

NOTE:

1. ALL THE ANALYSIS PERFORMED BY USING 100K REVIEWS

Steps to build co-occurrence matrix:(Overview)

1. pick top features based on their IDF values and create corpus of top features.
2. create a matrix of shape = (len(corpus), len(corpus))and initialize it with zeros .
3. now select a review and tokenize it.
4. now loop over all the columns and rows of matrix created in step-2 and update it for a review.
5. updation of a matrix:
 - a.) pick two words from corpus say word1 and word2:
 - [a.1] if word1 and word2 are equal and present in a review add zero to this cell.
 - [a.2] if word1!=word2 and present in a review , consider word1 as base word and find how much time word
2
occur in forward and backward window of base word=word1, and add count(word2) to this cell.

Load required libraries

```
In [1]: 1 import warnings
        2 warnings.filterwarnings("ignore")
```

```
In [26]: 1 import pickle
2 import seaborn as sns
3 from sklearn.decomposition import TruncatedSVD
4 import numpy as np
5 import pandas as pd
6 import matplotlib.pyplot as plt
7 import multiprocessing as mp
8 from scipy.sparse import csr_matrix, find
9 from sklearn.cluster import KMeans
10 from os import path
11 import os
12 import requests
13 from PIL import Image
14 from wordcloud import WordCloud
15 from tqdm import tqdm
16 from sklearn.metrics.pairwise import cosine_similarity
17 %matplotlib inline
```

Load required data:

```
In [3]: 1 #Functions to save objects for later use and retrieve it
2 def savetofile(obj, filename):
3     pickle.dump(obj, open(filename+".pkl", "wb"))
4 def openfromfile(filename):
5     temp = pickle.load(open(filename+".pkl", "rb"))
6     return temp
7
8 X=openfromfile('X')
9 X_tfidf=openfromfile('X_tfidf')
10 tf_idf_vect=openfromfile('tf_idf_vect')
11 #corpus=openfromfile('corpus')
12 #coo_all=openfromfile('coo_all')
13 #df=openfromfile('df')
```

Truncated-SVD

Function for constructing Co-Occurrence matrix:

```
In [4]: ▶ 1 def coo_mat(list_sent_word,corpus>window_len):
2         '''CO-OCCURENCE MATRIX CONSTRUCTION FOR GIVEN CORPUS AND WINDOW SIZE GIVEN'''
3         coo=np.zeros((len(corpus),len(corpus)))
4         i=j=0
5         for k in tqdm(range(len(list_sent_word))):
6             j=i=0
7             for word1 in corpus:
8                 if word1 in list_sent_word[k]:
9                     j=0
10                    for word2 in corpus:
11                        count=0
12                        if word1!=word2 and word2 in list_sent_word[k]:
13                            #FIND THE INDEX OF BASE WORD FROM WHICH WE CALCULATE WINDOW SIZE
14                            base_index=[ind for ind,w in enumerate(list_sent_word[k]) if word1==w]
15                            for index in base_index:
16                                #LIST TO STORE WORDS IN LEFT AND RIGHT SIDE OF BASE WORD WITHIN WINDOW DISTANCE
17                                list_word_window=[]
18                                if index+1<len(list_sent_word[k]):
19                                    list_word_window.extend(list_sent_word[k][max(index-window_len,0) : index] +\
20                                                            list_sent_word[k][index+1 : index + window_len+1])
21                                elif index==len(list_sent_word[k]):
22                                    list_word_window.extend(list_sent_word[k][index-window_len: index])
23                                #NO. OF TIME A WORD IN COLUMN OCCURS IN LEF AND RIGHT WINDOW OF BASE WORD
24                                count += list_word_window.count(word2)
25                                coo[i][j]+=count
26                                #IF BOTH WORDS ARE SAME THEN WE ADD 0 TO THAT CELL
27                            else:
28                                coo[i][j]+=0
29                                j+=1
30                            #i+=1
31                            #IF WORD FROM SENTENCE IS NOT IN ROW OF CO-OCCURENCE MATRIX THEN WE ADD 0 TO THAT ROW
32                        else:
33                            coo[i][:]+=0
34                            #i+=1
35                    i+=1
36            return coo
37
```

Function to draw co-occurrence matrix using heatmap:

```
In [5]: 1 def visualizeCooMat(coo_mat, corpus, dim=10):
2     corpus_row=corpus[:dim]
3     coo_mat_df=pd.DataFrame(coo_mat[:dim,:dim], columns=corpus_row, index=corpus_row)
4     plt.figure(1, figsize=(18,12))
5     sns.set(font_scale=1.2)
6     cmat=sns.heatmap(coo_mat_df, annot=True, cmap='RdBu', cbar=False)
7     cmat.set_yticklabels(corpus_row, rotation=0)
8     plt.title('Part of Co-occurrence matrix of dimension(%d * %d)'%(dim,dim), size=20)
9     plt.show()
```

Function for applying TSVD nd find the optimal no. of feature which can explained max variance:

```
In [6]: 1 def OptimalComponents_Mat(co_occ_mat, opt_component, transformed_mat=False):
2     '''FUNCTION TO FIND OPTIMA NO. OF COMPONENTS BASED ON ELBOW CURVE'''
3     t_svd = TruncatedSVD(n_components=opt_component)
4     t_svd_data = t_svd.fit_transform(co_occ_mat)
5     #IF TRANSFORMED_MAT =FALSE DRAW ELBOW CURVE ONLY
6     if transformed_mat==False:
7         cum_var_explained = np.cumsum(t_svd.explained_variance_ratio_)
8         plt.figure(1, figsize=(10, 6))
9         plt.clf()
10        plt.plot(cum_var_explained, linewidth=2)
11        plt.axis('tight')
12        plt.grid(True)
13        plt.xlabel('No. of Components')
14        plt.ylabel('Cumulative Explained Variance')
15        plt.title("Total-Variance-Explained VS No. of Components")
16        plt.show()
17    #IF TRANSFORMED_MAT=TRUE RETURNED TRANSFORMED MATRIX
18    else:
19        return t_svd_data
```

Function for finding top important features based on IDF values:

```

In [7]: 1 # Code From -> https://buhrmann.github.io/tfidf-analysis.html
2 def get_top_n_features(X_tfidf, features, top_n):
3     '''FUNCTION FOR FINDING TOP FEATURE'''
4     column_mean = np.mean(X_tfidf, axis = 0)
5     column_mean = np.array(column_mean)[0].tolist()
6     topn_ids = np.argsort(column_mean)[::-1][:top_n]
7     top_feats = [(features[i], column_mean[i]) for i in topn_ids]
8     df = pd.DataFrame(top_feats)
9     df.columns = ['feature', 'tfidf']
10    return df

```

Function for K-Means hyperparameter tuning using elbow method :

```

In [8]: 1 def clustering_model(data,params):
2     '''HYPERPARAM TUNNING AND DRAW ELBOW CURVE'''
3     inertia_value=[]
4     for k in tqdm(params['clusters']):
5         model=KMeans(n_clusters=k,init='k-means++',n_jobs=-1)
6         model.fit(data)
7         inertia_value.append(model.inertia_)
8     plt.figure(1,figsize=(10,6))
9     plt.plot(params['clusters'],inertia_value)
10    plt.title('ELBOW METHOD (OPTIMAL HYPERPARAM)')
11    plt.xlabel('K(No. of clusters): Hyperparam')
12    plt.ylabel('Inertia(Sum of intra cluster distance)')
13    plt.grid(True)
14    plt.show()
15

```

Function for labelling each value of a corpus to cluster label :

```

In [9]: 1 def corpus_labelling(data,corpus,params):
2         '''FUNCTION FOR LABELCORPUS WITH CLUSTER LABEL AND DRAW BAR PLOT OF NO. OF WORDS DISTRIBUTED TO EACH CLUSTER
3         model=KMeans(n_clusters=params,precompute_distances=True,init='k-means++',n_jobs=-1)
4         model.fit(data)
5         labels=model.labels_
6
7         cluster_label=np.c_[ corpus, labels ]
8         cluster_label_df=pd.DataFrame(cluster_label,columns=['word','label'])
9         print(cluster_label_df['label'].value_counts())
10        list_samples=[]
11        g=cluster_label_df.groupby(['label'])
12        for key in g.groups:
13            list_samples.append(len(g.groups[key]))
14        #BAR CHART
15        if len(list_samples)!=0:
16            plt.figure(1,figsize=(10,6))
17            sns.set(rc={'figure.figsize':(11.7,8.27)})
18            plt.title("NO. of Words per cluster")
19            # ADD BARS
20            plt.bar(range(len(list_samples)), list_samples) #plt.bar(range(params), list_samples)
21            # ADD LABELS
22            plt.xticks(range(len(list_samples)))
23            plt.xlabel('No. of Cluster')
24            plt.ylabel('No. of Words')
25            plt.show()
26        else:
27            print('All points are declared as noisy points ')
28        cluster_label_df=cluster_label_df.astype({'label':int})
29        return model,cluster_label_df

```

Function for wordcloud:

```

In [24]: 1 def wordcloud_each_cluster(optimal_cluster, cluster_label_df):
2         '''FUNCTION TO PLOT WORDS BELONGS TO A CLUSTER IN A WORDCLOUD'''
3         words_cluster={}
4         #PROPERTIES OF TITLE
5         title_font = {'size':'20', 'color':'black'}
6         #DOWNLOAD AND PUT MASKABLE IMAGE TO CWD
7         if not os.path.isfile('mask-cloud.png'):
8             url='http://www.shapecollage.com/shapes/mask-cloud.png'
9             r=requests.get(url)
10            with open('mask-cloud.png','wb') as f:
11                f.write(r.content)
12            #SHAPE WORDCLOUD ACCORDING TO MASKABLE IMAGE
13            mask_image_path = path.dirname(__file__) if "__file__" in locals() else os.getcwd()
14            mask = np.array(Image.open(path.join(mask_image_path, "mask-cloud.png")))
15            stopwords=set([])
16            for i in range(optimal_cluster):
17                #INITIALIZE WORDCLOUD OBJECT
18                wordcloud = WordCloud(max_font_size=120,\
19                                     stopwords=stopwords,\
20                                     colormap='winter',\
21                                     max_words=100,\
22                                     collocations=False,\
23                                     relative_scaling=1.0,\
24                                     mask=mask,\
25                                     background_color='white',\
26                                     contour_width=2,\
27                                     contour_color='royalblue')
28                #PICK WORDS BELONGS TO EACH CLUSTER
29                t=''
30                for token in list(cluster_label_df.loc[cluster_label_df['label']==i,\
31                                                         ['word']][ 'word']):
32                    t=t+' '+token
33                words_cluster[i]=t.strip().split()
34                wordcloud.generate(t)
35                #FIGURE SIZE ACCORDING TO NO. OF WORDS
36                if len(t.split())<=5:
37                    plt.figure(1,figsize=(5,5))
38                elif len(t.split())>=5:
39                    plt.figure(1,figsize=(14,13))
40                plt.title("WORDS BELONGS TO CLUSTER-%d"%i, **title_font )
41                plt.imshow(wordcloud, interpolation="bilinear")

```

```

42     plt.axis("off")
43     plt.show()
44     return words_cluster

```

Function which accept a word as input and return similar words to that word:

```

In [11]: ▶ 1 def similarWords(coo_all, corpus, word, words_in_cluster):
2     '''FUNCTION FOR CALCULATING SIMILARITY BETWEEN WORDS AND RETURN SIMILAR WORDS FROM CLUSTER'''
3     #DICTIONARY TO STORE WORDS AND ITS SIMILARITY WITH OTHER WORDS BELONGS TO PARTICULAR CLUSTER
4     similar_words={}
5     for key in words_in_cluster.keys():
6         if word in words_in_cluster[key]:
7             #FIND THE INDEX OF WORD IN CORPUS
8             word_index=corpus.index(word)
9             #WORD VECTOR OF A WORD GIVEN
10            given_word_vector=coo_all[word_index]
11            break
12
13    for w in words_in_cluster[key]:
14        #FIND INDEX OF WORD IN CORPUS
15        index=corpus.index(w)
16        #COSINE SIMILARITY OF A WORD WITH OTHER WORDS IN THAT CLUSTER
17        similar_words[w]=float(cosine_similarity(coo_all[index].reshape(1, -1),given_word_vector.reshape(1, -1))
18        similarity_df=pd.DataFrame(similar_words,index=[word]).round(3)
19    return similarity_df
20

```

Applying TruncateSVD :

[1.1] Taking top features from TFIDF


```
In [12]: 1 no_of_feat=2000
2 df=get_top_n_features(X_tfidf,tf_idf_vect.get_feature_names(),no_of_feat)
3 print(df.head(5))
4 df.shape
```

| | feature | tfidf |
|---|---------|----------|
| 0 | not | 0.044046 |
| 1 | like | 0.026099 |
| 2 | great | 0.025042 |
| 3 | good | 0.025020 |
| 4 | coffee | 0.024549 |

Out[12]: (2000, 2)

[1.2] Calulation of Co-occurrence matrix

[1.2.1]prepare input for creating C-occurence matrix:

```
In [13]: 1 #WINDOW LENGTH
2 window_len=5
3 #CORPUS/WORDS FOR WHICH WE HAVE TO CREATE CO-OCCURENCE MATRIX
4 corpus=list(df['feature'].values)
5 print('no. of words in a corpus:',len(corpus))
6 #LIST OF LIST OF WORDS IN A SENTENCE FOR ALL REVIEWS
7 list_sent_word=[]
8 for sent in X:
9     list_sent_word.append(sent.split())
```

no. of words in a corpus: 2000

[1.2.2] Create co-occurence matrix:

```
In [14]: ▶ 1 pool = mp.Pool(mp.cpu_count())
2 # Step 2: `pool.apply` the `howmany_within_range()`
3 %time coo_all = pool.apply(coo_mat, args=(list_sent_word, corpus, window_len))
4 # Step 3: Don't forget to close
5 pool.close()
6 #SAVE CO-OCCURENCE MATRIX TO PICKLE FILES
7 savetofile(coo_all, 'coo_all')
8 #CONVERT CO-OCCURENCE MATRIX TO SPARSE CO-OCCURENCE MATRIX
9 coo_matrix=csr_matrix(coo_all)
```

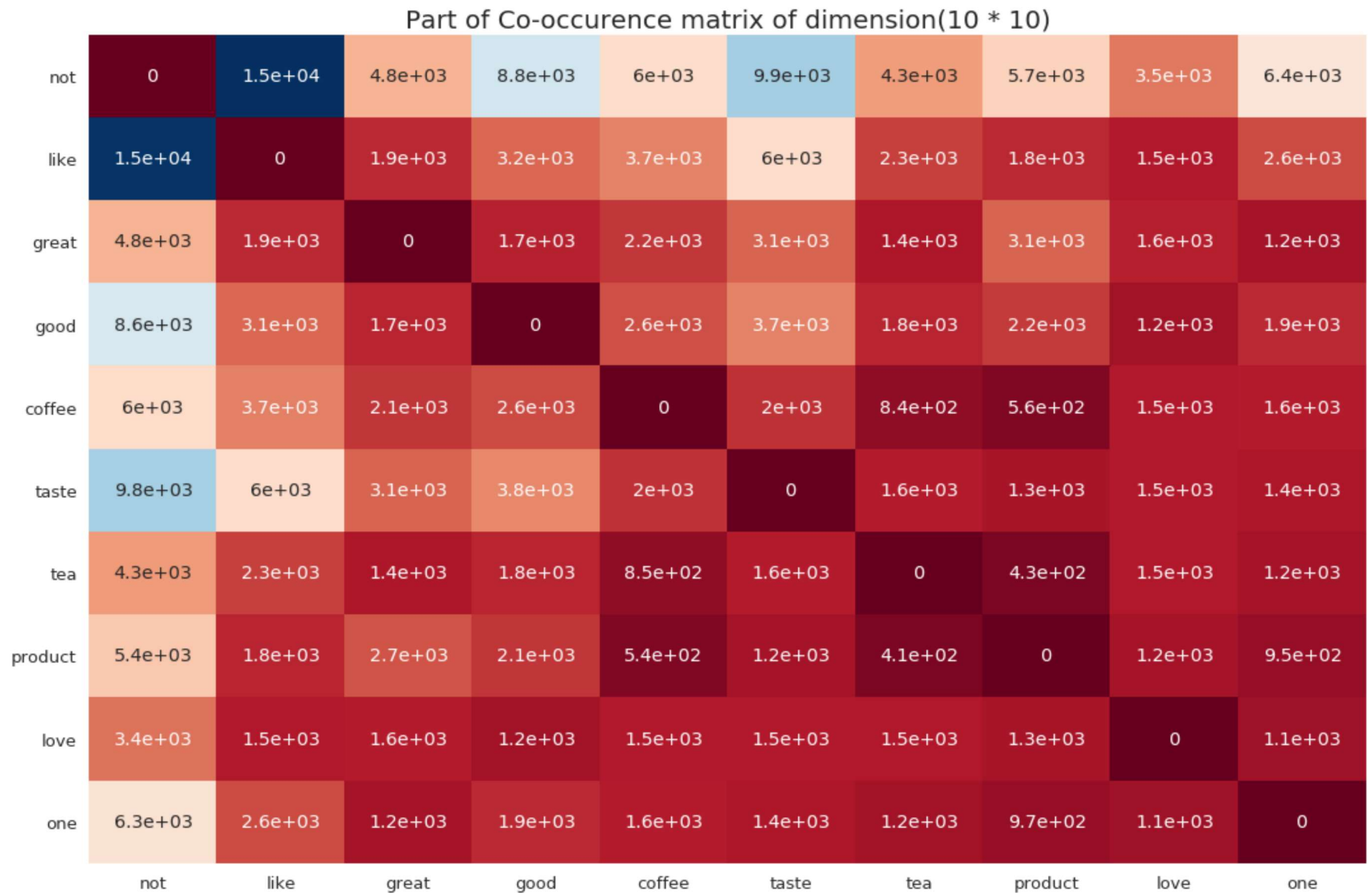
100%|██████████| 87773/87773 [2:39:38<00:00, 9.16it/s]

CPU times: user 34.5 s, sys: 20.7 s, total: 55.2 s

Wall time: 2h 39min 41s

[1.2.3] Visualization of co-occurrence matrix:

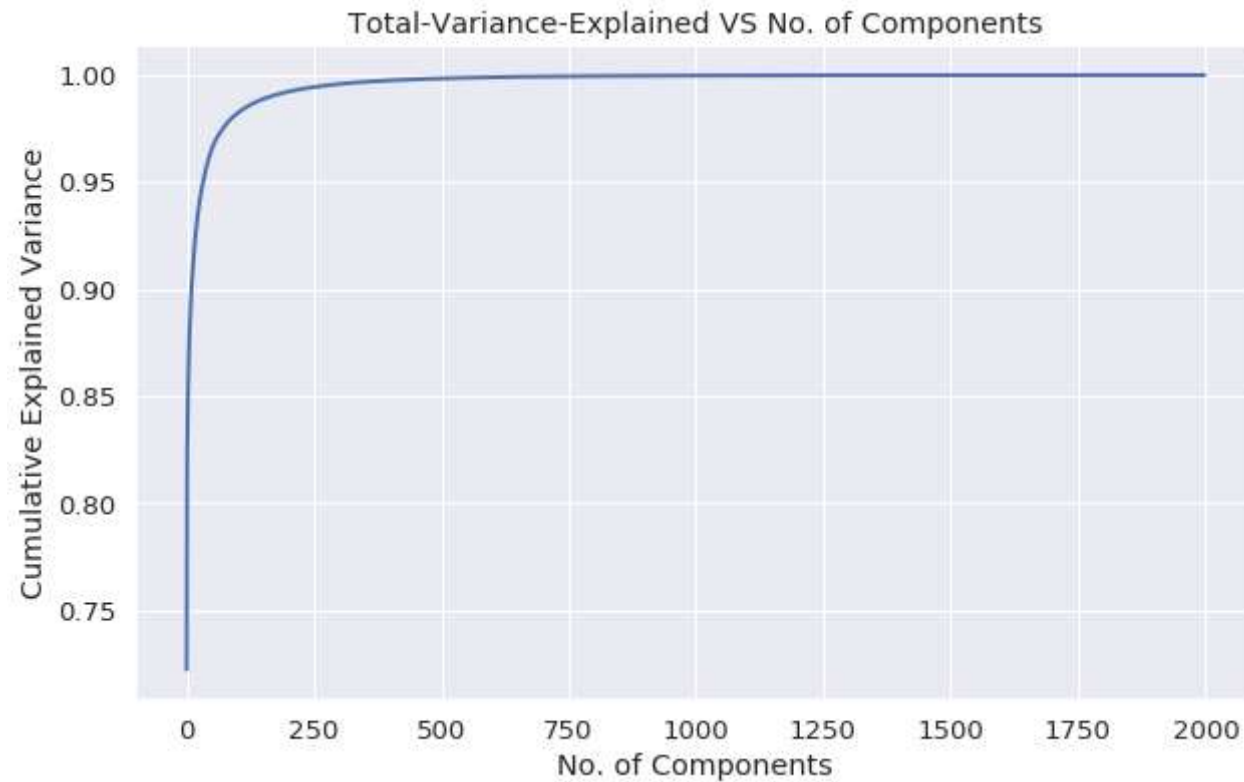
In [29]: 1 visualizeCooMat(coo_all,corpus,dim=10)



[1.3] Finding optimal value for number of components (n) to be retained

[1.3.1] find optimal no. of component and draw elbow curve:

```
In [16]: 1 opt_component=coo_all.shape[1]-1  
2 OptimalComponents_Mat(coo_all, opt_component, transformed_mat=False)
```

**Observation:**

1. from the above plot we observe that 98% of the variance is explained by the 250 features. So we select 250 as optimal no. of features.

[1.3.2] Retrained with optimal components:

```
In [17]: ▶ 1 optimal_components=250
          2 transformed_coo_mat = OptimalComponents_Mat(coo_all,optimal_components,transformed_mat=True)
          3 print(transformed_coo_mat.shape)
          4
```

(2000, 250)

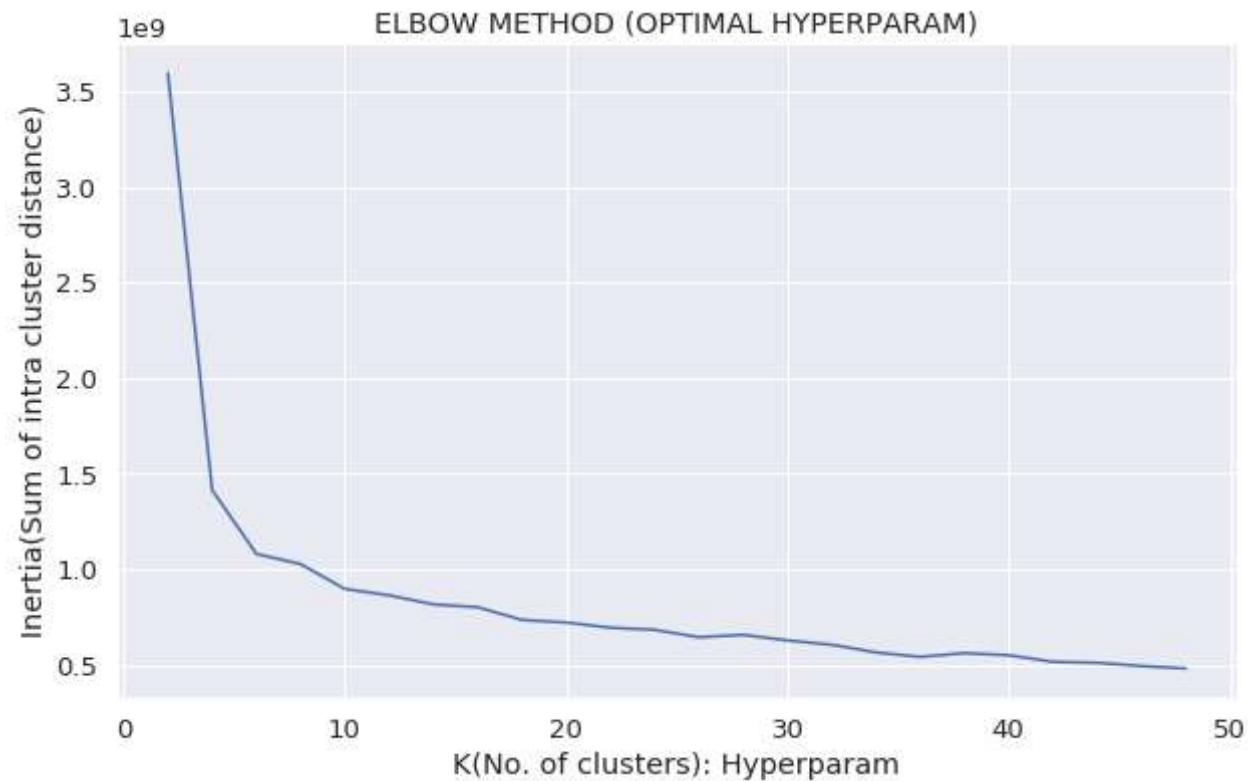
[1.4] Applying k-means clustering

```
In [18]: ▶ 1 params={'clusters':range(2,50,2)}
```

[1.4.1]Hyperparam Tunning using Elbow Method:

```
In [19]: 1 clustering_model(transformed_coo_mat,params)
```

```
100%|██████████| 24/24 [00:00<00:00, 42.81it/s]
```



Observation:

1. from the above plot we observe that after $k=6$, inertia reduces gradually so we choose $k=6$ as our knee point

[1.4.2] corpus labelling and distribution of words per cluster:

```
In [20]: 1 optimal_clusters=6  
        2 model,cluster_label_df=corpus_labelling(transformed_coo_mat,corpus,optimal_clusters)
```

```
0    1732
```

```
4     224
```

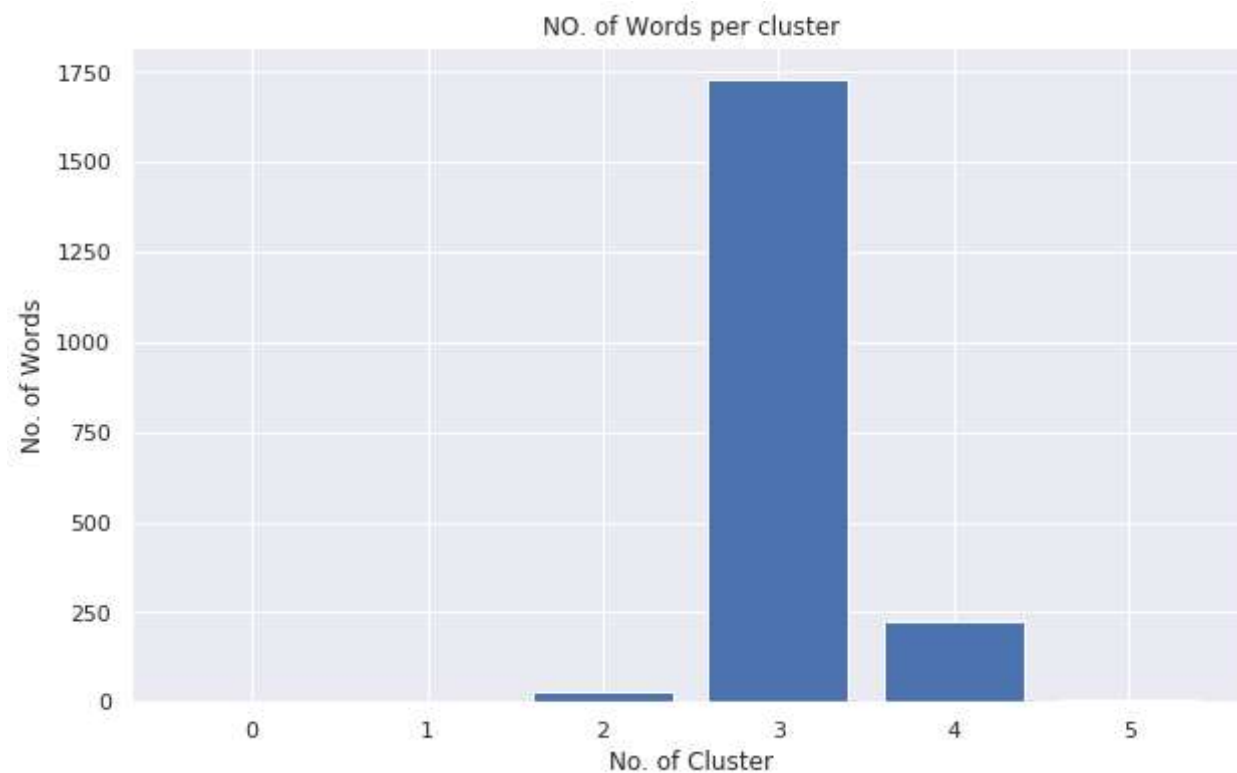
```
1      32
```

```
5      10
```

```
2       1
```

```
3       1
```

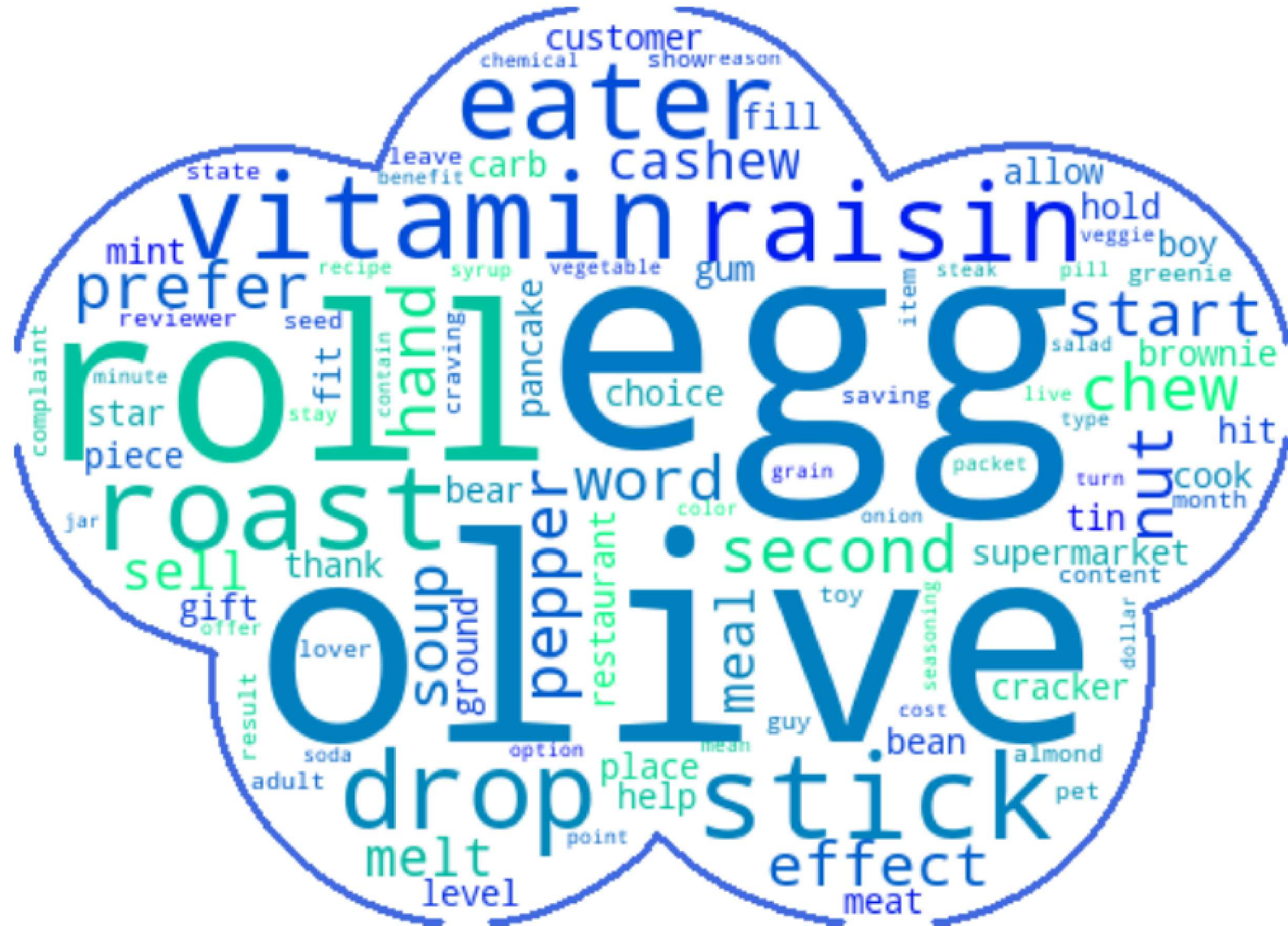
```
Name: label, dtype: int64
```



[1.5] Wordclouds of clusters obtained in the above section

```
1 words_cluster=wordcloud_each_cluster(optimal_clusters,cluster_label_df)
```

WORDS BELONGS TO CLUSTER-0



WORDS BELONGS TO CLUSTER-1



WORDS BELONGS TO CLUSTER-2



WORDS BELONGS TO CLUSTER-3



WORDS BELONGS TO CLUSTER-4



WORDS BELONGS TO CLUSTER-5



Observation:

1. From the above wordcloud we observe that :
 - a. cluster-0 contains eatable words like cashew, egg, raisin and olive etc.

[1.6] Function that returns most similar words for a given word.

```
In [32]: 1 word=input('Input word to find similar words to that:')#'great'
          2 similarity_df=similarWords(transformed_coo_mat, corpus, word, words_cluster)
          3 similarity_df
```

Input word to find similar words to that:eat

Out[32]:

| | also | amazon | best | better | buy | chocolate | could | cup | dog | drink | ... | price | sugar | sweet | tastes | think | time | tried | try | us |
|-----|-------|--------|-------|--------|-------|-----------|-------|------|-------|-------|-----|-------|-------|-------|--------|-------|-------|-------|------|------|
| eat | 0.873 | 0.701 | 0.707 | 0.803 | 0.877 | 0.738 | 0.889 | 0.58 | 0.783 | 0.783 | ... | 0.696 | 0.766 | 0.841 | 0.685 | 0.931 | 0.829 | 0.833 | 0.86 | 0.85 |

1 rows × 32 columns

Observation:

1. from the above analysis we observe that words in a cluster have high similarity with each other.

[6] Conclusions

1. TSVD is used for dimensionality reduction by discarding features which represent least amount of information.
2. Applying TSVD on Co-Occurrence matrix decomposes the co-occurrence matrix into 3- matrices:
 - a. left singular matrix
 - b. right singular matrix
 - c. diagonal matrix having singular values in diagonal
3. The product of left singular matrix and diagonal matrix of singular values shows word vector representation.
4. C-Occurrence matrix preserves the semantic relationship between words.

The optimal values for:

1. number of components = 250
2. number of clusters = 6

REFERENCE LINKS:

1. <https://stackoverflow.com/questions/41661801/python-calculate-the-co-occurrence-matrix>
(<https://stackoverflow.com/questions/41661801/python-calculate-the-co-occurrence-matrix>)
2. <https://www.analyticsvidhya.com/blog/2017/06/word-embeddings-count-word2veec/> (<https://www.analyticsvidhya.com/blog/2017/06/word-embeddings-count-word2veec/>)
3. <https://medium.com/data-science-group-iitr/word-embedding-2d05d270b285> (<https://medium.com/data-science-group-iitr/word-embedding-2d05d270b285>)
4. <https://stackoverflow.com/questions/37331708/nltk-find-occurrences-of-a-word-within-5-words-left-right-of-context-words-in>
(<https://stackoverflow.com/questions/37331708/nltk-find-occurrences-of-a-word-within-5-words-left-right-of-context-words-in>)

In []: ▶

1

