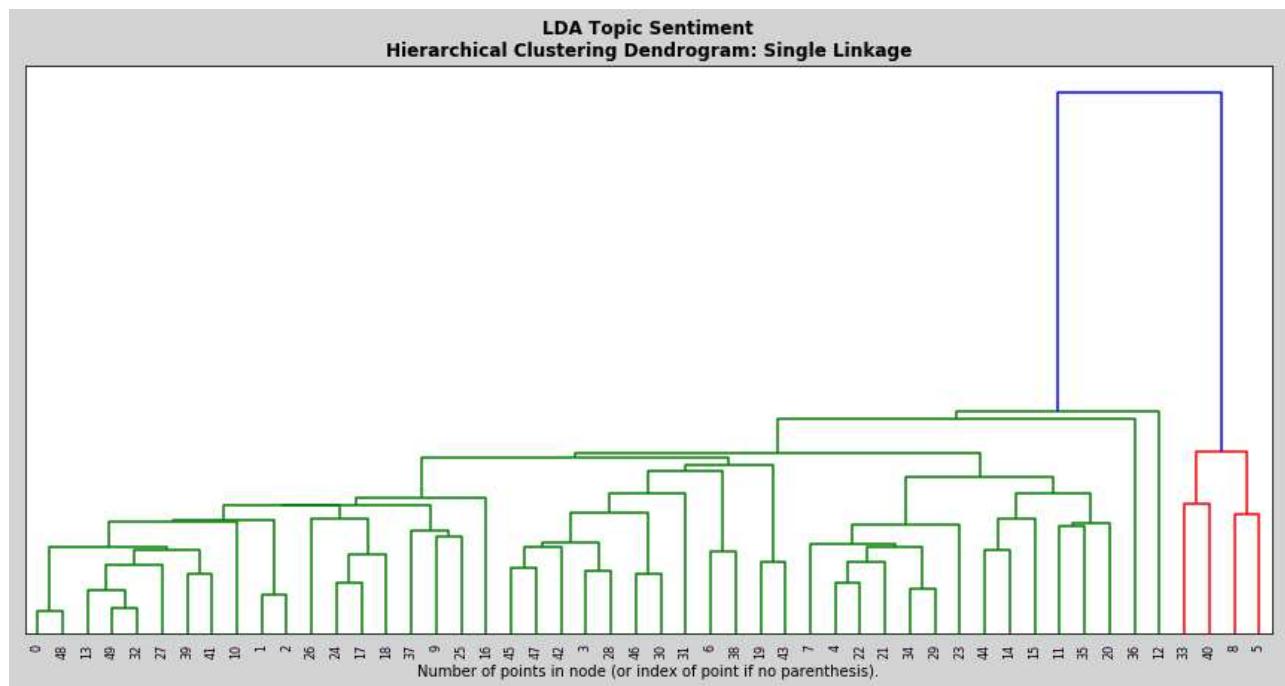
```
In [34]: 1 _ = plot_dendrogram(hc_lda_df[['neg', 'neu', 'pos']],
2                           'LDA Topic Sentiment\nHierarchical Clustering Dendrogram: Complete Linkage',
3                           linkage_ = 'complete')
```



```
In [35]: 1 _ = plot_dendrogram(hc_lda_df[['neg','neu','pos']],
2                         'LDA Topic Sentiment\nHierarchical Clustering Dendrogram: Average Linkage',
3                         linkage_ = 'average')
```

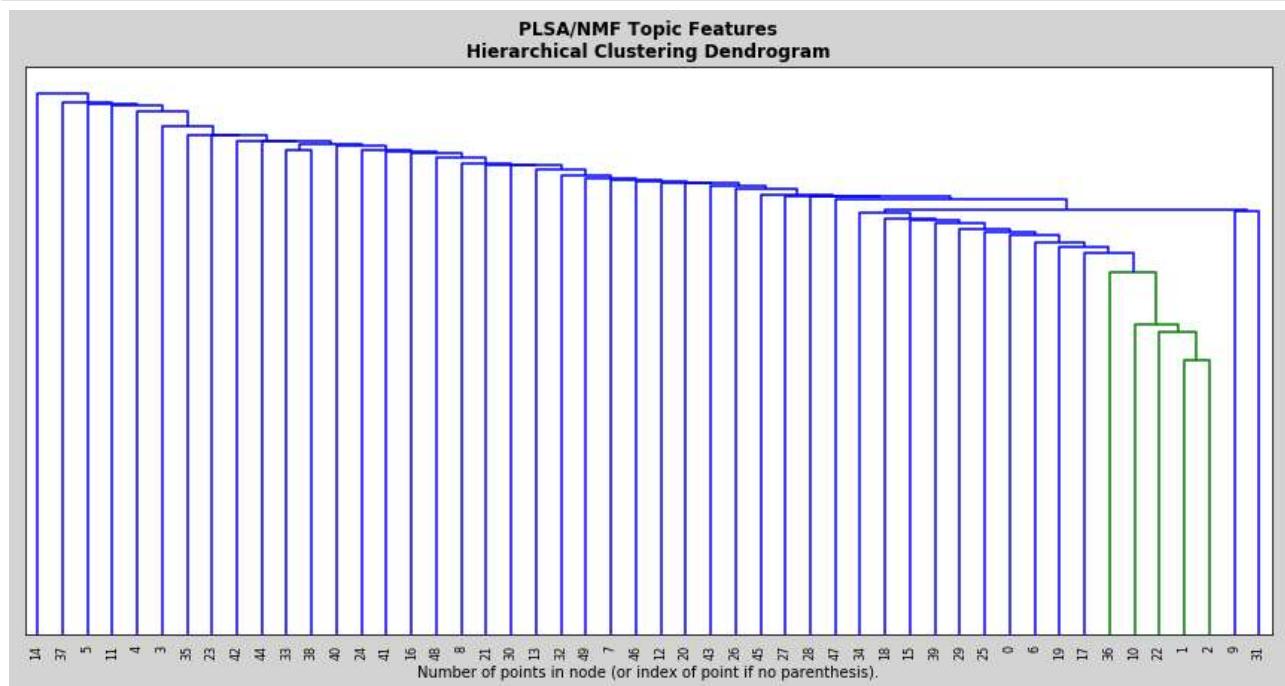


```
In [36]: 1 _ = plot_dendrogram(hc_lda_df[['neg', 'neu', 'pos']],
2                           'LDA Topic Sentiment\nHierarchical Clustering Dendrogram: Single Linkage',
3                           linkage_ = 'single')
```



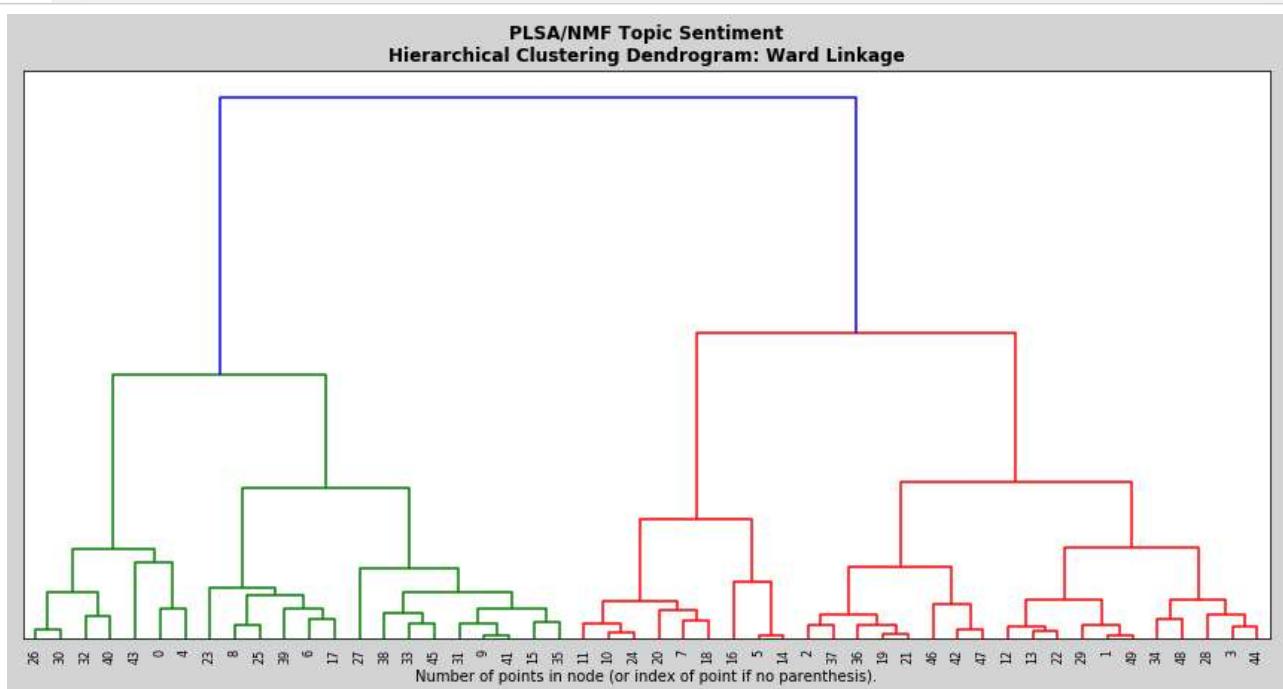
Hierarchical Clustering for NMF Topics

```
In [37]: 1 _ = plot_dendrogram(nmfKL.components_,
2                           'PLSA/NMF Topic Features\nHierarchical Clustering Dendrogram')
```

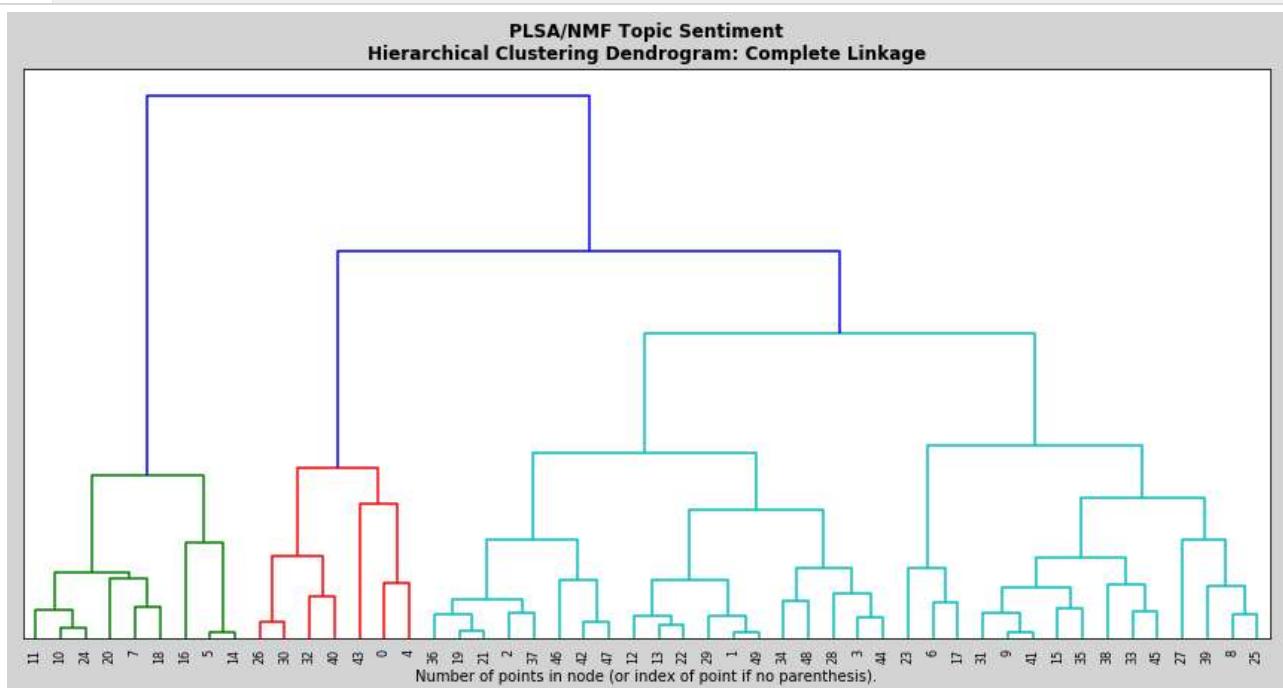


Hierarchical Clustering for NMF Topic Sentiments

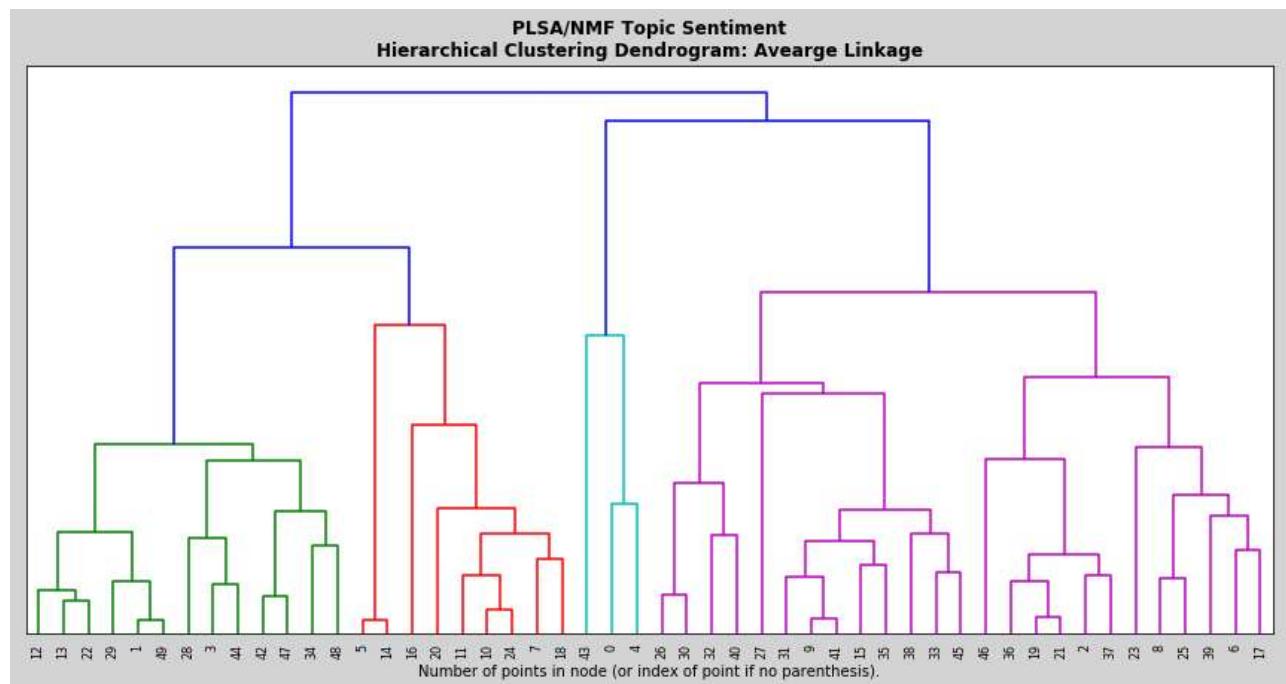
```
In [38]: 1 _ = plot_dendrogram(hc_nmf_df[['neg','neu','pos']],
2                           'PLSA/NMF Topic Sentiment\nHierarchical Clustering Dendrogram: Ward Linkage')
```



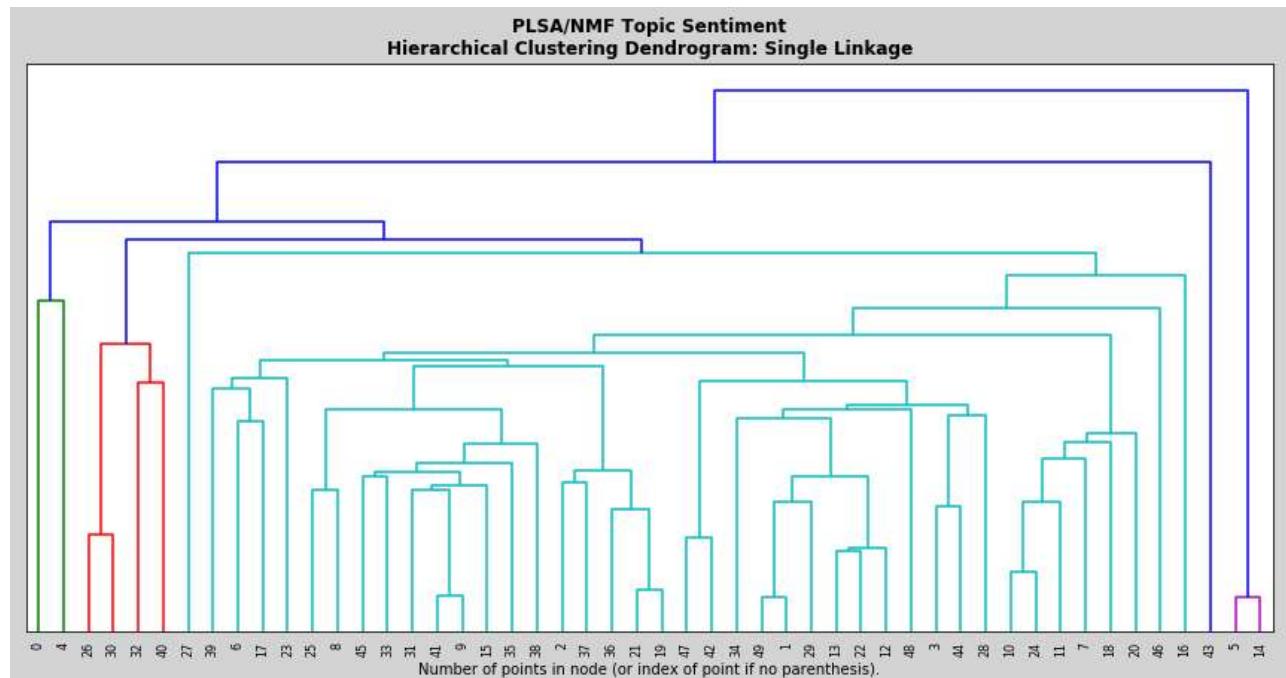
```
In [39]: 1 _ = plot_dendrogram(hc_nmf_df[['neg','neu','pos']],
2                           'PLSA/NMF Topic Sentiment\nHierarchical Clustering Dendrogram: Complete Linkage',
3                           linkage_='complete')
```



```
In [40]: 1 _ = plot_dendrogram(hc_nmf_df[['neg','neu','pos']],
2                         'PLSA/NMF Topic Sentiment\nHierarchical Clustering Dendrogram: Averge Linkage',
3                         linkage_ = 'average')
```



```
In [41]: 1 _ = plot_dendrogram(hc_nmf_df[['neg','neu','pos']],
2                         'PLSA/NMF Topic Sentiment\nHierarchical Clustering Dendrogram: Single Linkage',
3                         linkage_ = 'single')
```



Heirarchical Clustering with LDA and PLSA/NMF Topic Sentiments

```
In [42]: 1 def hierarchical_wrapper(df, n_clusters_, msg):
2     # Perform Clustering
3     t0 = time()
4     print('Beginning Hierarchical Clustering on %s Data with %d clusters...' % (msg,n_clusters_))
5     mdl = AgglomerativeClustering(n_clusters = n_clusters_)
6     mdl.fit(df[['neg','neu','pos']])
7     # Label Topics
8     df['Label'] = mdl.labels_
9     print("completed in %0.3fs" % (time() - t0))
10    return mdl, df
```

```
In [43]: 1 # Perform Clustering
2 hc_lda_clusters = 4
3 hc_nmf_clusters = 4
4
5 hc_lda, hc_lda_df = hierarchical_wrapper(hc_lda_df,hc_lda_clusters, 'LDA')
6 hc_nmf, hc_nmf_df = hierarchical_wrapper(hc_nmf_df,hc_nmf_clusters, 'PLSA/NMF')
7 done()
```

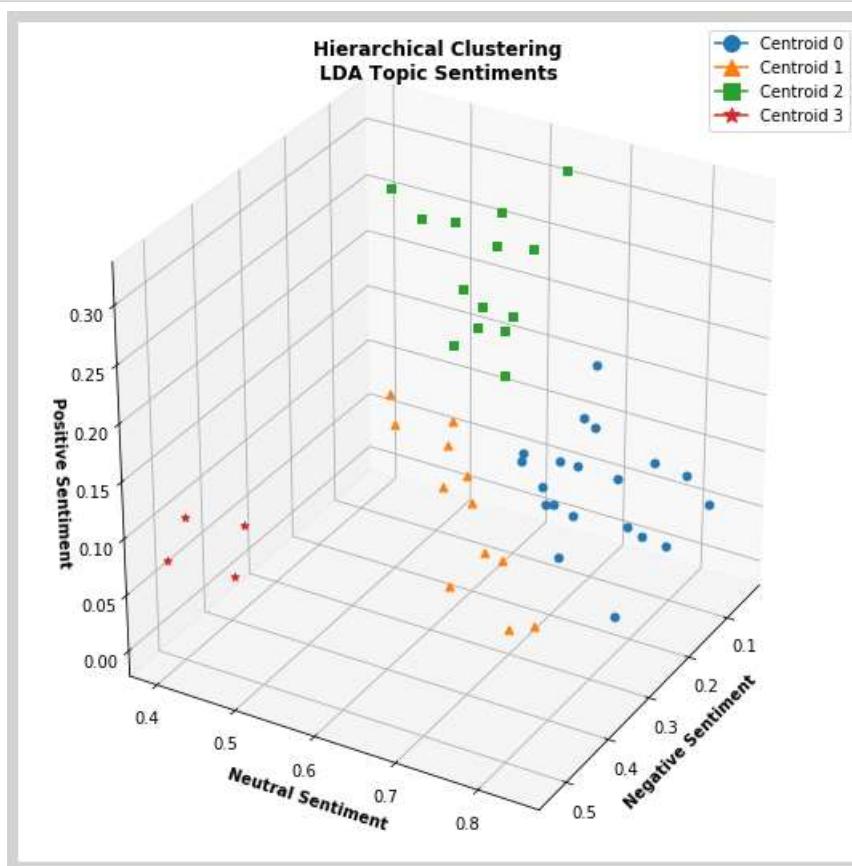
Beginning Hierarchical Clustering on LDA Data with 4 clusters...
completed in 0.004s
Beginning Hierarchical Clustering on PLSA/NMF Data with 4 clusters...
completed in 0.003s

0:04 / 0:04

Plot Results of Heirarchical Clustering

In [44]:

```
1 # Plot the results
2 ang = [(None,30)] #,(10,90),(100,0),(10,0)
3 plot_clustering(None,
4                  hc_lda_df,'Hierarchical Clustering\nLDA Topic Sentiments',
5                  ang)
```

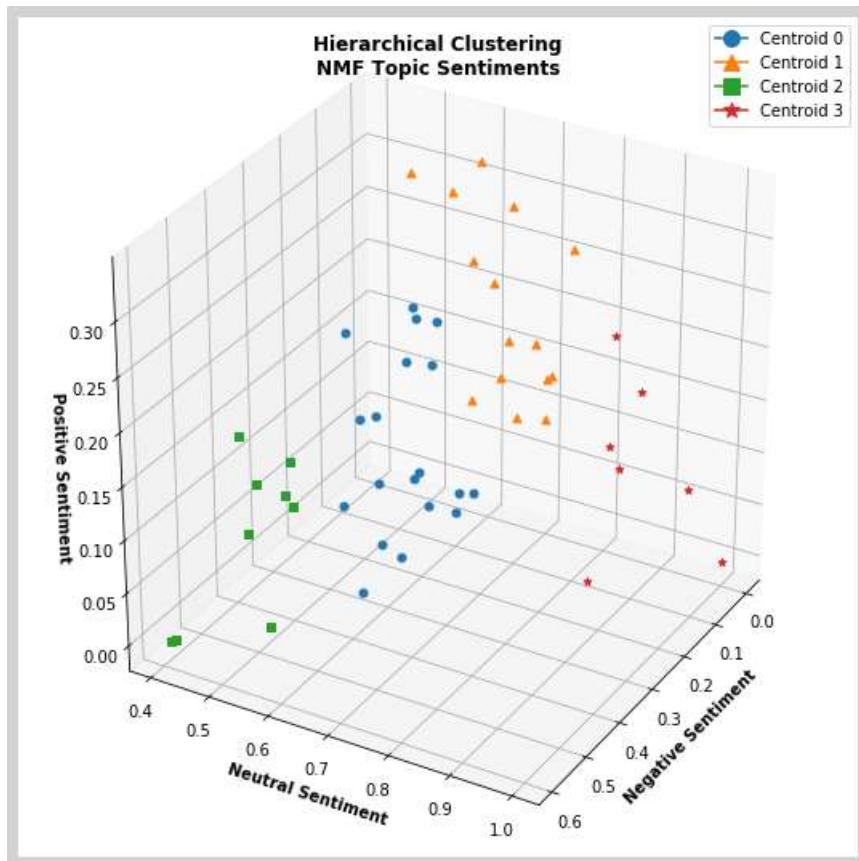


In [45]:

```

1 # Plot the results
2 ang = [(None,30)] #,(10,90),(100,0),(10,0)
3 plot_clustering(None,
4                  hc_nmf_df,'Hierarchical Clustering\nNMF Topic Sentiments',
5                  ang)

```



Score Hierarchical Clustering Models

In [46]:

```

1 def score_model(df):
2     ch_score = calinski_harabasz_score(df[['neg','neu','pos']],
3                                         df['Label'])
4     sil_score = silhouette_score(df[['neg','neu','pos']],
5                                   df['Label'],metric='euclidean')
6     return np.array([ch_score,sil_score])

```

In [47]:

```

1 # Create DataFrame to Tabulate Model Scores
2 model_scores_df = pd.DataFrame(np.vstack([np.append(score_model(hc_lda_df),hc_lda_clusters),
3                                           np.append(score_model(hc_nmf_df),hc_nmf_clusters)]).T,
4                                     index = ['CH Score', 'Silhouette Score','# Clusters'],
5                                     columns = ['Hierarchical LDA','Hierarchical PLSA/NMF'])
6 model_scores_df

```

Out[47]:

	Hierarchical LDA	Hierarchical PLSA/NMF
CH Score	63.920187	59.160354
Silhouette Score	0.464105	0.391698
# Clusters	4.000000	4.000000

Heirarchical Clustering Summary

From the dendograms, it's clear that clustering should be done on the sentiment scores for each topic. This makes sense as the dimensionality reduction to 50 topics is done to identify strictly unique sentiments/feature vectors. The complete linkage calculation creates the best looking clustering per the vertical distance between levels. Complete or maximum linkage uses the maximum distances between all observations of the two sets.

The hierarchical clustering shows that 2-4 clusters are appropriate. 4 was used as it seems the best choice.

KMeans Clustering

Perform KMeans Clustering and Record Sum of Squared Error for Several Numbers of Clusters on LDA and PLSA/NMF Topic Sentiments

In [48]:

```

1 # copy dataframe
2 km_lda_df = lda_sent_df.copy(deep=True)
3 km_nmf_df = nmf_sent_df.copy(deep=True)
4
5 def run_nclusters(fun_build_mdl, fun_mdl_score, lda_df,nmf_df,range_clusters = (1,10)):
6     score_nmf = []
7     score_lda = []
8     nclust = []
9     for i in range(*range_clusters,1):
10         nclust.append(i)
11         mdl = fun_build_mdl(i)
12         print("Clustering with n_clusters = %d and model %s" % (i,mdl))
13         t0 = time()
14         mdl.fit(nmf_df[['neg','neu','pos']])
15         score_nmf.append(fun_mdl_score(mdl,nmf_df))
16         mdl.fit(lda_df[['neg','neu','pos']])
17         score_lda.append(fun_mdl_score(mdl,lda_df))
18         print("completed in %0.3fs" % (time() - t0))
19     return nclust, score_lda, score_nmf
20
21 bld_km = lambda i: KMeans(n_clusters=i, init='k-means++', max_iter=100, n_init=10, algorithm='full',
22                           random_state=42, tol = 1e-8, verbose=0)
23 score_km = lambda km,df: km.inertia_
24
25 X = run_nclusters(bld_km, score_km, km_lda_df, km_nmf_df)
26 done()

```

```

Clustering with n_clusters = 1 and model KMeans(algorithm='full', max_iter=100, n_clusters=1, random_st
ate=42, tol=1e-08)
completed in 0.805s
Clustering with n_clusters = 2 and model KMeans(algorithm='full', max_iter=100, n_clusters=2, random_st
ate=42, tol=1e-08)
completed in 0.750s
Clustering with n_clusters = 3 and model KMeans(algorithm='full', max_iter=100, n_clusters=3, random_st
ate=42, tol=1e-08)
completed in 0.774s
Clustering with n_clusters = 4 and model KMeans(algorithm='full', max_iter=100, n_clusters=4, random_st
ate=42, tol=1e-08)
completed in 0.956s
Clustering with n_clusters = 5 and model KMeans(algorithm='full', max_iter=100, n_clusters=5, random_st
ate=42, tol=1e-08)
completed in 1.029s
Clustering with n_clusters = 6 and model KMeans(algorithm='full', max_iter=100, n_clusters=6, random_st
ate=42, tol=1e-08)
completed in 0.924s
Clustering with n_clusters = 7 and model KMeans(algorithm='full', max_iter=100, n_clusters=7, random_st
ate=42, tol=1e-08)
completed in 0.937s
Clustering with n_clusters = 8 and model KMeans(algorithm='full', max_iter=100, random_state=42, tol=1e
-08)
completed in 1.143s
Clustering with n_clusters = 9 and model KMeans(algorithm='full', max_iter=100, n_clusters=9, random_st
ate=42, tol=1e-08)
completed in 1.081s

```

0:04 / 0:04

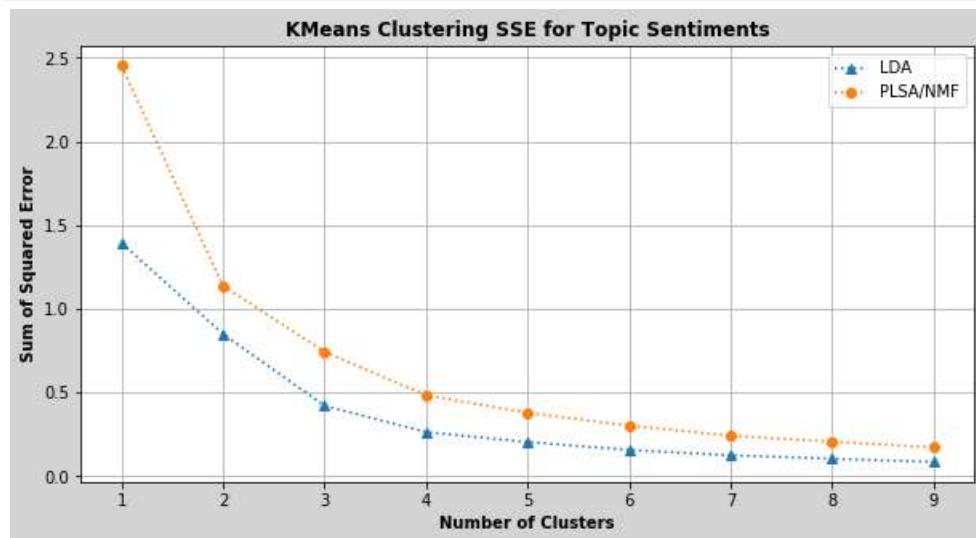
Plot SSE for LDA and PLSA/NMF Topics

In [49]:

```

1 # Plot Results to Find Knee if One Exists
2 plt.figure(figsize=(10,5),facecolor='#d3d3d3')
3 plt.plot(X[0],X[1],':^',label = 'LDA')
4 plt.plot(X[0],X[2],':o',label = 'PLSA/NMF')
5 plt.ylabel('Sum of Squared Error',fontweight='bold')
6 plt.xlabel('Number of Clusters',fontweight='bold')
7 plt.legend(loc='best')
8 plt.grid()
9 _ = plt.title('KMeans Clustering SSE for Topic Sentiments',fontweight='bold')

```



Initial Clustering Summary

For the quick clustering parameters, the LDA topic sentiments had a generally lower sum of squared errors than the PLSA topic model. Both models appear to have a knee in the 3 - 4 cluster range, as we saw duplicated in the hierarchical models.

KMeans 4 Clusters on LDA and PLSA/NMF Topic Sentiments and Plot

In [50]:

```

1 def kmeans_wrapper(df, n_clusters_, msg):
2     # Perform Kmeans
3     t0 = time()
4     print('Beginning Kmeans Clustering on %s Data with %d clusters...' % (msg,n_clusters_))
5     km = KMeans(n_clusters=n_clusters_, init='k-means++', max_iter=300, n_init=1000, algorithm='full',
6                  random_state=42, tol = 1e-8, verbose=0)
7     km.fit(df[['neg','neu','pos']])
8     # Label Topics
9     df['Label'] = km.labels_
10    print("completed in %0.3fs" % (time() - t0))
11    return km,df
12

```

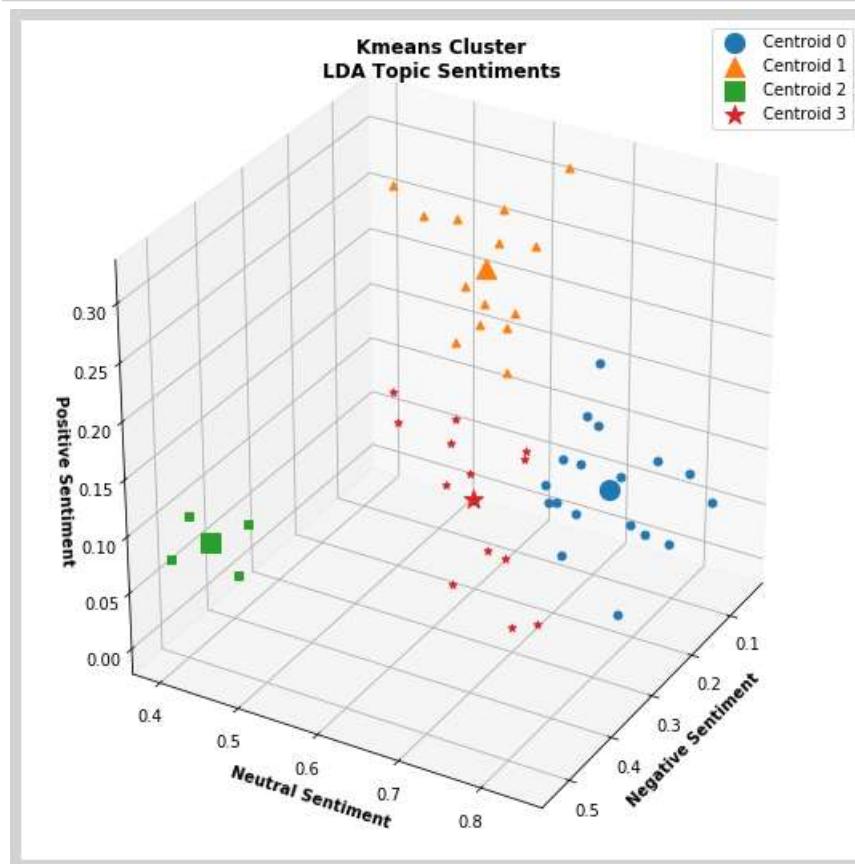
```
In [51]: 1 km_clusters_ = 4 # Number of clusters for Kmeans
2 km_nmf, km_nmf_df = kmeans_wrapper(km_nmf_df, km_clusters_, 'PLSA/NMF')
3 km_lda, km_lda_df = kmeans_wrapper(km_lda_df, km_clusters_, 'LDA')
4 done()
```

Beginning Kmeans Clustering on PLSA/NMF Data with 4 clusters...
 completed in 39.080s
 Beginning Kmeans Clustering on LDA Data with 4 clusters...
 completed in 39.663s

0:04 / 0:04

Plot Kmeans Results for LDA and PLSA/NMF

```
In [52]: 1 # Plot the results
2 ang = [(None,30), (10,90), (100,0), (10,0)]
3 plot_clustering(km_lda.cluster_centers_,
4                  km_lda_df, 'Kmeans Cluster\nnLDA Topic Sentiments', ang)
5
```

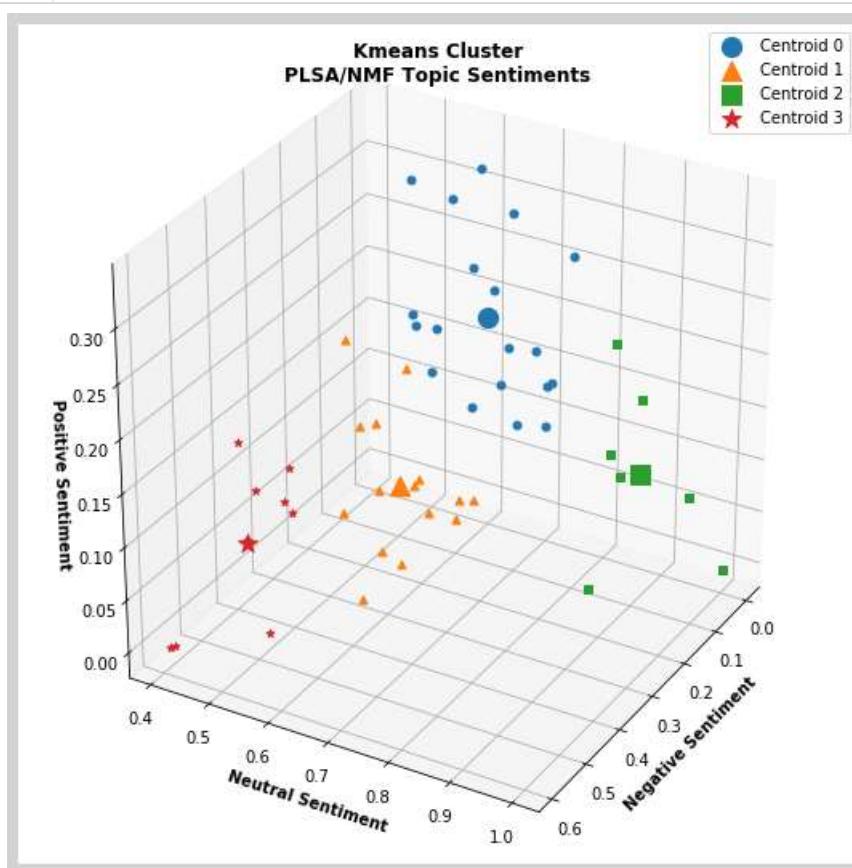


In [53]:

```

1 # Plot the results
2 ang = [(None,30), #,(10,90),(100,0),(10,0)]
3 plot_clustering(km_nmf.cluster_centers_,
4                  km_nmf_df, 'Kmeans Cluster\nPLSA/NMF Topic Sentiments', ang)
5

```



Score KMeans Models

In [54]:

```

1 # Create DataFrame to Tabulate Model Scores
2 model_scores_df['KMeans LDA'] = np.append(score_model(km_lda_df),km_clusters_).reshape(-1,1)
3 model_scores_df['KMeans PLSA/NMF'] = np.append(score_model(km_nmf_df),km_clusters_).reshape(-1,1)
4 model_scores_df

```

Out[54]:

	Hierarchical LDA	Hierarchical PLSA/NMF	KMeans LDA	KMeans PLSA/NMF
CH Score	63.920187	59.160354	65.412290	62.244816
Silhouette Score	0.464105	0.391698	0.451274	0.406453
# Clusters	4.000000	4.000000	4.000000	4.000000

GMM Clustering

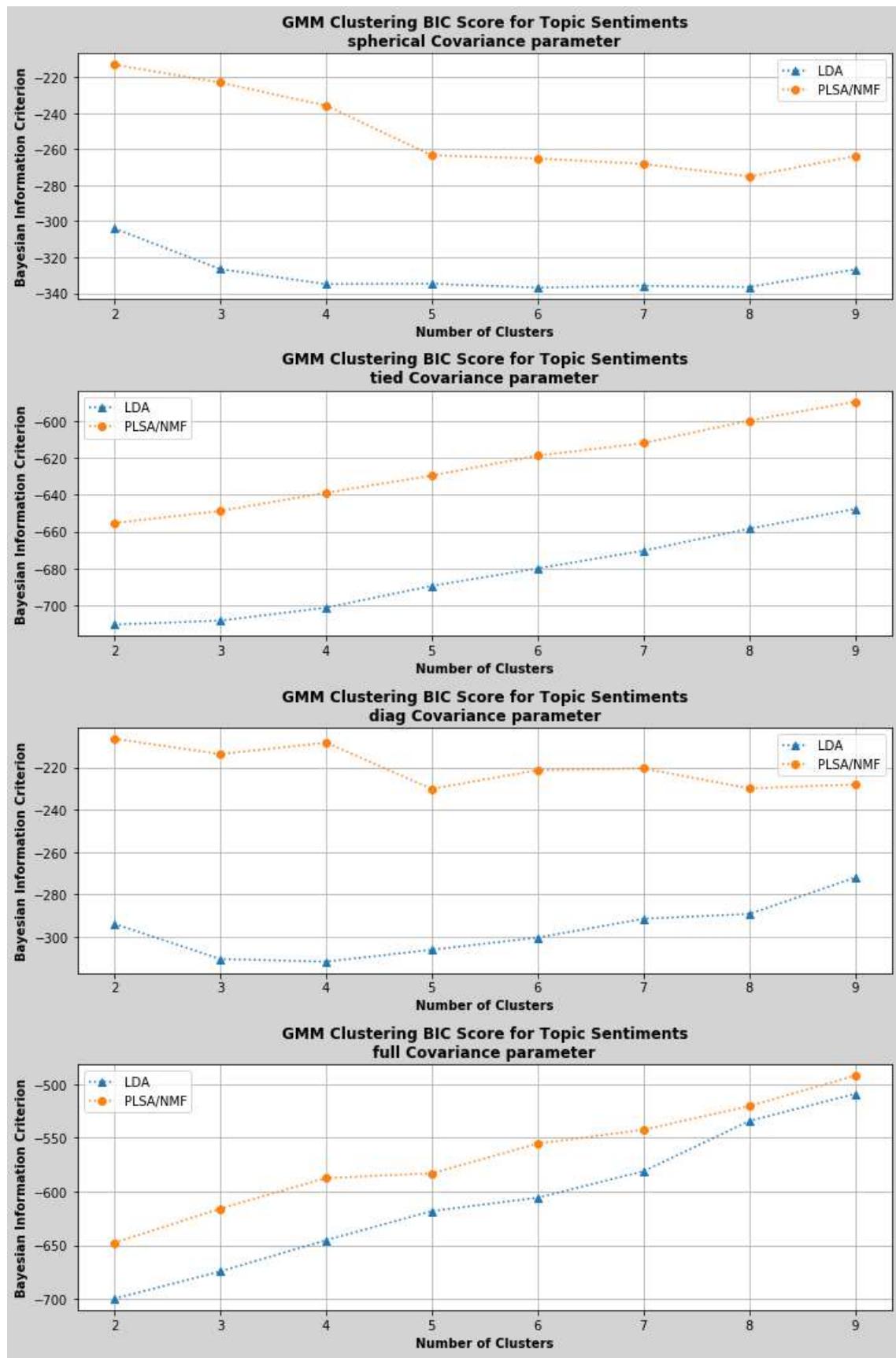
In [55]:

```
1 # copy dataframe
2 gmm_lda_df = lda_sent_df.copy(deep=True)
3 gmm_nmf_df = nmf_sent_df.copy(deep=True)
4
5
6 score_gmm = lambda gmm,df: gmm.bic(df[['neg','neu','pos']])
7 cv_type = ['spherical', 'tied', 'diag', 'full']
8 X = []
9 for cv in cv_type:
10     bld_gmm = lambda i: GMM(n_components=i, covariance_type=cv, max_iter=100, n_init = 10,
11                             random_state=42, verbose = 0)
12     X.append(run_nclusters(bld_gmm, score_gmm, gmm_lda_df, gmm_nmf_df, range_clusters=(2,10)))
13 done()
```

Clustering with n_clusters = 2 and model GaussianMixture(covariance_type='spherical', n_components=2, n_init=10,
random_state=42)
completed in 0.067s
Clustering with n_clusters = 3 and model GaussianMixture(covariance_type='spherical', n_components=3, n_init=10,
random_state=42)
completed in 0.076s
Clustering with n_clusters = 4 and model GaussianMixture(covariance_type='spherical', n_components=4, n_init=10,
random_state=42)
completed in 0.077s
Clustering with n_clusters = 5 and model GaussianMixture(covariance_type='spherical', n_components=5, n_init=10,
random_state=42)
completed in 0.087s
Clustering with n_clusters = 6 and model GaussianMixture(covariance_type='spherical', n_components=6, n_init=10,
random_state=42)
..

In [56]:

```
1 # Plot Results to Find Knee if One Exists
2 fig = plt.figure(tight_layout = True,figsize=(10,15),facecolor="#d3d3d3")
3 for i,x in enumerate(X):
4     ax = fig.add_subplot(4,1,i+1)
5     ax.plot(x[0],x[1],':^',label = 'LDA')
6     ax.plot(x[0],x[2],':o',label = 'PLSA/NMF')
7     ax.set_ylabel('Bayesian Information Criterion',fontweight='bold')
8     ax.set_xlabel('Number of Clusters',fontweight='bold')
9     ax.legend(loc='best')
10    ax.grid(True)
11    t = 'GMM Clustering BIC Score for Topic Sentiments\n'
12    t = t + '{} Covariance parameter'.format(cv_type[i])
13    ax.set_title(t,fontweight='bold')
```



BIC Results

It appears that the covariance parameter 'spherical' produces slightly better results (lower BIC is better). From that plot, it looks as though the LDA data should have 4 clusters and the PLSA/NMF 5. We'll use 4 for continuity.

Perform GMM Clustering and Plot Results

```
In [57]: 1 def gmm_wrapper(df, n_clusters_, msg):
2     # Perform GMM
3     t0 = time()
4     print('Beginning GMM Clustering on %s Data with %d clusters...' % (msg,n_clusters_))
5     mdl = GMM(n_components=n_clusters_, covariance_type='spherical', max_iter=300, n_init = 1000,
6                random_state=42, verbose = 0)
7     mdl.fit(df[['neg','neu','pos']])
8     # Label Topics
9     df['Label'] = mdl.predict(df[['neg','neu','pos']])
10    print("completed in %0.3fs" % (time() - t0))
11    return mdl,df
12
13 gmm_lda_clusters = 4 # Number of clusters
14 gmm_nmf_clusters = 4
15 gmm_nmf, gmm_nmf_df = gmm_wrapper(gmm_nmf_df, gmm_nmf_clusters, 'PLSA/NMF')
16 gmm_lda, gmm_lda_df = gmm_wrapper(gmm_lda_df, gmm_lda_clusters, 'LDA')
17 done()
```

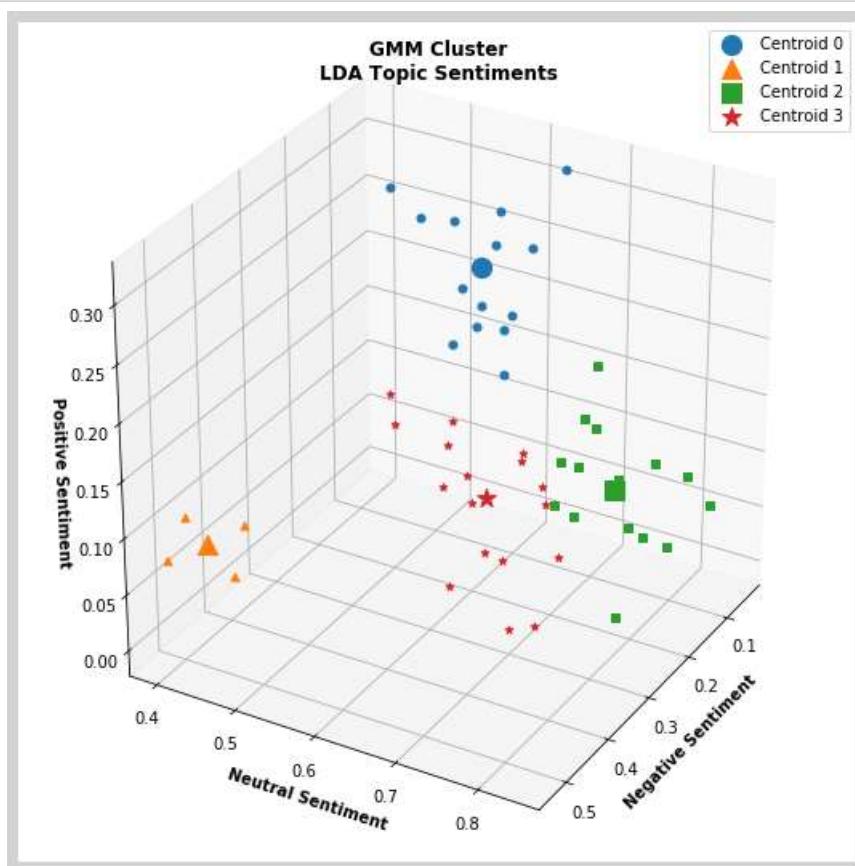
```
Beginning GMM Clustering on PLSA/NMF Data with 4 clusters...
completed in 4.083s
Beginning GMM Clustering on LDA Data with 4 clusters...
completed in 3.112s
```

0:04 / 0:04

Plot GMM Results for LDA and PLSA/NMF

In [58]:

```
1 # Plot the results
2 ang = [(None,30)] #[,(10,90),(100,0),(10,0)]
3 plot_clustering(gmm_lda.means_
4                 ,gmm_lda_df,'GMM Cluster\nLDA Topic Sentiments', ang)
5
```

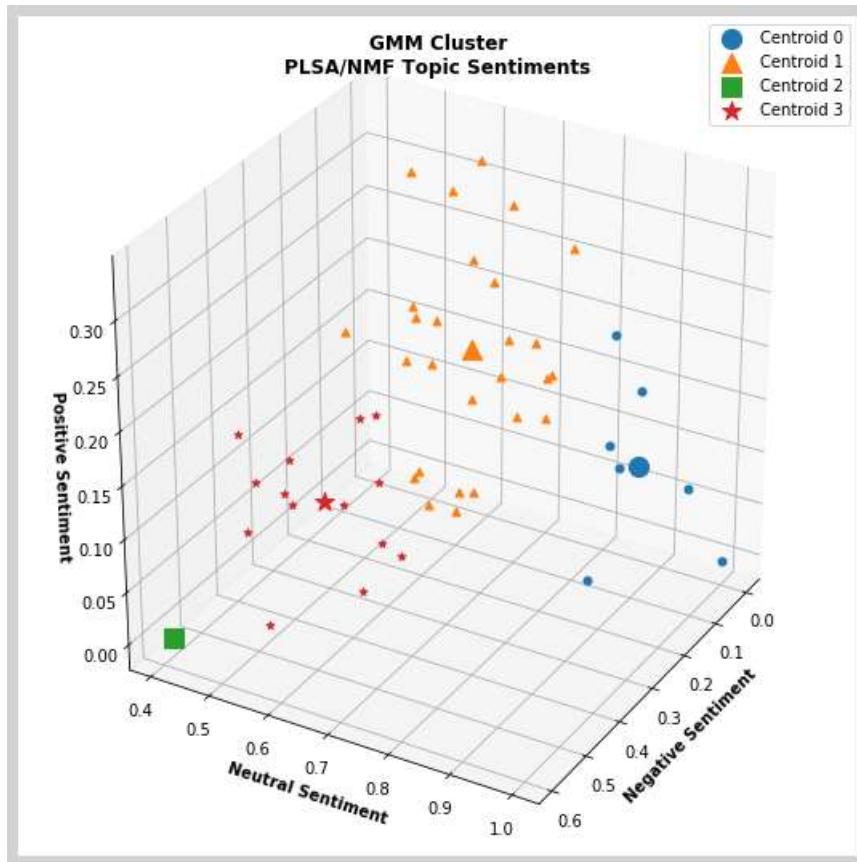


In [59]:

```

1 # Plot the results
2 ang = [(None,30), #(None,30), (10,90), (100,0), (10,0)]
3 plot_clustering(gmm_nmf.means_,
4                  gmm_nmf_df, 'GMM Cluster\nPLSA/NMF Topic Sentiments', ang)

```



Score GMM Results

In [60]:

```

1 # Create DataFrame to Tabulate Model Scores
2 model_scores_df['GMM LDA'] = np.append(score_model(gmm_lda_df),gmm_lda_clusters).reshape(-1,1)
3 model_scores_df['GMM PLSA/NMF'] = np.append(score_model(gmm_nmf_df),gmm_nmf_clusters).reshape(-1,1)
4 model_scores_df

```

Out[60]:

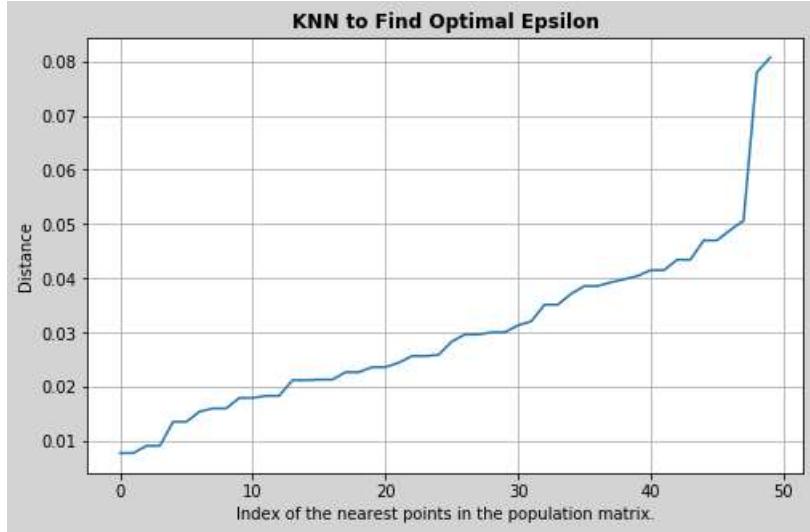
	Hierarchical LDA	Hierarchical PLSA/NMF	KMeans LDA	KMeans PLSA/NMF	GMM LDA	GMM PLSA/NMF
CH Score	63.920187	59.160354	65.412290	62.244816	64.913570	43.074089
Silhouette Score	0.464105	0.391698	0.451274	0.406453	0.436603	0.356773
# Clusters	4.000000	4.000000	4.000000	4.000000	4.000000	4.000000

DBSCAN Clustering

Perform Nearest Neighbors to Find Initial Value of Epsilon

Epsilon is then maximum distance between two samples for one to be considered as in the neighborhood of the other. For brevity, only search the LDA dataset as it has out performed the NMF thus far.

```
In [61]:  
1 db_lda_df = lda_sent_df.copy(deep=True)  
2 db_nmf_df = nmf_sent_df.copy(deep=True)  
3  
4 nn = NearestNeighbors(n_neighbors=6,algorithm='brute').fit(db_lda_df[['neg','neu','pos']])  
5 distances, idx = nn.kneighbors(db_lda_df[['neg','neu','pos']])  
6 distances = np.sort(distances, axis=0)  
7 distances = distances[:,1]  
8  
9 plt.figure(figsize=(8,5),facecolor="#d3d3d3")  
10 plt.title('KNN to Find Optimal Epsilon',fontweight='bold')  
11 plt.xlabel('Index of the nearest points in the population matrix.')  
12 plt.ylabel('Distance')  
13 plt.plot(distances)  
14 _ = plt.grid()
```



Grid Search to Fine Tune Epsilon and Minimum Samples

Minimumn samples is the minimum number of samples (or total weight) in a neighborhood for a point to be considered as a core point.

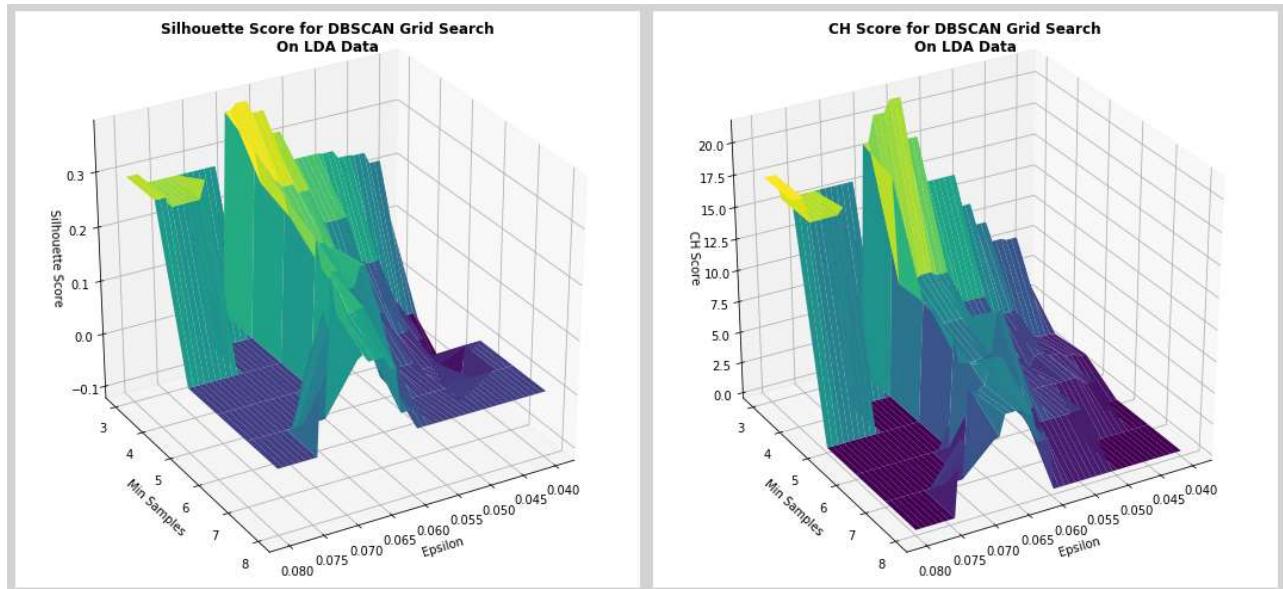
```
In [62]: 1  '''Setup Search Parameters'''
2  dbSCAN_eps = np.linspace(0.04, 0.08, 50)
3  dbSCAN_min_samp = np.arange(3,9)
4
5  EE, MS = np.meshgrid(dbSCAN_eps,dbSCAN_min_samp)
6  ee, ms = EE.flatten(), MS.flatten()
7  dbSCAN_no_clusters = np.zeros_like(ee)
8  dbSCAN_sil_scores = np.zeros_like(ee)
9  dbSCAN_ch_scores = np.zeros_like(ee)
10
11 # Build Small Test Set
12 np.random.seed(42)
13
14 '''Start Grid Search'''
15 print('Grid Search Sample Set Size: ',len(db_lda_df.index))
16 print('Beginning grid search on %d parameter combinations...' % (len(ee)))
17 results = []
18 t0 = time()
19 for i in range(len(ee)):
20     mdl = DBSCAN(eps=ee[i], min_samples=ms[i])
21     print('\nLoop %d, epsilon = %.3f, min samples = %d' % (i, ee[i], ms[i]))
22     mdl.fit(db_lda_df[['neg','neu','pos']])
23     dbSCAN_no_clusters[i] = len(set(mdl.labels_)) - (1 if -1 in mdl.labels_ else 0)
24     if dbSCAN_no_clusters[i] < 2:
25         continue
26     dbSCAN_sil_scores[i] = silhouette_score(db_lda_df[['neg','neu','pos']],
27                                            mdl.labels_,metric='euclidean')
28     dbSCAN_ch_scores[i] = calinski_harabasz_score(db_lda_df[['neg','neu','pos']],
29                                                mdl.labels_)
30
31 clear_output(wait=True)
32 print('Grid Search Sample Set Size: ',len(db_lda_df.index))
33 print('Beginning grid search on %d parameter combinations...' % (len(ee)))
34 print("Completed in %.3fs." % (time() - t0))
35
36 '''Display Results'''
37 df_db_grid = pd.DataFrame({'Epsilon':ee, 'Min Samples':ms,
38                            'No. Clusters':dbSCAN_no_clusters,
39                            'Silhouette Score':dbSCAN_sil_scores,
40                            'CH Score':dbSCAN_ch_scores})
41
42 best_sil = df_db_grid['Silhouette Score'].argmax()
43 best_ch = df_db_grid['CH Score'].argmax()
44 print('Best Silhouette Score:\n', df_db_grid.loc[best_sil])
45 print('Best CH Score:\n', df_db_grid.loc[best_ch])
46 done()
47
```

```
Grid Search Sample Set Size:  50
Beginning grid search on 300 parameter combinations...
Completed in 1.602s.
Best Silhouette Score:
    Epsilon          0.062041
    Min Samples      3.000000
    No. Clusters     3.000000
    Silhouette Score 0.382059
    CH Score        20.383818
Name: 27, dtype: float64
Best CH Score:
    Epsilon          0.059592
    Min Samples      3.000000
    No. Clusters     3.000000
    Silhouette Score 0.358113
    CH Score        21.293836
Name: 24, dtype: float64
```

0:04 / 0:04

Plot results of DBSCAN grid search

```
In [63]: 1 '''Plot Results of Grid Search'''
2 Z_sil = df_db_grid['Silhouette Score'].values.reshape(EE.shape)
3 Z_ch = df_db_grid['CH Score'].values.reshape(EE.shape)
4 fig = plt.figure(tight_layout=True, figsize=(15,7), facecolor='#d3d3d3')
5 titles = ['Silhouette Score for DBSCAN Grid Search\nOn LDA Data',
6           'CH Score for DBSCAN Grid Search\nOn LDA Data']
7 zlbl = ['Silhouette Score', 'CH Score']
8 for i,Z in enumerate([Z_sil,Z_ch]):
9     ax = fig.add_subplot(1,2,i+1, projection='3d')
10    ax.plot_surface(EE, MS, Z, cmap='viridis', edgecolor='none')
11    ax.set_xlabel('Epsilon')
12    ax.set_ylabel('Min Samples')
13    ax.set_zlabel(zlbl[i])
14    ax.set_title(titles[i], fontweight='bold')
15    ax.view_init(None, 60)
```



From the Grid Search, Best Min Samples is 3, and Best Epsilon is around 0.06

Perform DBSCAN Clustering and Plot Results

```
In [64]: 1 lda_dbSCAN = DBSCAN(eps=0.06, min_samples=3)
2 def dbSCAN_wrapper(df, msg):
3     # Perform DBSCAN
4     t0 = time()
5     print('Beginning DBSCAN Clustering on %s Data...' % (msg))
6     mdl = DBSCAN(eps=0.06, min_samples=3)
7     mdl.fit(df[['neg','neu','pos']])
8     # Label Topics
9     df['Label'] = mdl.labels_
10    print("completed in %0.3fs" % (time() - t0))
11    return mdl,df
12
13 db_nmf, db_nmf_df = dbSCAN_wrapper(db_nmf_df,'PLSA/NMF')
14 db_lda, db_lda_df = dbSCAN_wrapper(db_lda_df,'LDA')
15 done()
```

Beginning DBSCAN Clustering on PLSA/NMF Data...

completed in 0.004s

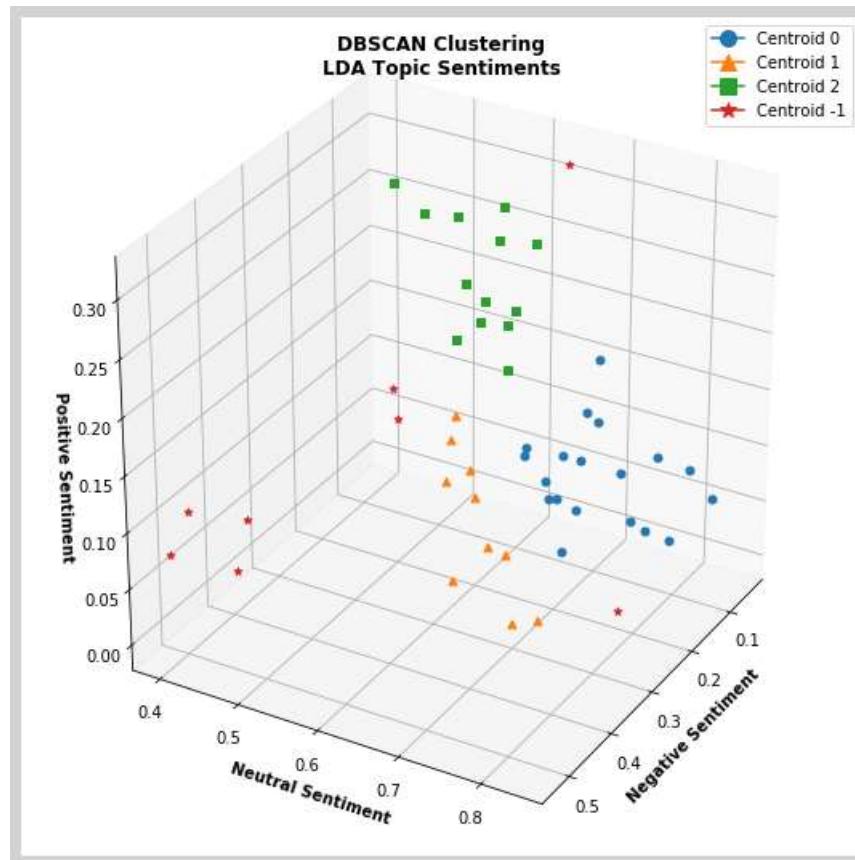
Beginning DBSCAN Clustering on LDA Data...

completed in 0.003s

0:04 / 0:04

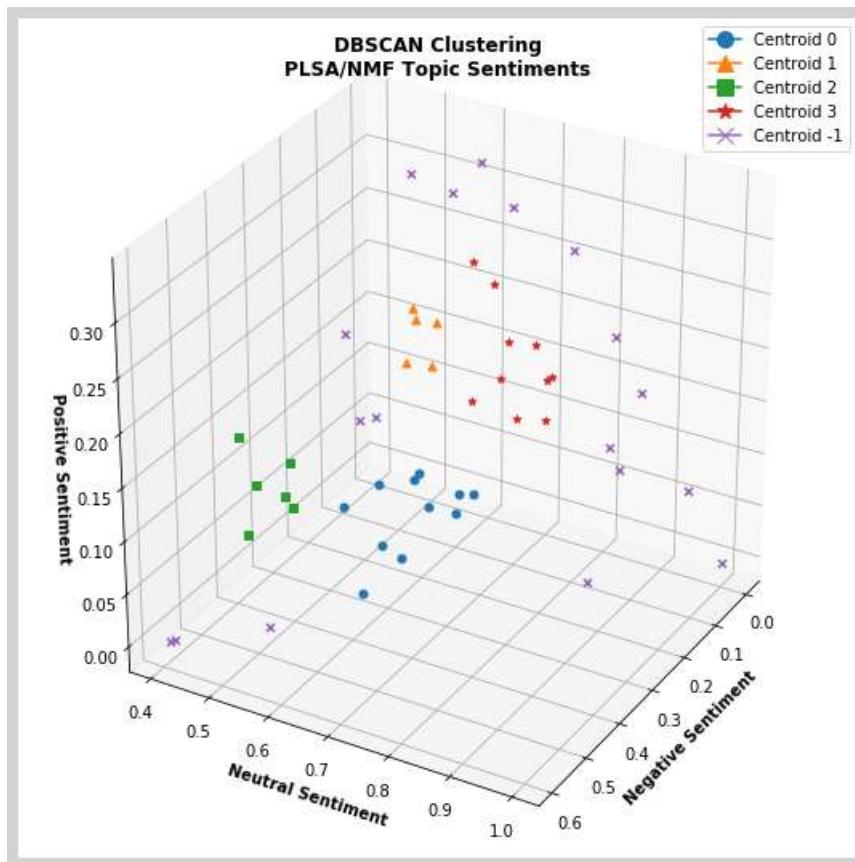
Plot DBSCAN Results

```
In [65]: 1 # Plot the results
2 ang = [(None,30)],[,(10,90),(100,0),(10,0)]
3 plot_clustering(None,db_lda_df,'DBSCAN Clustering\nnLDA Topic Sentiments', ang)
```



In [66]:

```
1 # Plot the results
2 ang = [(None,30)] #[, (10,90), (100,0), (10,0)]
3 plot_clustering(None,db_nmf_df,'DBSCAN Clustering\nnPLSA/NMF Topic Sentiments', ang)
```



Score DBSCAN Model Filtering out Noise Labels

In [67]:

```

1 # filter out noise and count clusters
2 lda_filter = db_lda_df.Label >= 0
3 db_lda_clusters = len(set(db_lda_df[lda_filter].Label))
4 nmf_filter = db_nmf_df.Label >= 0
5 db_nmf_clusters = len(set(db_nmf_df[nmf_filter].Label))
6
7 ds_len = len(db_lda_df)
8
9 print('DBSCAN Percent Noise for LDA: ', (ds_len - len(db_lda_df[lda_filter]))/ds_len*100)
10 print('DBSCAN Percent Noise for PLSA/NMF: ', (ds_len - len(db_nmf_df[nmf_filter]))/ds_len*100)
11
12
13 # Create DataFrame to Tabulate Model Scores
14 model_scores_df['DBSCAN LDA'] = np.append(score_model(db_lda_df[lda_filter]),db_lda_clusters).reshape(1,1)
15 model_scores_df['DBSCAN PLSA/NMF'] = np.append(score_model(db_nmf_df[nmf_filter]),db_nmf_clusters).reshape(1,1)
16 model_scores_df

```

DBSCAN Percent Noise for LDA: 16.0
DBSCAN Percent Noise for PLSA/NMF: 36.0

Out[67]:

	Hierarchical LDA	Hierarchical PLSA/NMF	KMeans LDA	KMeans PLSA/NMF	GMM LDA	GMM PLSA/NMF	DBSCAN LDA	DBSCAN PLSA/NMF
CH Score	63.920187	59.160354	65.412290	62.244816	64.913570	43.074089	49.645517	60.763639
Silhouette Score	0.464105	0.391698	0.451274	0.406453	0.436603	0.356773	0.456330	0.504851
# Clusters	4.000000	4.000000	4.000000	4.000000	4.000000	4.000000	3.000000	4.000000

Interpereting Results

The KMeans LDA model appears to be the best model. DBSCAN PLSA performed well too, but it classified 36% of the points as noise making it's results less desirable in context.

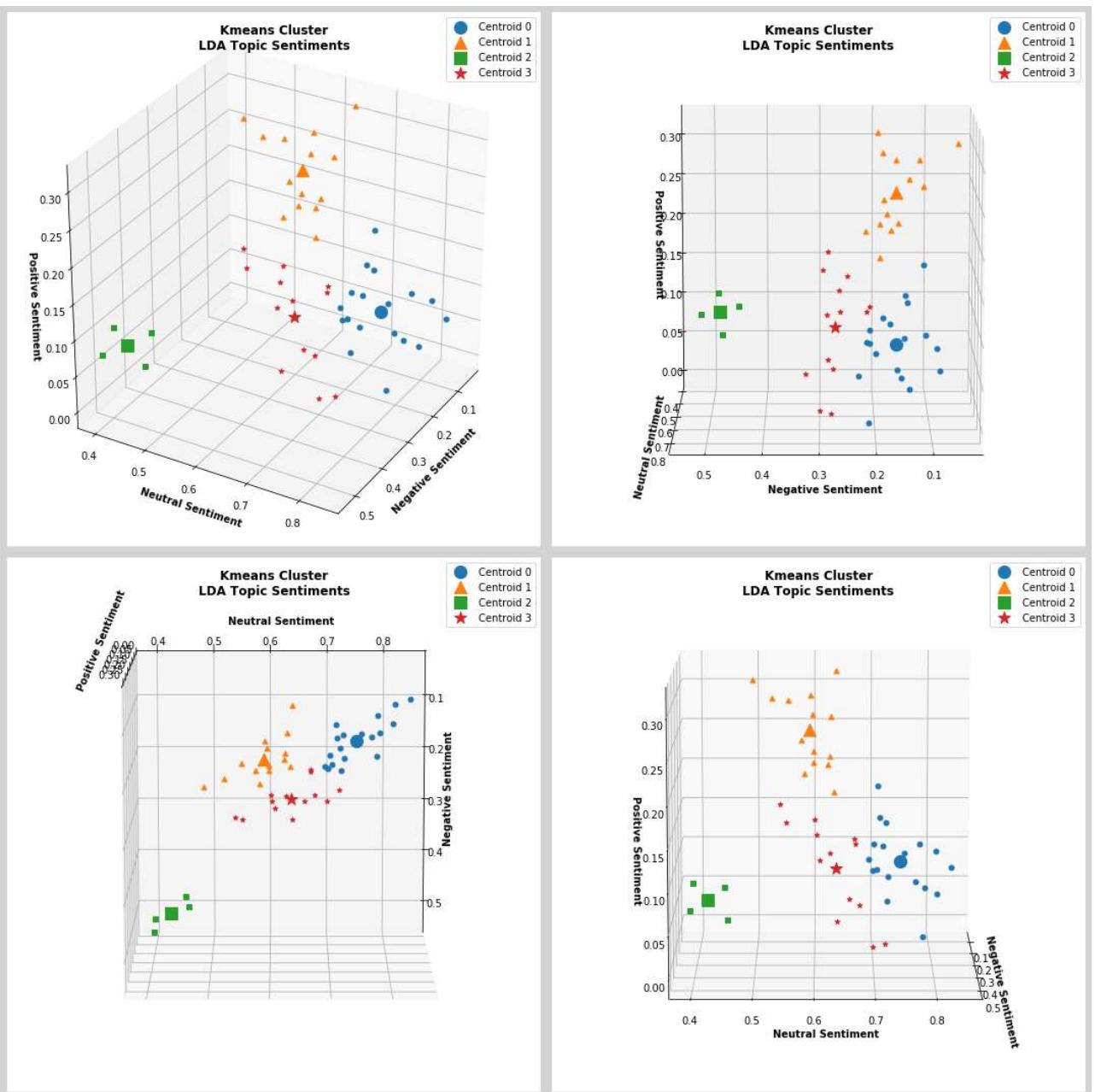
Plot Kmeans LDA Model From Several Angles to Better Inspect

In [68]:

```

1 # Plot the results
2 ang = [(None,30),(10,90),(100,0),(10,0)]
3 plot_clustering(km_lda.cluster_centers_,
4                  km_lda_df, 'Kmeans Cluster\nLDA Topic Sentiments', ang)

```



Process the Test Data by Scoring for Sentiment Only on The Features From the LDA Topic Model

Vectorize Test Data With LDA Topic Model Vocabulary

```
In [69]: 1 # Build A Vectorizer Based on Topic Model Vocabulary
2 lda_topic_vocab = {}
3 for i,token in enumerate(tf_feature_names):
4     # key must be feature name and value is index
5     lda_topic_vocab[token] = i
6
7 # Vectorize test data
8 tf_test_vectorizer = CountVectorizer(lowercase=False,
9                                     ngram_range=(1,2),
10                                    vocabulary=lda_topic_vocab)
11 tf_test = tf_test_vectorizer.fit_transform(totalDF.loc[idx_test].Cleaned)
12 tf_test_names = tf_test_vectorizer.get_feature_names()
13
14 # Check that vocabularies are the same
15 for i,key in enumerate(lda_topic_vocab.keys()):
16     if tf_test_names[i] != key:
17         print('\n'*25)
18         print('Index %d vocab not matching' % i)
19         print('LDA Vocab: {}'.format(tf_test_names[i]))
20         print('Test Vect Vocab: {}'.format(key))
```

Inspect the Test Data

```
In [70]: 1 def display_vect_words(vect,names,n_display):
2     for i in range(n_display):
3         print("Test Sample {:d} =====".format(i))
4         idx = vect[i,:].nonzero()[1]
5         print("|".join([names[k] for k in idx]))
6
7 display_vect_words(tf_test,tf_test_names,10)
```

```
Test Sample 0 =====
hand|korea|park|south|south korea|verdict
Test Sample 1 =====
australia|save
Test Sample 2 =====
suspect|theft
Test Sample 3 =====
bomber|return|uk
Test Sample 4 =====
saturday|wrap
Test Sample 5 =====
announced|asbestos|panel|response
Test Sample 6 =====
drought|music|relief
Test Sample 7 =====
join|olympics|woman
Test Sample 8 =====
could|model|water
Test Sample 9 =====
adelaide|birthday|festival
```

```
In [71]: 1 # Build Test DataFrame with Sentiment Score for Each Entry
2 def sent_score_test_data(tf_matrix, tf_names, sentAnalyzer):
3     t0 = time()
4     print('Beginning Sentiment Scoring on %d entries...' % tf_matrix.shape[0])
5     sent_list = []
6     for i in range(tf_matrix.shape[0]):
7         idx = np.array(list(tf_test[i,:].nonzero()[1])), dtype = 'int32')
8         quant = [tf_matrix[i,j] for j in idx]
9         words = [tf_names[j] for j in idx]
10        s = []
11        for k in range(len(quant)):
12            for l in range(quant[k]):
13                s = s + [words[k]]
14        sent_list.append(sentAnalyzer.polarity_scores(' '.join(s)))
15    print("Completed in %0.3fs." % (time() - t0))
16    return pd.DataFrame(sent_list)
17
```

```
In [72]: 1  '''Score the Test Data'''
2  vsia = vader.SentimentIntensityAnalyzer()
3  testDF = sent_score_test_data(tf_test, tf_test_names, vsia)
4  display(testDF.head())
5  done()
```

Beginning Sentiment Scoring on 314907 entries...
Completed in 97.210s.

	compound	neg	neu	pos
0	0.5859	0.000	0.510	0.490
1	0.4939	0.000	0.238	0.762
2	-0.2960	0.688	0.312	0.000
3	0.0000	0.000	1.000	0.000
4	0.0000	0.000	1.000	0.000

0:04 / 0:04

Test KMeans Model by Bootstrapping CH and Silhouette Scores

```
In [73]: 1  def boot_bootstrap_results(test_df, mdl, n_folds = 2000 ,n_samples = 50):
2      t0 = time()
3      print('Beginning testing clustering model with %d samples and %d folds'\
4          % (n_samples,n_folds))
5      print('With model {}'.format(mdl))
6      results = []
7      for k in range(n_folds):
8          df = test_df.sample(n=n_samples, replace = False,
9              random_state = k)
10         df['Label'] = mdl.predict(df[['neg','neu','pos']])
11         if len(set(df.Label)) > 1:
12             ch,sil = score_model(df)
13             results.append({'CH Score':ch, 'Silhouette Score':sil})
14     print("Completed in %.3fs." % (time() - t0))
15  return pd.DataFrame(results)
```

In [74]:

```
1 test_results_df = boot_strap_results(testDF,km_lda)
2 display(test_results_df.head(10))
3 display(test_results_df.describe())
4 done()
```

Beginning testing clustering model with 50 samples and 2000 folds
With model KMeans(algorithm='full', n_clusters=4, n_init=1000, random_state=42, tol=1e-08)
Completed in 30.304s.

	CH Score	Silhouette Score
0	109.253418	0.741201
1	71.549669	0.662645
2	107.332098	0.646145
3	61.531661	0.364559
4	107.650287	0.666073
5	67.004170	0.516775
6	104.068073	0.672356
7	85.222724	0.676210
8	48.699477	0.550615
9	69.871843	0.639992

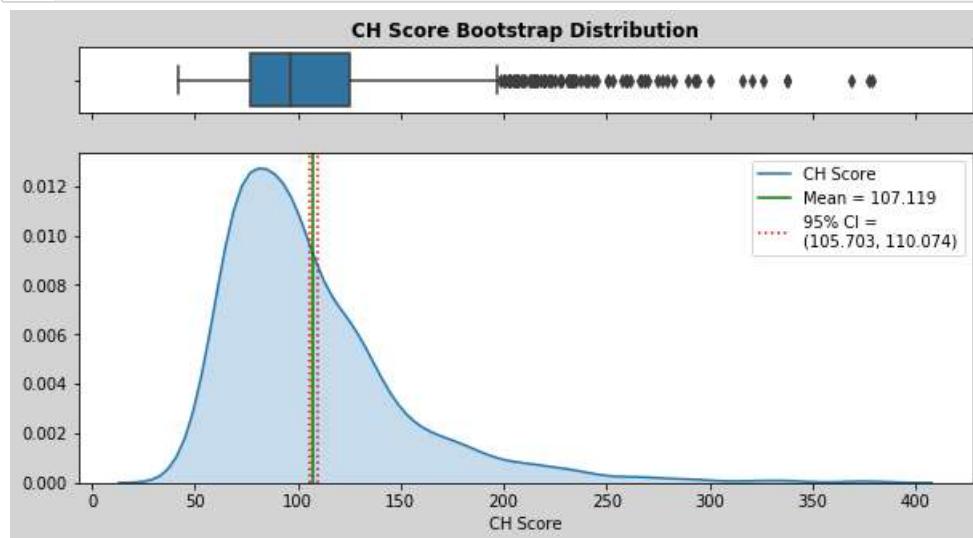
	CH Score	Silhouette Score
count	2000.000000	2000.000000
mean	107.119269	0.659873
std	44.175669	0.076415
min	42.242472	0.364559
25%	76.818871	0.608043
50%	96.253290	0.660825
75%	125.180427	0.711906
max	379.050637	0.893939

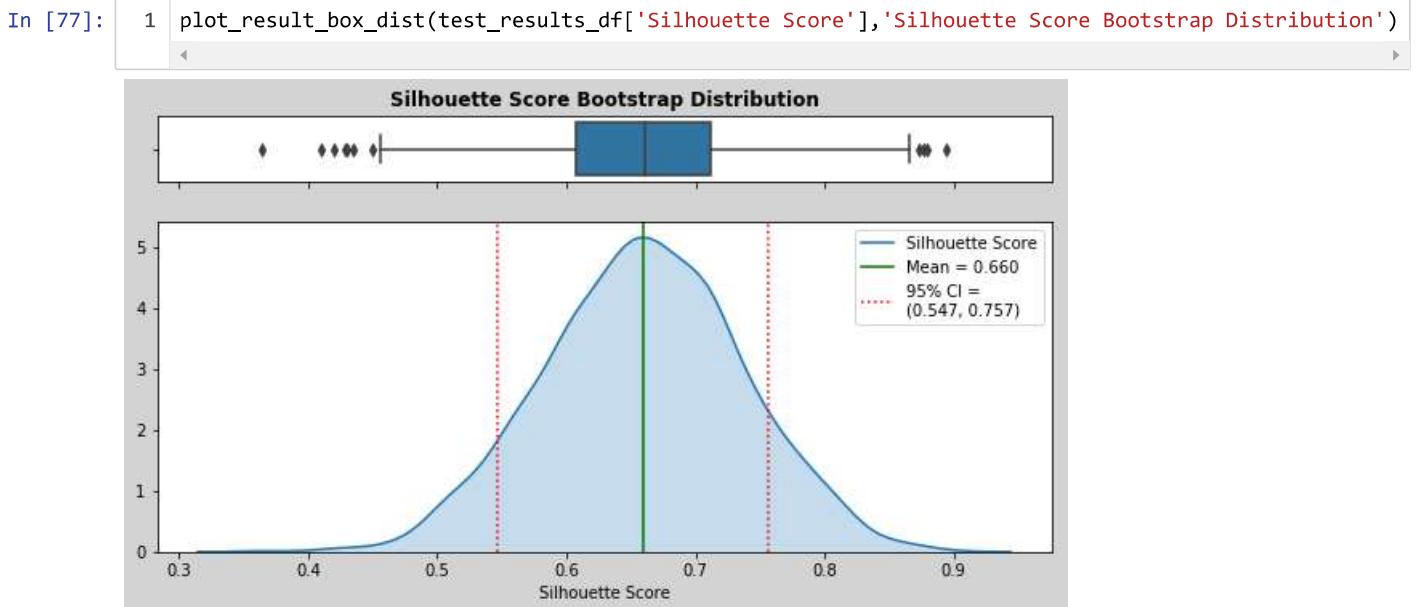
0:04 / 0:04

```
In [75]: 1 def plot_result_box_dist(data, title):
2     f, (ax_box, ax_hist) = plt.subplots(2, sharex=True,
3                                         facecolor='#d3d3d3',
4                                         figsize=(10,5),
5                                         gridspec_kw= {"height_ratios": (0.2, 1)})
6     z95 = 1.96      # zscore for 95%
7     N = len(data)
8     d = round(N*0.95/2)      # get middle 95% of values
9
10    # Get Bootstrapped Confidence Interval and Add to Plot
11    ci_values = data.sort_values()
12    ci = (ci_values[0+d],ci_values[N-1-d])      # store 95% confidence interval
13    md = ci_values.mean()
14
15    # plot box plot
16    sns.boxplot(data, ax=ax_box)
17    # plot distribution
18    sns.kdeplot(data, shade=True, ax=ax_hist)
19    ax_hist.axvline(md,0,10,color='g',linestyle='solid',label='Mean = {:.3f}'.format(md))
20    ax_hist.axvline(ci[0],0,10,color='r',linestyle='dotted')
21    ax_hist.axvline(ci[1],0,10,color='r',linestyle='dotted',
22                     label='95% CI =\n{:0.3f}, {:0.3f}'.format(*ci))
23    ax_hist.set(xlabel=data.name)
24    ax_hist.legend(loc='best')
25    ax_box.set(xlabel='')
26    ax_box.set_title(title,fontweight='bold')
```

Plot Results of Bootstrapping

```
In [76]: 1 '''Plot Results'''
2 plot_result_box_dist(test_results_df['CH Score'],'CH Score Bootstrap Distribution')
```





Compare Other LDA Models Against Kmeans LDA as "Ground Truth"

Find Closest Point to Kmeans Center for Each Cluster

In [78]:

```
1 # Align Labels Between Models to Kmeans Nomenclature
2 conf_df = lda_sent_df.copy(deep=True)
3
4 # Find Closest Point to Each Cluster Center and Use to Find Matching Label in
5 #...other models
6 label_dict = {}
7 for center in km_lda.cluster_centers_:
8     closest_pnt = cosine_similarity(center.reshape(1,-1),conf_df[['neg','neu','pos']]).argmax()
9     label_dict[km_lda.predict(center.reshape(1,-1))[0]] = [center,closest_pnt]
10 print('Kmeans Label: Kmeans Center, Closest Data Point to Center')
11 display(label_dict)
12
13
```

Kmeans Label: Kmeans Center, Closest Data Point to Center

```
{0: [array([0.16377778, 0.74905556, 0.08705556]), 26],
 1: [array([0.15964286, 0.59057143, 0.24985714]), 21],
 2: [array([0.4905 , 0.42675, 0.08325]), 8],
 3: [array([0.27185714, 0.63628571, 0.09171429]), 47]}
```

Inspect Topics Closest to Each Center

In [95]:

```

1 def display_topic(model, comp_index, label, feature_names, no_top_words):
2     print("Cluster {:d}, Centroid Topic {:d} =====".format(label,comp_index))
3     idx = (-model.components_[comp_index,:]).argsort()[:no_top_words]
4     print("|".join([feature_names[k] for k in idx]))
5     print()
6
7
8
9 for key in label_dict.keys():
10     display_topic(lda, label_dict[key][-1],key, tf_feature_names, 10)
11

```

Cluster 0, Centroid Topic 26 =====
plan|australia|job|fear|test|nt|leader|mine|opposition|river

Cluster 1, Centroid Topic 21 =====
election|help|law|act|air|keep|light|syria|northern|problem

Cluster 2, Centroid Topic 8 =====
woman|accused|dead|melbourne|lead|sex|study|chinese|spark|teacher

Cluster 3, Centroid Topic 47 =====
china|case|centre|move|mayor|korea|warns|uk|stand|north korea

Align the Labeling Between Algorithms for Comparison

In [79]:

```

1 # Build Master DataFrame with Labels from All Models
2 cols = ['Kmeans Label', 'Heirarchical Label', 'GMM Label', 'DBSCAN Label']
3 for i,mdl_df in enumerate([km_lda_df, hc_lda_df, gmm_lda_df, db_lda_df]):
4     conf_df[cols[i]] = mdl_df['Label']
5
6 closest_points_to_centers = [label_dict[key][1] for key in label_dict.keys()]
7 rosetta_df = conf_df.loc[closest_points_to_centers][cols].set_index(cols[0])
8 print('Label Rosetta Stone ' + '='*30)
9 display(rosetta_df)

```

Label Rosetta Stone =====

Heirical Label GMM Label DBSCAN Label

Kmeans Label

	Heirical Label	GMM Label	DBSCAN Label
0	0	2	0
1	2	0	2
2	3	1	-1
3	1	3	1

DBSCAN Noise Label Changed to Cluster 3 to align with label space of the other 3 models

```
In [80]: 1  '''Fix Labeling'''
2  adj_lbls = ['KMeans Adj.', 'Heirarchical Adj.', 'GMM Adj.', 'DBSCAN Adj.']
3  conf_df[adj_lbls[0]] = conf_df[cols[0]] # copy over kmeans Labels
4
5  # change -1 label to 3 for DBSCAN
6  conf_df[cols[-1]][conf_df[cols[-1]] < 0] = 3
7  rosetta_df.loc[2,cols[-1]] = 3
8
9  for i in range(1,len(cols),1):
10    # build converter function
11    converter = rosetta_df[cols[i]].values.argsort() # array of indices that would sort the labels
12    fconvert = lambda x:converter[x]
13    conf_df[adj_lbls[i]] = conf_df[cols[i]].apply(fconvert)
14    rosetta_df[adj_lbls[i]] = rosetta_df[cols[i]].apply(fconvert)
15
16 display(rosetta_df)
17 display(conf_df.loc[closest_points_to_centers])
```

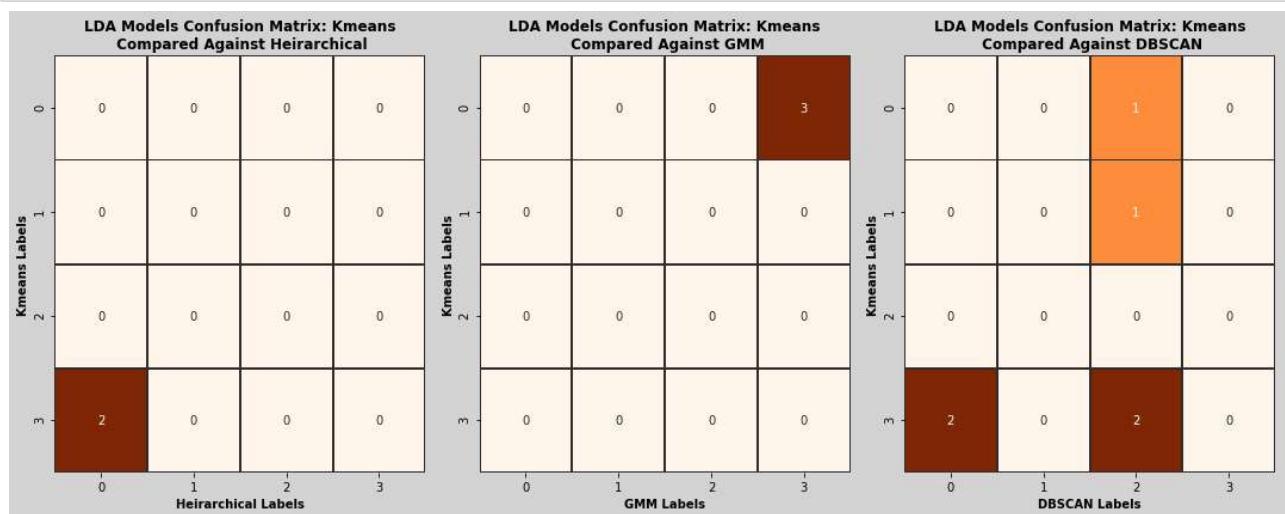
	Heirarchical Label	GMM Label	DBSCAN Label	Heirarchical Adj.	GMM Adj.	DBSCAN Adj.
Kmeans Label	0	0	2	0	0	0
0	0	2	0	1	1	1
1	2	0	2	1	1	1
2	3	1	3	2	2	2
3	1	3	1	3	3	3

	compound	neg	neu	pos	Kmeans Label	Heirarchical Label	GMM Label	DBSCAN Label	KMeans Adj.	Heirarchical Adj.	GMM Adj.	DBSCAN Adj.
26	-0.4939	0.149	0.756	0.095	0	0	2	0	0	0	0	0
21	0.5106	0.183	0.577	0.239	1	2	0	2	1	1	1	1
8	-0.9786	0.495	0.401	0.104	2	3	1	3	2	2	2	2
47	-0.8555	0.272	0.637	0.091	3	1	3	1	3	3	3	3

In []:

```
In [81]: 1  '''Display Confusion Matrix'''
2  def add_confusion_matrix(df, ref_col, compare_col, ax):
3      conf_mx = confusion_matrix(df[ref_col],df[compare_col])
4      np.fill_diagonal(conf_mx,0)
5      sns.heatmap(conf_mx, annot=True, ax=ax, cmap='Oranges',
6                  linewidths = 1, linecolor='#2b2d2f', cbar=False)
7      r = ref_col.split(' ')[0]
8      c = compare_col.split(' ')[0]
9      ax.set_ylabel(r + ' Labels', fontweight='bold')
10     ax.set_xlabel(c + ' Labels', fontweight = 'bold')
11     _= ax.set_title('LDA Models Confusion Matrix: {} \n Compared Against {}'.format(r,c),
12                      fontweight='bold')
13
14 return None
```

```
In [82]: 1 fig, axes = plt.subplots(1,3, facecolor='#d3d3d3', figsize=(15,6),  
2                      tight_layout = True)  
3 axes = axes.ravel()  
4 for i,ax in enumerate(axes):  
5     add_confusion_matrix(conf_df, cols[0],adj_lbls[i+1],ax)
```



In [83]:

```
1 #Show Points In Disagreement
2 miss_labeled = []
3 for lbl in adj_lbls:
4     if not('kmeans' in lbl.lower()):
5         print('='*40)
6         s = [adj_lbls[0],lbl]
7         miss_labeled.append(conf_df[conf_df[lbl] != conf_df[adj_lbls[0]]][s])
8         display(miss_labeled[-1])
9
=====

```

KMeans Adj.	Heirarchical Adj.
0	3
48	0

	KMeans Adj.	Heirarchical Adj.
0	3	0
48	0	0

KMeans Adj.	GMM Adj.
10	0
13	0
49	0

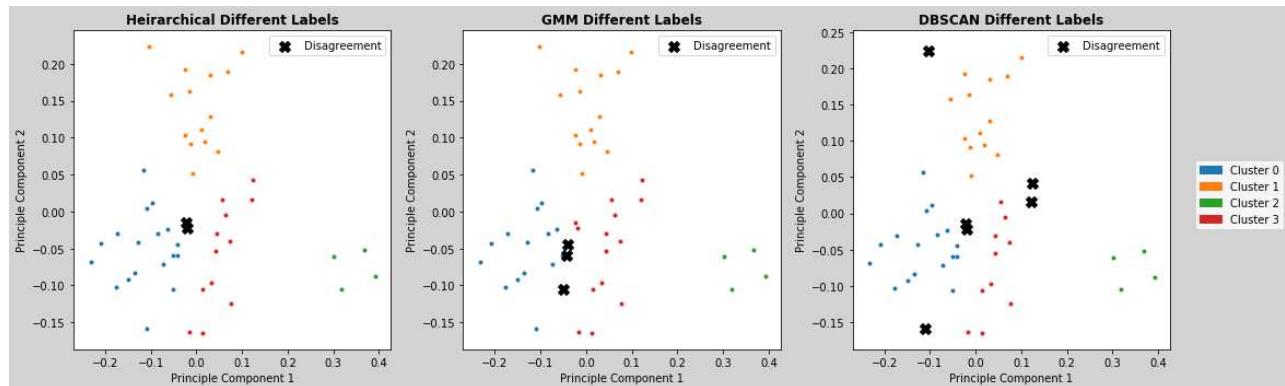
	KMeans Adj.	GMM Adj.
10	0	3
13	0	3
49	0	3

KMeans Adj.	DBSCAN Adj.
0	3
12	1
19	3
36	0
43	3
48	3

	KMeans Adj.	DBSCAN Adj.
0	3	0
12	1	2
19	3	2
36	0	2
43	3	2
48	3	0

Perfrom PCA and Graph Points in Disagreement

```
In [88]: 1 pca = PCA(n_components=2)
2 X = pca.fit_transform(conf_df[['neg','neu','pos']])
3
4 # build custom legend
5 le = []
6 for i in range(4):
7     le.append(Patch(facecolor=COLORS[i],
8                      edgecolor=COLORS[i],
9                      label='Cluster {:d}'.format(i)))
10
11 fig, axes = plt.subplots(nrows=1,ncols=3,tight_layout=True,
12                         facecolor = '#d3d3d3',figsize=(15,5))
13 axes = axes.ravel()
14 for j,ml in enumerate(miss_labeled):
15     for i in range(4):
16         idx = conf_df['KMeans Adj.'] == i
17         axes[j].scatter(X[idx,0],X[idx,1],marker='.',color = COLORS[i])
18         idx = ml.index
19         s = ml.columns[-1].split(' ')[0] + ' Different Labels'
20         axes[j].scatter(X[idx,0],X[idx,1],marker='X',color = 'k',label='Disagreement',s = 100)
21         axes[j].set_title(s,fontweight='bold')
22         axes[j].set_ylabel('Principle Component 2')
23         axes[j].set_xlabel('Principle Component 1')
24         _ = axes[j].legend(loc='upper right')
25
26 fig.subplots_adjust(top=0.9, left=0.1, right=0.9, bottom=0.12)
27 fig.legend(handles=le,loc='upper left',bbox_to_anchor = (1,.6),ncol=1)
28 fig.tight_layout()
```



In []:

1