# APPLY KMEANS , AGGLOMERATIVE & DBSCAN CLUSTERING ALGORITHM on Fine Food Reviews

The Amazon Fine Food Reviews dataset consists of reviews of fine foods from Amazon.

Number of reviews: 568,454

Number of users: 256,059

Number of products: 74,258

Timespan: Oct 1999 - Oct 2012

Number of Attributes/Columns in data: 10

Attribute Information:

1. Id
2. ProductId - unique identifier for the product
3. UserId - unqiue identifier for the user
4. ProfileName
5. HelpfulnessNumerator - number of users who found the review helpful
6. HelpfulnessDenominator - number of users who indicated whether they found the review helpful or not
7. Score - rating between 1 and 5
8. Time - timestamp for the review
9. Summary - brief summary of the review
10. Text - text of the review

Objective:

Label the review as Positive or negative

## ▾ Loading,Cleaning & Preprocessing the data

The dataset is available in two forms

1. .csv file
2. SQLite Database

   In order to load the data, We have used the SQLITE dataset as it easier to query the data and visualise the data efficiently.

   Here as we only want to get the global sentiment of the recommendations (positive or negative), we will purposefully ignc equal to 3. If the score id above 3, then the recommendation wil be set to "positive". Otherwise, it will be set to "negative".

```
 1 %matplotlib inline
 2 import warnings
 3 import sqlite3
 4 import pandas as pd
 5 import numpy as np
 6 import nltk
 7 import string
 8 import matplotlib.pyplot as plt
 9 import seaborn as sns
10 from sklearn.feature_extraction.text import TfidfTransformer
11 from sklearn.feature_extraction.text import TfidfVectorizer
12 from sklearn.feature_extraction.text import CountVectorizer
13 from sklearn.metrics import confusion_matrix
14 from sklearn import metrics
15 from sklearn.metrics import roc_curve, auc
```

```
15 |from sklearn.metrics import roc_curve, auc
16 from nltk.stem.porter import PorterStemmer
17
18 import re
19 import string
20 from nltk.corpus import stopwords
21 from nltk.stem import PorterStemmer
22 from nltk.stem.wordnet import WordNetLemmatizer
23
24 from gensim.models import Word2Vec
25 from gensim.models import KeyedVectors
26 import pickle
27
28 from tqdm import tqdm
29 import os
30
31 warnings.filterwarnings("ignore")
```

> /usr/local/lib/python3.6/dist-packages/statsmodels/tools/_testing.py:19: FutureWarning: pandas.util.testing is d
>     import pandas.util.testing as tm

```
1 from google.colab import drive
2 drive.mount('/content/drive')
```

> Go to this URL in a browser: https://accounts.google.com/o/oauth2/auth?client_id=947318989803-6bn6qk8qdgf4n4g3pf
>
>    Enter your authorization code:
>    ..........
>    Mounted at /content/drive

```
1 con = sqlite3.connect("/content/drive/My Drive/Colab Notebooks/database.sqlite")
2
3 filtered_data=pd.read_sql_query("""SELECT * FROM Reviews WHERE Score != 3""",con);
4 filtered_data.head(3)
```

> | | Id | ProductId | UserId | ProfileName | HelpfulnessNumerator | HelpfulnessDenominator |
> |---|----|-----------|--------|-------------|----------------------|------------------------|
> | **0** | 1 | B001E4KFG0 | A3SGXH7AUHU8GW | delmartian | 1 | 1 |
> | **1** | 2 | B00813GRG4 | A1D87F6ZCVE5NK | dll pa | 0 | 0 |
> | **2** | 3 | B000LQOCH0 | ABXLMWJIXXAIN | Natalia Corres "Natalia Corres" | 1 | 1 |

```
1 filtered_data.shape
2
3 def partition(x):
4   if x < 3 :
5     return 'negative'
6   return 'positive'
7
8 actualScore=filtered_data['Score']
9 positive_negative=actualScore.map(partition)
10 filtered_data['Score']=positive_negative
11 print("Number of datapoints",filtered_data.shape)
12 filtered_data.head(3)
```

>

Number of datapoints (525814, 10)

| | Id | ProductId | UserId | ProfileName | HelpfulnessNumerator | HelpfulnessDenominator |
|---|---|---|---|---|---|---|
| **0** | 1 | B001E4KFG0 | A3SGXH7AUHU8GW | delmartian | 1 | 1 |

```
1 display = pd.read_sql_query("""
2 SELECT UserId, ProductId, ProfileName, Time, Score, Text, COUNT(*)
3 FROM Reviews
4 GROUP BY UserId
5 HAVING COUNT(*)>1
6 """, con)
```

```
1 print(display.shape)
2 display.head(3)
```

(80668, 7)

| | UserId | ProductId | ProfileName | Time | Score | |
|---|---|---|---|---|---|---|
| **0** | #oc-R115TNMSPFT9I7 | B007Y59HVM | Breyton | 1331510400 | 2 | Overall its just OK when considering t |
| **1** | #oc-R11D9D7SHXIJB9 | B005HG9ET0 | Louis E. Emory "hoppy" | 1342396800 | 5 | My wife has recurring extreme muscle sp |
| **2** | #oc-R11DNU2NBKQ23Z | B007Y59HVM | Kim Cieszykowski | 1348531200 | 1 | This coffee is horrible and unfortuna |

```
1 display= pd.read_sql_query("""
2 SELECT *
3 FROM Reviews
4 WHERE Score != 3 AND UserId="AR5J8UI46CURR"
5 ORDER BY ProductID
6 """, con)
7 display.head()
```

| | Id | ProductId | UserId | ProfileName | HelpfulnessNumerator | HelpfulnessDenominator | Score | |
|---|---|---|---|---|---|---|---|---|
| **0** | 78445 | B000HDL1RQ | AR5J8UI46CURR | Geetha Krishnan | 2 | 2 | 5 | 119 |
| **1** | 138317 | B000HDOPYC | AR5J8UI46CURR | Geetha Krishnan | 2 | 2 | 5 | 119 |
| **2** | 138277 | B000HDOPYM | AR5J8UI46CURR | Geetha Krishnan | 2 | 2 | 5 | 119 |
| **3** | 73791 | B000HDOPZG | AR5J8UI46CURR | Geetha Krishnan | 2 | 2 | 5 | 119 |
| **4** | 155049 | B000PAQ75C | AR5J8UI46CURR | Geetha Krishnan | 2 | 2 | 5 | 119 |

```
1 sorted_data=filtered_data.sort_values('ProductId',axis=0,ascending=True,inplace=False,kind='quicksort',na_position
2
3 final_data=sorted_data.drop_duplicates(subset={"UserId","ProfileName","Time","Text"},keep='first',inplace=False)
4 final_data.shape
```

(364173, 10)

```
1 (final_data['Id'].size*1.0)/(filtered_data['Id'].size*1.0)*100
```

69.25890143662969

```
1 display= pd.read_sql_query("""
2 SELECT *
3 FROM Reviews
```

```
4 WHERE Score != 3 AND Id=44737 OR Id=64422
5 ORDER BY ProductID
6 """, con)
7
8 display.head()
```

| | Id | ProductId | UserId | ProfileName | HelpfulnessNumerator | HelpfulnessDenominator | Score |
|---|---|---|---|---|---|---|---|
| 0 | 64422 | B000MIDROQ | A161DK06JJMCYF | J. E. Stephens "Jeanne" | 3 | 1 | |
| 1 | 44737 | B001EQ55RW | A2V0I904FH7ABY | Ram | 3 | 2 | |

```
1 final_data=final_data[final_data.HelpfulnessNumerator<=final_data.HelpfulnessDenominator]
```

```
1 nltk.download('stopwords')
```

```
[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data]   Unzipping corpora/stopwords.zip.
True
```

```
1 stopping_words = set(stopwords.words('english'))
2 print(stopping_words)
```

```
{'mightn', 'wouldn', 'i', 'because', 'shan', 'ours', 'few', 'about', "that'll", 'now', 'myself', 'all', 'such',
```

```
1 def clean_html(text):
2     clean_r = re.compile('<,*?>')
3     clean_text = re.sub(clean_r,'',text)
4     return clean_text
5
6 def Clean_punc(text):
7     clean_sentence = re.sub(r'[?|!|\'|"|#]',r' ',text)
8     clean_data = re.sub(r'[.|,|)|(|\|/)]',r' ',clean_sentence)
9     return clean_data
```

```
1 from tqdm import tqdm
2 import os
3 import pdb
4 import pickle
5
6 from tqdm import tqdm
7 import os
8 import pdb
9 import pickle
10
11 stem_no = nltk.stem.SnowballStemmer('english')
12
13 if not os.path.isfile('final_data.sqlite'):
14     final_string=[]
15     all_positive_words=[]
16     all_negative_words=[]
17     for i,sentence in enumerate(tqdm(final_data['Text'].values)):
18         filtered_sentence=[]
19         sent_without_html_tags=clean_html(sentence)
20         #pdb.set_trace()
21         for w in sent_without_html_tags.split():
22             for cleaned_words in Clean_punc(w).split():
23                 if ((cleaned_words.isalpha()) & (len(cleaned_words) > 2)):
24                     if(cleaned_words.lower() not in stopping_words) :
25                         stemming=(stem_no.stem(cleaned_words.lower())).encode('utf8')
```

```
26                        filtered_sentence.append(stemming)
27                        if(final_data['Score'].values)[i]=='positive':
28                            all_positive_words.append(stemming)
29                        if(final_data['Score'].values)[i]=='negative':
30                            all_negative_words.append(stemming)
31        str1 = b" ".join(filtered_sentence)
32        final_string.append(str1)
33
34    final_data['Cleaned_text']=final_string
35    final_data['Cleaned_text']=final_data['Cleaned_text'].str.decode("utf-8")
36
37    conn = sqlite3.connect('final_data.sqlite')
38    cursor=conn.cursor
39    conn.text_factory = str
40    final_data.to_sql('Reviews',conn,schema=None,if_exists='replace',index=True,index_label=None,chunksize=None,dt
41    conn.close()
42
43
44    with open('positive_words.pkl','wb') as f :
45        pickle.dump(all_positive_words,f)
46    with open('negative_words.pkl','wb') as f :
47        pickle.dump(all_negative_words,f)
```

```
100%|██████████| 364171/364171 [06:11<00:00, 980.74it/s]
```

```
1 final_data['total_words'] = [len(x.split()) for x in final_data['Cleaned_text'].tolist()]
```

```
1 final_data.sort_values(by=['Time'], inplace=True, ascending=True)
```

```
1 final_data['Score'].value_counts()
```

```
positive    307061
negative     57110
Name: Score, dtype: int64
```

```
1 final_data=final_data[0:50000]
2 final_data.head(2)
```

| | Id | ProductId | UserId | ProfileName | HelpfulnessNumerator | HelpfulnessDenominator | Score |
|---|---|---|---|---|---|---|---|
| **138706** | 150524 | 0006641040 | ACITT7DI6IDDL | shari zychinski | 0 | 0 | positive |
| **138683** | 150501 | 0006641040 | AJ46FKXOVC7NR | Nicholas A Mesiano | 2 | 2 | positive |

## KMEAN CLUSTERING ALGORITHM

### KMEANS ON AMAZON FINE FOOD REVIEW TO FIND THE CLUSTER AND PLOT
### OF CLUSTERS WITH INERTIA(INTRA CLUSTER DISTANCE) -----> BAG OF WORD
### VECTORIZATION TECHNIQUE

```
1 final_data_bow=[]
2 print(final_data.shape)
```

```
 5
 4
 5 print("="*100)
 6
 7 count_vector=CountVectorizer(min_df=100)
 8 final_data_bow=(count_vector.fit_transform(final_data['Cleaned_text'].values))
 9
10 print("After vectorizations")
11 print(final_data_bow.shape)
12
13 print("="*100)
```

```
(50000, 12)
====================================================================================================
After vectorizations
(50000, 1750)
====================================================================================================
```

```
1 pd.DataFrame(final_data_bow.toarray(),columns=count_vector.get_feature_names())
```

|  | abl | absolut | absorb | accept | accord | acid | acquir | across | act | activ | actual | ad | add | addict | addit | a |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 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 | |
| 2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | |
| 3 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | |
| 4 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | |
| 49995 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | |
| 49996 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | |
| 49997 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 3 | 0 | 0 | |
| 49998 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | |
| 49999 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | |

50000 rows × 1750 columns

▾ Apply KMeans to Find the Hyperparameter i.e Best No of Cluster (K)

```
 1 from sklearn.model_selection import GridSearchCV
 2 from sklearn.model_selection import cross_val_score
 3 from sklearn.cluster import KMeans
 4
 5 n_clusters = [1,2,3,4,5,6,7,8,9,10]
 6 plot_inertia=[]
 7 def K_Means_Clustering(x_training_data):
 8   for i  in n_clusters:
 9     KMeans_Clustering = KMeans(init='k-means++',max_iter=300,n_clusters=i)
10     KMeans_Clustering.fit(x_training_data)
11     plot_inertia.append(KMeans_Clustering.inertia_)
12   plt.plot(n_clusters, plot_inertia)
13   plt.xlabel("Number of Clusters: Hyperparameter")
14   plt.ylabel("Inertia(Intra Cluster Distance")
15   plt.title("Plot Number of Clusters v/s Inertia(Intra Cluster Distance)")
16   plt.grid()
17   plt.show()
```

```
1 K_Means_Clustering(final_data_bow)
```



Plot Number of Clusters v/s Inertia(Intra Cluster Distance)

- ▾ By looking at the above Plot. Looks like the BEST K is 3. At K=3 , there is an inflexion

```
1 from sklearn.cluster import KMeans
2 K_Mean_Model=KMeans(init='k-means++',max_iter=300,n_clusters=3)
3 K_Mean_Model.fit(final_data_bow)
```

```
KMeans(algorithm='auto', copy_x=True, init='k-means++', max_iter=300,
       n_clusters=3, n_init=10, n_jobs=None, precompute_distances='auto',
       random_state=None, tol=0.0001, verbose=0)
```

- ▾ Find the Top 10 Words for each Cluster

```
1 print("Top terms per cluster")
2 centroids=K_Mean_Model.cluster_centers_.argsort()[:,::-1]
3
4 features=count_vector.get_feature_names()
5 top_words=[]
6 for i in range(3):
7     top_ten_words = [features[ind] for ind in centroids[i,:10]]
8     top_words.append(top_ten_words)
9     print("Clusters {} : {} ".format(i,' '.join(top_ten_words)))
```

```
Top terms per cluster
Clusters 0 : tast great love good like product flavor use one tri
Clusters 1 : tea flavor green like tast bag drink use make one
Clusters 2 : like tast flavor one use product good tri food make
```

```
1 top_words[0]
```

```
['tast',
 'great',
 'love',
 'good',
 'like',
 'product',
 'flavor',
 'use',
 'one',
 'tri']
```

- ▾ GENERATE THE TOP 10 WORDS FROM EACH CLUSTER USING WORD CLOUD

```
1 from wordcloud import WordCloud
```

```
1  from wordcloud import WordCloud
2  import matplotlib.pyplot as plt
3  %matplotlib inline
4
5  for i in range(3) :
6    wordcloud = WordCloud(background_color='white').generate(str(top_words[i]))
7    plt.figure(figsize=(8,8),facecolor=None)
8    plt.imshow(wordcloud)
9    plt.axis("off")
10   plt.tight_layout(pad = 0)
11   plt.show()
12   print('--------------------------------------------------------------------------')
```

# KMEANS ON AMAZON FINE FOOD REVIEW TO FIND THE CLUSTER AND PLOT
▾ OF CLUSTERS WITH INERTIA(SUM OF ALL INTRA CLUSTER DISTANCE) ----> tf-
VECTORIZATION TECHNIQUE

```
1 from sklearn.feature_extraction.text import TfidfVectorizer
2 final_data_tfidf=[]
3 print(final_data.shape)
4
5
6 print("="*100)
7
8 tfidf_vector=TfidfVectorizer(min_df=100)
9 final_data_tfidf=(tfidf_vector.fit_transform(final_data['Cleaned_text'].values))
10
11 print("After vectorizations")
12 print(final_data_tfidf.shape)
13
14 print("="*100)
```

⤷　(50000, 12)
　　====================================================================================================
　　After vectorizations
　　(50000, 1750)
　　====================================================================================================

```
1 pd.DataFrame(final_data_tfidf.toarray(),columns=tfidf_vector.get_feature_names())
```

⤷

|  | abl | absolut | absorb | accept | accord | acid | acquir | across | act | activ | actual | ad | a |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0.144869 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.00000 | 0.0 | 0.000000 | 0.000000 | 0.0000 |
| 1 | 0.000000 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.00000 | 0.0 | 0.000000 | 0.000000 | 0.0000 |
| 2 | 0.000000 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.35954 | 0.0 | 0.000000 | 0.000000 | 0.0000 |
| 3 | 0.000000 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.00000 | 0.0 | 0.000000 | 0.000000 | 0.0000 |
| 4 | 0.000000 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.00000 | 0.0 | 0.000000 | 0.000000 | 0.0000 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | |
| 49995 | 0.000000 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.00000 | 0.0 | 0.000000 | 0.109359 | 0.0000 |
| 49996 | 0.000000 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.00000 | 0.0 | 0.000000 | 0.000000 | 0.0000 |
| 49997 | 0.000000 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.00000 | 0.0 | 0.000000 | 0.000000 | 0.2861 |
| 49998 | 0.000000 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.00000 | 0.0 | 0.000000 | 0.000000 | 0.0000 |
| 49999 | 0.000000 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.00000 | 0.0 | 0.194975 | 0.000000 | 0.0000 |

50000 rows × 1750 columns

▾ Apply KMeans to find the Hyperparameter i.e Best No of Cluster(K)

```
1 K_Means_Clustering(final_data_tfidf)
```

⤷

**BY LOOKING AT THE ABOVE PLOT , LOOKS IKE THE BEST K IS 5.**

```
1 from sklearn.cluster import KMeans
2 K_Mean_Model=KMeans(init='k-means++',max_iter=300,n_clusters=5)
3 K_Mean_Model.fit(final_data_tfidf)
```

```
KMeans(algorithm='auto', copy_x=True, init='k-means++', max_iter=300,
       n_clusters=5, n_init=10, n_jobs=None, precompute_distances='auto',
       random_state=None, tol=0.0001, verbose=0)
```

**FIND TOP 10 WORDS FOR EACH CLUSTER**

```
1 print("Top terms per cluster")
2 centroids=K_Mean_Model.cluster_centers_.argsort()[:,::-1]
3
4 features=tfidf_vector.get_feature_names()
5 top_words=[]
6 for i in range(5):
7   top_ten_words = [features[ind] for ind in centroids[i,:10]]
8   top_words.append(top_ten_words)
9   print("Clusters {} : {} ".format(i,' '.join(top_ten_words)))
```

```
Top terms per cluster
Clusters 0 : tast like good flavor great love use tri one make
Clusters 1 : product order amazon store price find ship great buy love
Clusters 2 : tea green flavor drink tast like bag cup love good
Clusters 3 : coffe cup flavor pod tast like roast strong good tri
Clusters 4 : dog cat treat food love chew one eat like get
```

**GENERATE THE TOP 10 WORDS FROM EACH CLUSTER USING WORD CLOUD**

```
 1 from wordcloud import WordCloud
 2 import matplotlib.pyplot as plt
 3 %matplotlib inline
 4
 5 for i in range(5) :
 6   wordcloud = WordCloud(background_color='white').generate(str(top_words[i]))
 7   print("="*100)
 8   plt.figure(figsize=(8,8),facecolor=None)
 9   plt.imshow(wordcloud)
10   plt.axis("off")
11   plt.tight_layout(pad = 0)
12   plt.show()
13   print("="*100)
```

==========================================================================



==========================================================================
==========================================================================



==========================================================================
==========================================================================



==========================================================================
==========================================================================

# KMEANS ON AMAZON FINE FOOD REVIEW TO FIND THE CLUSTER AND PLOT OF CLUSTERS WITH INERTIA(SUM OF ALL INTRA CLUSTER DISTANCE) --> AV WORD2VEC

```
1 from gensim.models import Word2Vec
2 from gensim.models import KeyedVectors
3 import pickle
4
5 list_of_sent_train_avgw2v=[]
6
7 for sent_train_avgw2v in tqdm(final_data['Cleaned_text'].values):
8     list_of_sent_train_avgw2v.append(sent_train_avgw2v.split())
```

```
100%|██████████| 50000/50000 [00:00<00:00, 83276.96it/s]
```

```
1 w2v_model_train = Word2Vec(list_of_sent_train_avgw2v,min_count=5,size=50,workers=4)
2 w2v_words_svm_train=list(w2v_model_train.wv.vocab)
```

```
1 train_vectors=[];
2 for sent in tqdm(list_of_sent_train_avgw2v):
3     sent_vec=np.zeros(50)
4     cnt_words=0;
5     for word in sent:
6         if word in w2v_words_svm_train:
7             vec=w2v_model_train.wv[word]
8             sent_vec+=vec
9             cnt_words+=1
10    if cnt_words !=0:
11        sent_vec/=cnt_words
12    train_vectors.append(sent_vec)
13 print(len(train_vectors))
14 print(len(train_vectors[0]))
```

```
100%|██████████| 50000/50000 [00:54<00:00, 911.49it/s]50000
50
```

## ▾ Apply KMeans to find the Hyperparameter i.e Best No of Clusters(K)

```
1 K_Means_Clustering(train_vectors)
```

⤷



## ▾ By Looking at the Plot. Looks like No Of Cluster is 5 i.e K=5

```
1 K_Mean_Model=KMeans(init='k-means++',max_iter=300,n_clusters=5)
2 predict_cluster_index=K_Mean_Model.fit_predict(train_vectors)
```

```
1 cluster_index =[]
2
3 for i in range(len(predict_cluster_index)):
4   if predict_cluster_index[i] == 2:
5     cluster_index.append(i)
```

```
1 text_features=[]
2
3 for i in range(len(cluster_index)) :
4   text_features.append(list_of_sent_train_avgw2v[cluster_index[i]])
```

```
1 from wordcloud import WordCloud
2 import matplotlib.pyplot as plt
3 %matplotlib inline
4
5 text_var = ''
6 for j in range(len(text_features)):
7     for i in range(len(text_features[j])):
8         text_var = text_var + text_features[j][i] + ' '
9 wordcloud = WordCloud(background_color='white').generate(str(text_var))
10 print("="*100)
11 plt.figure(figsize=(8,8),facecolor=None)
12 plt.imshow(wordcloud)
13 plt.axis("off")
14 plt.tight_layout(pad = 0)
15 plt.show()
16 print("="*100)
```

⤷

=============================================================================================



# KMEANS ON AMAZON FINE FOOD REVIEW TO FIND THE CLUSTER AND PLOT

▾ OF CLUSTERS WITH INERTIA(SUM OF ALL INTRA CLUSTER DISTANCE) ----> TI
WORD 2VEC VECTORIZATION TECHNIQUE

```
1 model_Avgw2v = TfidfVectorizer()
2 X_Train_Avgw2v=model_Avgw2v.fit_transform(final_data['Cleaned_text'].values)
```

```
1 dictionary = dict(zip(model_Avgw2v.get_feature_names(), list(model_Avgw2v.idf_)))
```

```
 1 tfidf_feature=model_Avgw2v.get_feature_names()
 2
 3 tfidf_sent_vectors_train=[];
 4 row=0;
 5
 6 for sent in tqdm(list_of_sent_train_avgw2v):
 7     sent_vec=np.zeros(50)
 8     weight_sum=0;
 9     for word in sent :
10         if word in w2v_words_svm_train and word in tfidf_feature :
11             vec=w2v_model_train.wv[word]
12             tf_idf=dictionary[word]*(sent.count(word)/len(sent))
13             sent_vec+=(vec*tf_idf)
14             weight_sum+=tf_idf
15
16     if weight_sum!=0:
17         sent_vec/=weight_sum
18     tfidf_sent_vectors_train.append(sent_vec)
19     row+=1
```

⊳ 100%|██████████| 50000/50000 [10:50<00:00, 76.90it/s]

▾ Apply KMeans to find the Best Hyperparameter i.e Best No of Clusters K

```
1 K_Means_Clustering(tfidf_sent_vectors_train)
```

⊳

Plot Number of Clusters v/s Inertia(Intra Cluster Distance)



▾ By looking at the plot. It looks like the Best K will be 3 i.e K=3

## No of Clusters = 3

```
1 K_Mean_Model=KMeans(init='k-means++',max_iter=300,n_clusters=3)
2 fit_model=K_Mean_Model.fit(tfidf_sent_vectors_train)
```

```
1 print("Top terms per cluster")
2 centroids=fit_model.cluster_centers_.argsort()[:,::-1]
3
4 tfidf_feature=model_Avgw2v.get_feature_names()
5
6 top_words=[]
7 for i in range(3):
8    top_ten_words = [tfidf_feature[ind] for ind in centroids[i,:10]]
9    top_words.append(top_ten_words)
10   print("Clusters {} : {} ".format(i,' '.join(top_ten_words)))
```

```
☐→  Top terms per cluster
    Clusters 0 : aberr aad aafco abound aboth aback abigirl abc abour aand
    Clusters 1 : abigirl aaaaah aand aberr aad aboth abondanza abour abound aback
    Clusters 2 : aboard aah abj aardvark abi abigirl aaa aad aand aberr
```

```
1 from wordcloud import WordCloud
2 import matplotlib.pyplot as plt
3 %matplotlib inline
4
5 for i in range(3) :
6    wordcloud = WordCloud(background_color='white').generate(str(top_words[i]))
7    print("="*100)
8    plt.figure(figsize=(8,8),facecolor=None)
9    plt.imshow(wordcloud)
10   plt.axis("off")
11   plt.tight_layout(pad = 0)
12   plt.show()
13   print("="*100)
```

☐→

=============================================================================================



=============================================================================================
=============================================================================================



## ▾ AGGLOMERATIVE CLUSTERING ALGORITHM

## ▾ USING AVGWORD2VEC FEATURIZATION TECHNIQUE

```
1 final_data_new=final_data[:5000]
```

```
1 final_data_new=final_data_new['Cleaned_text']
```

```
1 final_data_new_array=final_data_new.to_numpy()
```

```
1 final_data_new_array[66]
```

    'tim burton start take serious michael keaton unleash unus restrain alec baldwin support cast known looni danc p

```
1 from gensim.models import Word2Vec
2 from gensim.models import KeyedVectors
3 import pickle
4
5 list_of_sent_avgw2v_agglo=[]
6
7 for sent_train_avgw2v in tqdm(final_data_new):
8     list_of_sent_avgw2v_agglo.append(sent_train_avgw2v.split())
```

```
100%|████████| 5000/5000 [00:00<00:00, 134248.66it/s]
```

```
1 w2v_model_train = Word2Vec(list_of_sent_avgw2v_agglo,min_count=5,size=50,workers=4)
2 w2v_words_svm_train=list(w2v_model_train.wv.vocab)
```

```
1 train_vectors=[];
2 for sent in tqdm(list_of_sent_avgw2v_agglo):
3     sent_vec=np.zeros(50)
4     cnt_words=0;
5     for word in sent:
6         if word in w2v_words_svm_train:
7             vec=w2v_model_train.wv[word]
8             sent_vec+=vec
9             cnt_words+=1
10    if cnt_words !=0:
11        sent_vec/=cnt_words
12    train_vectors.append(sent_vec)
13 print(len(train_vectors))
14 print(len(train_vectors[0]))
```

```
100%|████████| 5000/5000 [00:05<00:00, 999.47it/s] 5000
50
```

## APPLY AGGLOMERATIVE CLUSTERING TO FIND THE BEST HYPERPARAMETER I.E BEST NO OF CLUS

### ▾ For K = 2

```
1 from sklearn.cluster import AgglomerativeClustering
2
3 k = 2
4 cluster_label=[]
5 Agglo_Clustering = AgglomerativeClustering(n_clusters=k)
6 cluster_label=Agglo_Clustering.fit_predict(train_vectors)
```

```
1 C1,C2 =[],[]
2
3 for i in range(Agglo_Clustering.labels_.shape[0]):
4   if Agglo_Clustering.labels_[i]==0 :
5     C1.append(final_data_new_array[i])
6   else:
7     C2.append(final_data_new_array[i])
```

### ▾ CLUSTER 1 FEATURES

```
1 display_data=''
2
3 for i in C1:
4   display_data+=str(i)
5
6 from wordcloud import WordCloud
7 import matplotlib.pyplot as plt
8 wordcloud = WordCloud(background_color='white').generate(display_data)
9 print("="*100)
10 plt.figure(figsize=(8,8),facecolor=None)
11 plt.imshow(wordcloud)
12 plt.axis("off")
13 plt.tight_layout(pad = 0)
```

```
14 plt.show()
15 print("="*100)
```

⎘    =======================================================================================



     =======================================================================================

## ▾ CLUSTER 2 FEATURES

```
1 display_data_C2=''
2
3 for i in C2:
4   display_data_C2+=str(i)
5 wordcloud = WordCloud(background_color='white').generate(display_data_C2)
6 print("="*100)
7 plt.figure(figsize=(8,8),facecolor=None)
8 plt.imshow(wordcloud)
9 plt.axis("off")
10 plt.tight_layout(pad = 0)
11 plt.show()
12 print("="*100)
```

⎘    =======================================================================================



     =======================================================================================

## ▾ For K=5

```
1 from sklearn.cluster import AgglomerativeClustering
2 k = 5
3 cluster_label=[]
```

```
3 cluster_label=[]
4 Agglo_Clustering = AgglomerativeClustering(n_clusters=k)
5 cluster_label=Agglo_Clustering.fit_predict(train_vectors)
```

```
1 def display_data_agglo(cluster):
2   display_data=''
3   for i in cluster:
4     display_data+=str(i)
5   wordcloud = WordCloud(background_color='white').generate(display_data)
6   print("="*100)
7   plt.figure(figsize=(8,8),facecolor=None)
8   plt.imshow(wordcloud)
9   plt.axis("off")
10  plt.tight_layout(pad = 0)
11  plt.show()
12  print("="*100)
```

```
1 C1,C2,C3,C4,C5 =[],[],[],[],[]
2
3 for i in range(Agglo_Clustering.labels_.shape[0]):
4   if Agglo_Clustering.labels_[i]==0 :
5     C1.append(final_data_new_array[i])
6   elif Agglo_Clustering.labels_[i]==1 :
7     C2.append(final_data_new_array[i])
8   elif Agglo_Clustering.labels_[i]==2 :
9     C3.append(final_data_new_array[i])
10  elif Agglo_Clustering.labels_[i]==3 :
11    C4.append(final_data_new_array[i])
12  else:
13    C5.append(final_data_new_array[i])
```

## ▾ CLUSTER 1 FEATURES

```
1 from wordcloud import WordCloud
2 import matplotlib.pyplot as plt
3 display_data_agglo(C1)
```

⊳  ====================================================================================================



====================================================================================================

## ▾ CLUSTER 2 FEATURES

```
1 display_data_agglo(C2)
```

⊳    ================================================================================



================================================================================

▾ CLUSTER 3 FEATURES

```
1 display_data_agglo(C3)
```

⊳    ================================================================================



================================================================================

▾ CLUSTER 4 FEATURES

```
1 display_data_agglo(C4)
```

⊳

```
============================================================================================
```



## ▾ CLUSTER 5 FEATURES



```
1 display_data_agglo(C5)
```

```
============================================================================================
```



```
============================================================================================
```

## ▾ USING TF-IDF WORD2VEC FEATURIZATION TECHNIQUE

```
1 model_Avgw2v = TfidfVectorizer()
2 X_Train_Avgw2v=model_Avgw2v.fit_transform(final_data_new)
```

```
1 dictionary = dict(zip(model_Avgw2v.get_feature_names(), list(model_Avgw2v.idf_)))
```

```
1 tfidf_feature=model_Avgw2v.get_feature_names()
2
3 tfidf_sent_vectors_train=[];
4 row=0;
5
6 for sent in tqdm(list_of_sent_avgw2v_agglo):
7     sent_vec=np.zeros(50)
8     weight_sum=0;
9     for word in sent :
10         if word in w2v_words_svm_train and word in tfidf_feature :
11             vec=w2v_model_train.wv[word]
12             tf_idf=dictionary[word]*(sent.count(word)/len(sent))
13             sent_vec+=(vec*tf_idf)
14             weight_sum+=tf_idf
15
16     if weight_sum!=0:
17         sent_vec/=weight_sum
18     tfidf_sent_vectors_train.append(sent_vec)
19     row+=1
```

```
100%|██████████| 5000/5000 [00:33<00:00, 150.70it/s]
```

```
1 k = 2
2 cluster_label=[]
3 Agglo_Clustering = AgglomerativeClustering(n_clusters=k)
4 cluster_label=Agglo_Clustering.fit_predict(tfidf_sent_vectors_train)
```

```
1 C1,C2 =[],[]
2
3 for i in range(Agglo_Clustering.labels_.shape[0]):
4   if Agglo_Clustering.labels_[i]==0 :
5     C1.append(final_data_new_array[i])
6   else:
7     C2.append(final_data_new_array[i])
```

## ▾ CLUSTER 1 FEATURES

```
1 display_data_agglo(C1)
```

[→  =====================================================================================



=====================================================================================

## ▾ CLUSTER 2 FEATURES

```
1 display_data_agglo(C2)
```

[→  =====================================================================================



=====================================================================================

```
1 k = 5
2 cluster_label=[]
3 Agglo_Clustering = AgglomerativeClustering(n_clusters=k)
4 cluster_label=Agglo_Clustering.fit_predict(tfidf_sent_vectors_train)
```

```
1  C1,C2,C3,C4,C5 =[],[],[],[],[]
2
3  for i in range(Agglo_Clustering.labels_.shape[0]):
4    if Agglo_Clustering.labels_[i]==0 :
5      C1.append(final_data_new_array[i])
6    elif Agglo_Clustering.labels_[i]==1 :
7      C2.append(final_data_new_array[i])
8    elif Agglo_Clustering.labels_[i]==2 :
9      C3.append(final_data_new_array[i])
10   elif Agglo_Clustering.labels_[i]==3 :
11     C4.append(final_data_new_array[i])
12   else:
13     C5.append(final_data_new_array[i])
```

▾ CLUSTER 1 FEATURES

```
1  display_data_agglo(C1)
```



```
1  display_data_agglo(C2)
```



```
1  display_data_agglo(C3)
```

=========================================================================================



```
1 display_data_agglo(C4)
```

=========================================================================================



=========================================================================================

```
1 display_data_agglo(C5)
```

=========================================================================================



=========================================================================================

## ▾ DBSCAN CLUSTERING ALGORITHM

## USING AVGWORD2VEC FEATURIZATION TECHNIQUE APPLY DBSCAN CLUSTERING TECHNIQUE TO F BEST HYPERPARAMETER

```
1 from sklearn.cluster import DBSCAN
2
3
4 cluster_label=[]
5 dbscan_Clustering = DBSCAN(eps=0.5,n_jobs=-1,algorithm='kd_tree')
6 dbscan_Clustering.fit(train_vectors)
7 print('No Of Clusters' ,len(set(dbscan_Clustering.labels_)))
```

```
No Of Clusters 2
```

## ▾ No of Cluster for eps = 0.5 is 2

```
1 print('Cluster are ignored',set(dbscan_Clustering.labels_))
```

```
Cluster are ignored {0, -1}
```

```
1 C1,C2 =[],[]
2
3 for i in range(dbscan_Clustering.labels_.shape[0]):
4   if dbscan_Clustering.labels_[i]==0 :
5     C1.append(final_data_new_array[i])
6   else:
7     C2.append(final_data_new_array[i])
```
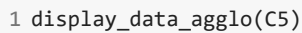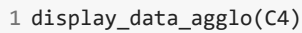
## ▾ For K=2

## ▾ CLUSTER 1 FEATURES

```
 1 display_data=''
 2
 3 for i in C1:
 4   display_data+=str(i)
 5
 6 from wordcloud import WordCloud
 7 import matplotlib.pyplot as plt
 8 wordcloud = WordCloud(background_color='white').generate(display_data)
 9 print("="*100)
10 plt.figure(figsize=(8,8),facecolor=None)
11 plt.imshow(wordcloud)
12 plt.axis("off")
13 plt.tight_layout(pad = 0)
14 plt.show()
15 print("="*100)
```

==================================================================================================



## ▼ CLUSTER 2 FEATURES

```
1 display_data=''
2
3 for i in C2:
4   display_data+=str(i)
5
6 from wordcloud import WordCloud
7 import matplotlib.pyplot as plt
8 wordcloud = WordCloud(background_color='white').generate(display_data)
9 print("="*100)
10 plt.figure(figsize=(8,8),facecolor=None)
11 plt.imshow(wordcloud)
12 plt.axis("off")
13 plt.tight_layout(pad = 0)
14 plt.show()
15 print("="*100)
```

[→  ==================================================================================================



==================================================================================================

## ▼ USING TF-IDF FEATURIZATION TECHNIQUE APPLY DBSCAN CLUSTERING TECHNIQUE TO FIND THE HYPERPARAMETER

```
1 dbscan_Clustering = DBSCAN(eps=0.1,n_jobs=-1,algorithm='kd_tree')
2 dbscan_Clustering.fit(tfidf_sent_vectors_train)
3 print('No Of Clusters' ,len(set(dbscan_Clustering.labels_)))
```

[→  No Of Clusters 10

```
1 dbscan_Clustering = DBSCAN(eps=0.5,n_jobs=-1,algorithm='kd_tree')
2 dbscan_Clustering.fit(tfidf_sent_vectors_train)
3 print('No Of Clusters' ,len(set(dbscan_Clustering.labels_)))
```

> No Of Clusters 2

```
1 dbscan_Clustering = DBSCAN(eps=1,n_jobs=-1,algorithm='kd_tree')
2 dbscan_Clustering.fit(tfidf_sent_vectors_train)
3 print('No Of Clusters' ,len(set(dbscan_Clustering.labels_)))
```

> No Of Clusters 2

```
1 C1,C2 =[],[]
2
3 for i in range(dbscan_Clustering.labels_.shape[0]):
4   if dbscan_Clustering.labels_[i]==0 :
5     C1.append(final_data_new_array[i])
6   else:
7     C2.append(final_data_new_array[i])
```

```
1 display_data=''
2
3 for i in C1:
4   display_data+=str(i)
5
6 from wordcloud import WordCloud
7 import matplotlib.pyplot as plt
8 wordcloud = WordCloud(background_color='white').generate(display_data)
9 print("="*100)
10 plt.figure(figsize=(8,8),facecolor=None)
11 plt.imshow(wordcloud)
12 plt.axis("off")
13 plt.tight_layout(pad = 0)
14 plt.show()
15 print("="*100)
```

> ====================================================================================================



====================================================================================================
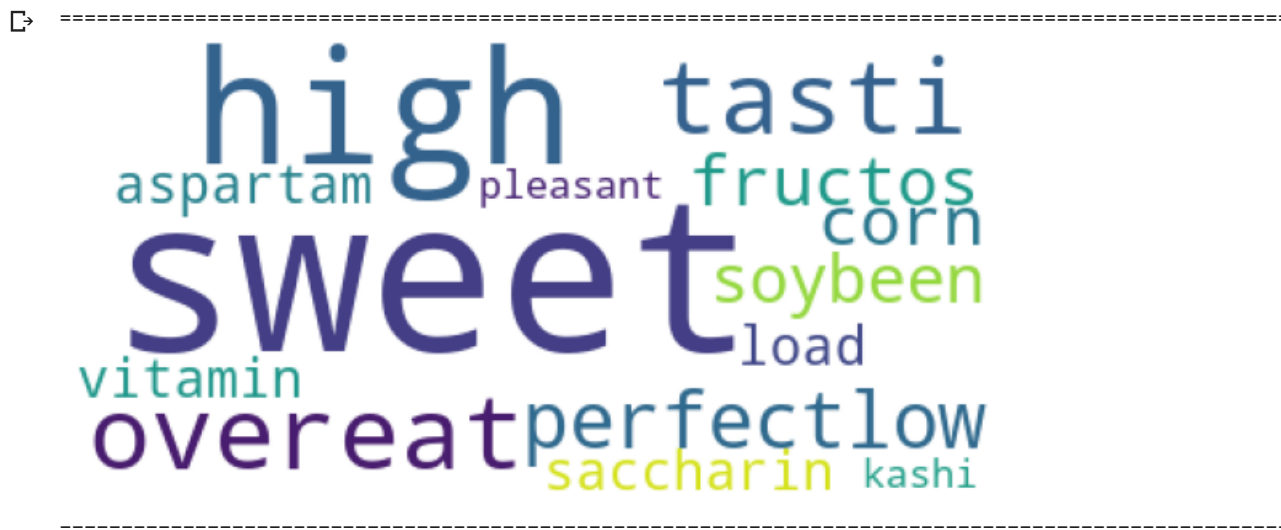
```
1 display_data=''
2
3 for i in C2:
4   display_data+=str(i)
5
6 from wordcloud import WordCloud
7 import matplotlib.pyplot as plt
8 wordcloud = WordCloud(background_color='white').generate(display_data)
9 print("="*100)
10 plt.figure(figsize=(8,8),facecolor=None)
11 plt.imshow(wordcloud)
12 plt.axis("off")
```

```
13 plt.tight_layout(pad = 0)
14 plt.show()
15 print("="*100)
```

====================================================================================



====================================================================================

```
 1 from prettytable import PrettyTable
 2
 3 x = PrettyTable()
 4 x.field_names = ["Clustering Algorithm Technique","Vectorizer","Best k"]
 5 x.add_row(['KMeans','BOW','3'])
 6 x.add_row(['KMeans','TFIDF','5'])
 7 x.add_row(['KMeans','AVG W2vec','5'])
 8 x.add_row(['KMeans','TFIDF W2vec','3'])
 9 x.add_row(['AGGLOMERATIVE','AVGWORD2VEC','(2,5)'])
10 x.add_row(['AGGLOMERATIVE','TFIDFW2vec','(2,5)'])
11 x.add_row(['DBSCAN','AVGW2vec(eps =0.5)','2'])
12 x.add_row(['DBSCAN','TFIDFW2vec(eps=0.1)','10'])
13 x.add_row(['DBSCAN','TFIDFW2vec(eps=0.5)','2'])
14 x.add_row(['DBSCAN','TFIDFW2vec(eps=1)','2'])
15 print(x)
```

```
+-------------------------------+---------------------+--------+
| Clustering Algorithm Technique |      Vectorizer     | Best k |
+-------------------------------+---------------------+--------+
|            KMeans             |         BOW         |   3    |
|            KMeans             |        TFIDF        |   5    |
|            KMeans             |      AVG W2vec      |   5    |
|            KMeans             |     TFIDF W2vec     |   3    |
|         AGGLOMERATIVE         |     AVGWORD2VEC     | (2,5)  |
|         AGGLOMERATIVE         |      TFIDFW2vec     | (2,5)  |
|            DBSCAN             |  AVGW2vec(eps =0.5) |   2    |
|            DBSCAN             | TFIDFW2vec(eps=0.1) |   10   |
|            DBSCAN             | TFIDFW2vec(eps=0.5) |   2    |
|            DBSCAN             |  TFIDFW2vec(eps=1)  |   2    |
+-------------------------------+---------------------+--------+
```

```
 1
```