

▼ [1]. Reading Data

Let me copy code from my previous assignment where we are doing preprocessing steps like deduplication, decontraction of words, lemmatization etc.

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

↳ Go to this URL in a browser: https://accounts.google.com/o/oauth2/auth?client_id=9473

```
Enter your authorization code:
.....
Mounted at /content/drive
```

```
cd drive/My\ Drive
```

↳ /content/drive/My Drive

```
%matplotlib inline
import warnings
warnings.filterwarnings("ignore")
import nltk
nltk.download('stopwords')

import sqlite3
import pandas as pd
import numpy as np
import nltk
import string
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.feature_extraction.text import TfidfTransformer
from sklearn.feature_extraction.text import TfidfVectorizer

from sklearn.feature_extraction.text import CountVectorizer
from sklearn.metrics import confusion_matrix
from sklearn import metrics
from sklearn.metrics import roc_curve, auc
from nltk.stem.porter import PorterStemmer
from nltk.stem import PorterStemmer
from nltk.stem.snowball import SnowballStemmer

import re
# Tutorial about Python regular expressions: https://pymotw.com/2/re/
import string
```

```
import nltk
from nltk.corpus import stopwords
from nltk.stem import PorterStemmer
from nltk.stem.wordnet import WordNetLemmatizer

from gensim.models import Word2Vec
from gensim.models import KeyedVectors
import pickle
from sklearn.model_selection import cross_val_score
from sklearn.metrics import accuracy_score
from sklearn.metrics import roc_auc_score

from tqdm import tqdm
import os
```

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

```
#con = sqlite3.connect(r"D:\AppliedAI\AAIC_Course_handouts\11_Amazon Fine Food Reviews\amazon_fine_food_reviews.sqlite")
con = sqlite3.connect(r"database.sqlite")
data = pd.read_sql_query(""" SELECT * FROM Reviews WHERE Score != 3""",con)
#data = pd.read_csv(""" SELECT * FROM Reviews WHERE Score != 3""",con)

# Change Score with 1 n 2 as -ve and 4 n 5 as +ve

def chng_to_0_or_1 (Score):
    if Score ==4 or Score ==5:
        return 1
    elif Score ==1 or Score ==2:
        return 0
    else:# Thus in case by some mistake any data is their with rating 6 or 7 etc due to some error
        pass
currentScore = data["Score"]
new_Score = currentScore.map(chng_to_0_or_1)
data["Score"] = new_Score
print ("Number of data points available")
print (data.shape)#Gives original number of data points available
```

```
#2 Data Cleaning a.) Getting rid of duplicates and b.) if helpfulnessdenominator < helpfulnessNumerator
```

```
data = data.drop_duplicates(subset = ["UserId","ProfileName","HelpfulnessNumerator","HelpfulnessDenominator"])
print ("Number of data points after removing duplicates")
print (data.shape)#Gives data points are deduplication
```

```
# Reference: Copied from above cell  final=final[final.HelpfulnessNumerator<=final.HelpfulnessDenominator]
data=data[data.HelpfulnessNumerator<=data.HelpfulnessDenominator]
print ("Number of data points after removing where HelpfulnessNumerator is more than HelpfulnessDenominator")
print (data.shape)
```

#3 Preprocessing begins

```
#Convert to lower case, convert shortcut words to proper words, remove Special Character
```

```
#i) Convert to lower case:
```

```
data["Text"] = (data["Text"].str.lower())
data["Summary"] = (data["Summary"].str.lower())
```

```
#ii) Convert Shortcuts words to proper words
```

```
#List of Words are:https://en.wikipedia.org/wiki/Wikipedia:List\_of\_English\_contractions
#Reference:https://stackoverflow.com/questions/39602824/pandas-replace-string-with-another
data['Text'] = data['Text'].replace({"ain't": "am not", "amn't": "am not", "aren't": "are not", "can't": "cannot", "cause": "because", "could've": "could have", "couldn't": "could not", "couldn't": "dare not", "daresn't": "dare not", "dasn't": "dare not", "didn't": "did not", "doesn't": "do not", "e'er": "ever", "everyone's": "everyone is", "finna": "fixing to", "gimme": "gimme", "gonna": "going to", "gon't": "go not", "gotta": "got to", "hadn't": "had not", "hasn't": "has not", "he'd": "he had", "he'll": "he shall", "he's": "he has", "he've": "he have", "how'd": "how did", "how're": "how are", "how's": "how has", "I'd": "I had", "I'll": "I shall", "I'm": "I am", "I'm'a": "I'm'o": "I am going to", "I've": "I have", "isn't": "is not", "it'd": "it would", "it'll": "it shall", "let's": "let us", "mayn't": "may not", "may've": "may have", "mightn't": "might not", "might've": "mustn't": "must not", "mustn't've": "must not have", "must've": "must have", "needn't": "need not", "o'clock": "of the clock", "o'er": "", "ol'": "old", "oughtn't": "ought not", "shalln't": "shall not", "she'd": "she had", "she'll": "she shall", "she's": "she is", "should've": "should have", "shouldn't've": "should not have", "somebody's": "somebody has", "someone's": "someone has", "so": "that'll": "that will", "that're": "that are", "that's": "that is", "that'd": "that would", "there": "there'll": "there shall", "there're": "there are", "there's": "there is", "these're": "these are", "they'll": "they will", "they're": "they are", "they've": "they have", "this's": "", "those're": "those are", "twas": "it was", "wasn't": "was not", "we'd": "we had", "we'd've": "we would have", "we'll": "we will", "we've": "we have", "weren't": "were not", "what'd": "what did", "what'll": "what will", "what're": "what've": "what have", "when's": "when is", "where'd": "where did", "where're": "where are", "which's": "which has", "who'd": "who would", "who'd've": "who would have", "who'll": "who shall", "who's": "who has", "who've": "who have", "why'd": "why did", "why're": "why are", "why's": "why have", "would've": "would have", "wouldn't": "would not", "y'all": "you all", "you'd": "you had", "you'll": "you've": "you have"})
```

```
##### Lets do the same for summary Text#####
```

```
data['Summary'] = data['Summary'].replace({"ain't": "am not", "amn't": "am not", "aren't": "are not", "can't": "cannot", "cause": "because", "could've": "could have", "couldn't": "could not", "couldn't": "dare not", "daresn't": "dare not", "dasn't": "dare not", "didn't": "did not", "doesn't": "do not", "e'er": "ever", "everyone's": "everyone is", "finna": "fixing to", "gimme": "gimme", "gonna": "going to", "gon't": "go not", "gotta": "got to", "hadn't": "had not", "hasn't": "has not", "he'd": "he had", "he'll": "he shall", "he's": "he has", "he've": "he have", "how'd": "how did", "how're": "how are", "how's": "how has", "I'd": "I had", "I'll": "I shall", "I'm": "I am", "I'm'a": "I'm'o": "I am going to", "I've": "I have", "isn't": "is not", "it'd": "it would", "it'll": "it shall", "let's": "let us", "mayn't": "may not", "may've": "may have", "mightn't": "might not", "might've": "mustn't": "must not", "mustn't've": "must not have", "must've": "must have", "needn't": "need not", "o'clock": "of the clock", "o'er": "", "ol'": "old", "oughtn't": "ought not", "shalln't": "shall not", "she'd": "she had", "she'll": "she shall", "she's": "she is", "should've": "should have", "shouldn't've": "should not have", "somebody's": "somebody has", "someone's": "someone has", "so": "that'll": "that will", "that're": "that are", "that's": "that is", "that'd": "that would", "there": "there'll": "there shall", "there're": "there are", "there's": "there is", "these're": "these are", "they'll": "they will", "they're": "they are", "they've": "they have", "this's": "", "those're": "those are", "twas": "it was", "wasn't": "was not", "we'd": "we had", "we'd've": "we would have", "we'll": "we will", "we've": "we have", "weren't": "were not", "what'd": "what did", "what'll": "what will", "what're": "what've": "what have", "when's": "when is", "where'd": "where did", "where're": "where are", "which's": "which has", "who'd": "who would", "who'd've": "who would have", "who'll": "who shall", "who's": "who has", "who've": "who have", "why'd": "why did", "why're": "why are", "why's": "why have", "would've": "would have", "wouldn't": "would not", "y'all": "you all", "you'd": "you had", "you'll": "you've": "you have"})
```

```

twas : it was , wasn t : was not , we a : we naa , we a ve : we would nave , we ll : we
"we've": "we have", "weren't": "were not", "what'd": "what did", "what'll": "what will", "what're"
"what've": "what have", "when's": "when is", "where'd": "where did", "where're": "where are", "wh
"which's": "which has", "who'd": "who would", "who'd've": "who would have", "who'll": "who shall"
"who's": "who has", "who've": "who have", "why'd": "why did", "why're": "why are", "why's": "why ha
"would've": "would have", "wouldn't": "would not", "y'all": "you all", "you'd": "you had", "you'l
"you've": "you have"})
#####
# iii) Remove Special Characters except alphabets and numbers
# The reason i dont want to remove number people might write got five eggs as 5 eggs or vice
# that information which could be useful
#Ref:https://stackoverflow.com/questions/33257344/how-to-remove-special-characters-from-a-c
data["Text"] = data["Text"].map(lambda x: re.sub(r'[^a-zA-Z_0-9 -]', ' ', x))
data["Summary_copy"] = data["Summary"].map(lambda x: re.sub(r'[^a-zA-Z_0-9 -]', ' ', x))

```

```

#The Summary are usually so small if we remove few stopwords the meaning itself would be lost
# So let us see what all stopwords we have
#Ref:https://stackoverflow.com/questions/5511708/adding-words-to-nltk-stoplist
#https://chrisalbon.com/machine\_learning/preprocessing\_text/remove\_stop\_words/

```

```

stopwords = nltk.corpus.stopwords.words('english')
newStopWords = ['would', 'could', 'br', '<br>', '<', '>']
notstopwords = ['not', 'no', 'nor']
stopwords.extend(newStopWords)
stopwords = [word for word in stopwords if word not in notstopwords]

```

```

# iv) For now let us just go with flow will use default stopwords as creating our own stop
# Rather will use n-gram strategy to get rid of problem of stopwords removal changing the
#Ref:https://stackoverflow.com/questions/43184364/python-remove-stop-words-from-pandas-data
data["New_Text"] = data['Text'].apply(lambda x: [item for item in str.split(x) if item not in stopwords])
data["Summary"] = data['Summary_copy'].apply(lambda x: [item for item in str.split(x) if item not in stopwords])

```

```

#Ref:https://stackoverflow.com/questions/37347725/convert-a-pandas-df-list-into-a-string
# we are creating new column New_summary so in case in future we need summary it is intact
data["New_Text"] = data["New_Text"].apply(' '.join)
data["Summary"] = data["Summary"].apply(' '.join)

```

```

# v) Now lets do Stemming
#https://stackoverflow.com/questions/48617589/beginner-stemming-in-pandas-produces-letter
english_stemmer = SnowballStemmer('english', ignore_stopwords=True)
data["New_Text"] = data["New_Text"].apply(english_stemmer.stem)
data["Summary"] = data["Summary"].apply(english_stemmer.stem)
data["New_Text"] = data["New_Text"].astype(str)
data["Summary"] = data["Summary"].astype(str)

```

```

#vi) stemming without removing stop words
english_stemmer = SnowballStemmer('english', ignore_stopwords=True)
#https://stackoverflow.com/questions/34724246/attributeerror-float-object-has-no-attribute
data["Text_with_stop"] = data["Text"].astype(str)
data["Summary"] = data["Summary"].astype(str)
data["Text_with_stop"] = data["Text_with_stop"].str.lower().map(english_stemmer.stem)
data["Summary"] = data["Summary"].str.lower().map(english_stemmer.stem)

```

```

data["Text_with_stop"] = data["Text_with_stop"].apply(''.join)
data["Summary"] = data["Summary"].apply(''.join)
data["Text_with_stop"] = data["Text_with_stop"].astype(str)
data["Summary"] = data["Summary"].astype(str)
print(data["Score"].value_counts())
print ("Thus we see there are 85% and 15% positive and negative reviews, thus a unbalanced dataset we first copy negative dataset 6 times than we sample with same number of times as positive reviews")
# Let's include another feature which is the length of the text
data_neg = data[data["Score"] == 0]
data_pos = data[data["Score"] == 1]
data = pd.concat([data_pos, data_neg])
#https://stackoverflow.com/questions/46429033/how-do-i-count-the-total-number-of-words-in-a-string
data["Text_length"] = (data["New_Text"].str.count(' ') + 1)
data["Summary_length"] = (data["Summary"].str.count(' ') + 1)
data["Time_formatted"] = pd.to_datetime(data["Time"])
data.sort_values(by=['Time_formatted'], inplace=True)

```

→ Number of data points available
 (525814, 10)
 Number of data points after removing duplicates
 (366392, 10)
 Number of data points after removing where HelpfulnessNumerator is more than Helpful
 (366390, 10)
 1 308679
 0 57711
 Name: Score, dtype: int64
 Thus we see there are 85% and 15% positive and negative reviews, thus a unbalanced dataset

Let us keep "data" frame intact so incase something gets corrupted we can reuse it. We will create a reduced copy of the data frame and name it "newdata"

```
newdata = data
```

Double-click (or enter) to edit

```

newdata_50K = newdata.tail(50000)
newdata_50K.sort_values(by=['Time_formatted'], inplace=True)

```

```

X_no_stop_50k_reviews = newdata_50K['New_Text'].values
print ((X_no_stop_50k_reviews).shape)

```

→ (50000,)

```
sent= (X_no_stop_50k_reviews.tolist())
```

▼ [4] Featurization

▼ [4.3] TF-IDF

```
tf_idf_vect = TfidfVectorizer(ngram_range=(1,2), min_df=10, max_features=9000)

tfidf_X_no_stop_50k_reviews=tf_idf_vect.fit_transform(X_no_stop_50k_reviews)
print("some sample features(unique 100 words in the corpus)",tf_idf_vect.get_feature_names()
print('*'*50)
print("the shape of out text TFIDF vectorizer ",tfidf_X_no_stop_50k_reviews.get_shape())
print("the number of unique words including both unigrams and bigrams ", tfidf_X_no_stop_50k_reviews.shape[1])

→ some sample features(unique 100 words in the corpus) ['10', '10 12', '10 days', '10 in
=====
the shape of out text TFIDF vectorizer (50000, 9000)
the number of unique words including both unigrams and bigrams 9000
```

▼ [5] Assignment 11: Truncated SVD

1. Apply Truncated-SVD on only this feature set:

- **SET 2:** Review text, preprocessed one converted into vectors using (TFIDF)
- **Procedure:**
 - Take top 2000 or 3000 features from tf-idf vectorizers using idf_ score.
 - You need to calculate the co-occurrence matrix with the selected features (Note: $X \cdot X^T$ returns the covariance matrix, check these blogs [blog-1](#), [blog-2](#) for more information)
 - You should choose the n_components in truncated svd, with maximum explained variance that and implement them. (hint: plot of cumulative explained variance ratio)
 - After you are done with the truncated svd, you can apply K-Means clustering and check the elbow method.
 - Print out wordclouds for each cluster, similar to that in previous assignment.
 - You need to write a function that takes a word and returns the most similar words using vectors(vector: a row in the matrix after truncatedSVD)

▼ Truncated-SVD

▼ [5.1] Taking top features from TFIDF, SET 2

```

features_name_tfidf = tf_idf_vect.get_feature_names() #https://www.kaggle.com/premvardhan,
idf_value_tfidf_vectors = tf_idf_vect.idf_


# sorted_features = np.argsort(tf_idf_vect.idf_[::-1]) #https://github.com/niketan108/Tru

print (features_name_tfidf)
print ("#####")
print (idf_value_tfidf_vectors)
#lst= list(zip(features_name_tfidf,idf_value_tfidf_vectors))
feat_vect_tfidf=np.column_stack((features_name_tfidf, idf_value_tfidf_vectors))

feat_vect_tfidf_sorted = feat_vect_tfidf[feat_vect_tfidf[:,1].argsort()[:-1]]
print ("#####")
print (feat_vect_tfidf_sorted)
print (feat_vect_tfidf_sorted.shape)
feat_vect_tfidf_sorted_new_with_values = feat_vect_tfidf_sorted[0 : 2000, 0 : 2]
feat_vect_tfidf_sorted_new = feat_vect_tfidf_sorted[0 : 2000, 0 : 1]

lst_with_values=(feat_vect_tfidf_sorted_new_with_values[:, [0]]).tolist()
lst=(feat_vect_tfidf_sorted_new[:, [0]]).tolist()

↳ ['10', '10 12', '10 days', '10 minutes', '10 seconds', '10 year', '10 years', '100',
#####
[4.73055513 8.38581108 8.18221212 ... 7.74226084 8.29343776 8.23627935]
#####
[['cocaine' '9.334891634422286']
 ['tabanero' '9.180740954595027']
 ['gbr' '9.180740954595027']
 ...
 ['good' '2.342642032463707']
 ['like' '2.2482930630225186']
 ['not' '2.0246200351810093']]
(9000, 2)

print (feat_vect_tfidf_sorted_new_with_values)
print ("#####")
print (feat_vect_tfidf_sorted_new)

```

↳

```
[['cocaine' '9.334891634422286'],
 ['tabanero' '9.180740954595027'],
 ['gbr' '9.180740954595027'],
 ...
 ['gradually' '8.106226217505977'],
 ['made whole' '8.106226217505977'],
 ['massive' '8.106226217505977']]
#####
[['cocaine'],
 ['tabanero'],
 ['gbr'],
 ...
 ['gradually'],
 ['made whole'],
 ['massive']]
```

```
print (lst_with_values)
print ("#####")
print (lst)
print ("#####")
print( ", ".join( repr(e) for e in lst ))
lst_2000 = [i[0] for i in lst_with_values]
print(lst_2000)
```

```
↳  [['cocaine'], ['tabanero'], ['gbr'], ['nbr'], ['amount per'], ['adult dogs'], ['wind
#####
[['cocaine'], ['tabanero'], ['gbr'], ['nbr'], ['amount per'], ['adult dogs'], ['wind
#####
[['cocaine'], ['tabanero'], ['gbr'], ['nbr'], ['amount per'], ['adult dogs'], ['wind
['cocaine', 'tabanero', 'gbr', 'nbr', 'amount per', 'adult dogs', 'windows', 'sucros
```

```
len(lst_2000)
```

```
↳  2000
```

▼ [5.2] Calculation of Co-occurrence matrix

```
"""def con_mat(word1,word2,sentence):
    str_list = str(sentence)
    c=0
    words = str_list.split()
    indices = [i for (i, x) in enumerate(words) if x == word1]
    #print (indices)

    for i in indices:
        #print ("The Value of index is "+str(i))
        if i >= 5:
```

```

    left_context_win = 5
    #print ("The Value of left context windows is "+str(left_context_win))
    words_sub=words[i-5:i]
    #print (words_sub)
    c1=words_sub.count(word2)
    c=c+c1
    #return (c)

else:
    left_context_win = i
    #print ("The Value of left context windows is "+str(left_context_win))
    words_sub=words[0:i]
    #print (words_sub)
    c1=words_sub.count(word2)
    c=c+c1
    #return (c)

for j in indices:
    #print ("The Value of index is "+str(j))
    #if j >= 5:
    right_context_win = 5
    #print ("The Value of right context windows is "+str(right_context_win))
    words_sub=words[j+1:j+6]
    #print (words_sub)
    c2=words_sub.count(word2)
    c=c+c2
    #return (c)
return (c)
"""

```

```

"""from tqdm import tqdm
occ_matrix =np.zeros(shape=(len(lst), len(lst)), dtype=int)
for rows in tqdm(range(0,10)):
    for columns in (range(0,10)):
        if rows==columns:
            occ_matrix[rows,columns]= 0
        else:
            occ_matrix[rows,columns]= con_mat(lst,lst,testsent)
"""

```

100%|██████████| 10/10 [00:00<00:00, 292.92it/s]

```

from tqdm import tqdm
n_neighbor = 5
occ_matrix = np.zeros((2000,2000))
for row in tqdm(sent):
    words_in_row = row.split()
    for index,word in enumerate(words_in_row):
        if word in lst_2000:
            for j in range(max(index-n_neighbor,0),min(index+n_neighbor,len(words_in_row))-1):

```

```

if words_in_row[j] in lst_2000:
    occ_matrix[lst_2000.index(word),lst_2000.index(words_in_row[j])] += 1
else:
    continue
else:
    continue

```

↳ 100% |██████████| 50000/50000 [00:57<00:00, 875.37it/s]

occ_matrix

```

↳ array([[70.,  0.,  0., ...,  0.,  0.,  0.],
       [ 0., 41.,  0., ...,  0.,  0.,  0.],
       [ 0.,  0., 85., ...,  0.,  0.,  0.],
       ...,
       [ 0.,  0.,  0., ..., 40.,  0.,  0.],
       [ 0.,  0.,  0., ...,  0.,  0.,  0.],
       [ 0.,  0.,  0., ...,  0.,  0., 41.]])

```

▼ [5.3] Finding optimal value for number of components (n) to be retained.

```

from sklearn.decomposition import TruncatedSVD
n_components=[10,20,50,60,100,200,300,400,500,600,700,800,850,900,1000,1100,1200,1500,1999]
for n in n_components:
    tsvd=TruncatedSVD(n_components=n)
    tsvd.fit(occ_matrix)
    expvar=tsvd.explained_variance_ratio_.sum()

    print('n_components=',n,'Explained variance=',expvar)

↳ n_components= 10 Explained variance= 0.07653152891364591
    n_components= 20 Explained variance= 0.11527798377341393
    n_components= 50 Explained variance= 0.20258979265759927
    n_components= 60 Explained variance= 0.2250785958725258
    n_components= 100 Explained variance= 0.2997247223305621
    n_components= 200 Explained variance= 0.4399296943283674
    n_components= 300 Explained variance= 0.5523726212923781
    n_components= 400 Explained variance= 0.6512358389841189
    n_components= 500 Explained variance= 0.7406269341767835
    n_components= 600 Explained variance= 0.8226896582274044
    n_components= 700 Explained variance= 0.897921529430273
    n_components= 800 Explained variance= 0.9659923902719085
    n_components= 850 Explained variance= 0.99522151068381
    n_components= 900 Explained variance= 1.00000000000000135
    n_components= 1000 Explained variance= 1.00000000000000138
    n_components= 1100 Explained variance= 1.00000000000000135
    n_components= 1200 Explained variance= 1.00000000000000133
    n_components= 1500 Explained variance= 1.00000000000000133
    n_components= 1999 Explained variance= 1.00000000000000133

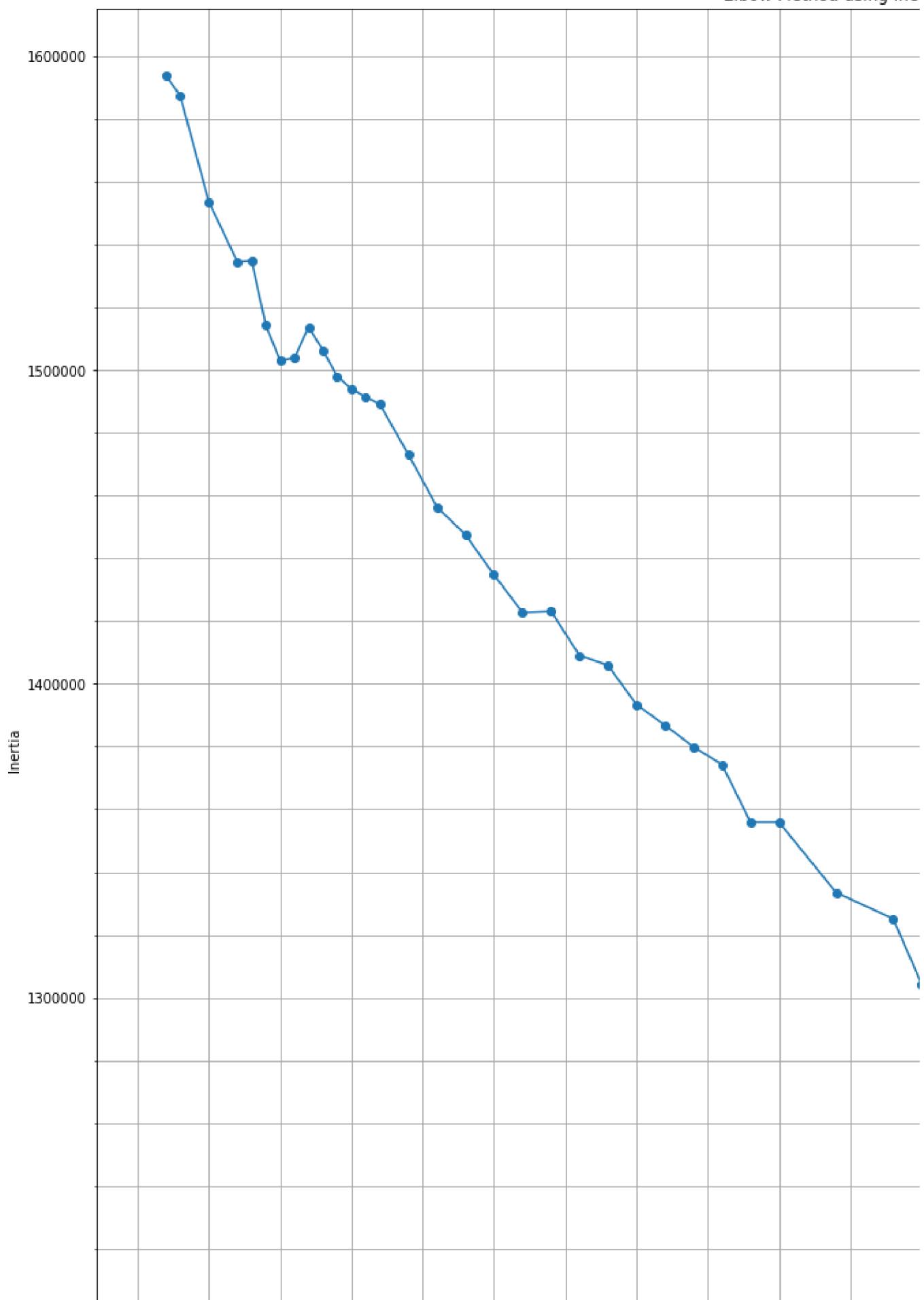
```

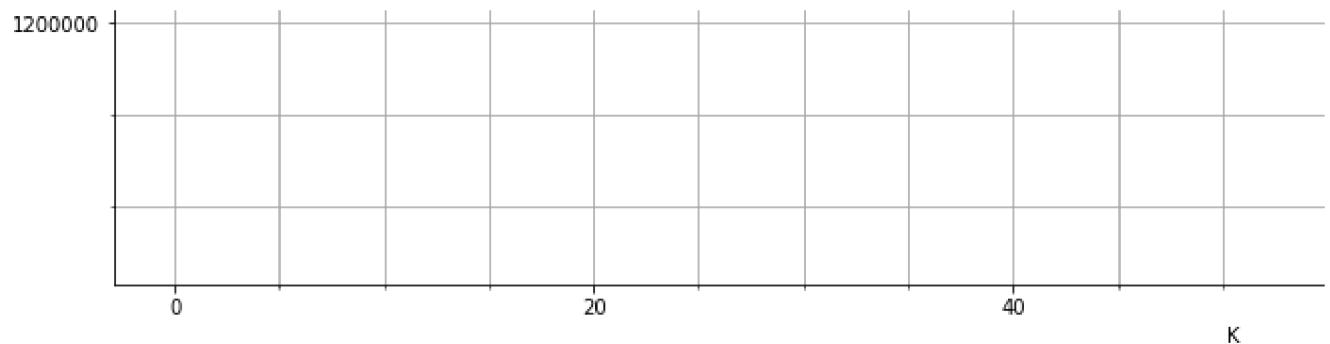
So we can see 850 top components or features almost explain the 99% of variability or the information

▼ [5.4] Applying k-means clustering

```
svd_850 = TruncatedSVD(n_components = 850)
svd_train = svd_850.fit_transform(occ_matrix)

from sklearn.cluster import KMeans
inertia_kmeans = []
n_clusters = [2,3,5,7,8,9,10,11,12,13,14,15,16,17,19,21,23,25,27,29,31,33,35,37,39,41,43,45,47,49,51,53,55,57,59,61,63,65,67,69,71,73,75,77,79,81,83,85,87,89,91,93,95,97,99,101,103,105,107,109,111,113,115,117,119,121,123,125,127,129,131,133,135,137,139,141,143,145,147,149,151,153,155,157,159,161,163,165,167,169,171,173,175,177,179,181,183,185,187,189,191,193,195,197,199,201,203,205,207,209,211,213,215,217,219,221,223,225,227,229,231,233,235,237,239,241,243,245,247,249,251,253,255,257,259,261,263,265,267,269,271,273,275,277,279,281,283,285,287,289,291,293,295,297,299,301,303,305,307,309,311,313,315,317,319,321,323,325,327,329,331,333,335,337,339,341,343,345,347,349,351,353,355,357,359,361,363,365,367,369,371,373,375,377,379,381,383,385,387,389,391,393,395,397,399,401,403,405,407,409,411,413,415,417,419,421,423,425,427,429,431,433,435,437,439,441,443,445,447,449,451,453,455,457,459,461,463,465,467,469,471,473,475,477,479,481,483,485,487,489,491,493,495,497,499,501,503,505,507,509,511,513,515,517,519,521,523,525,527,529,531,533,535,537,539,541,543,545,547,549,551,553,555,557,559,561,563,565,567,569,571,573,575,577,579,581,583,585,587,589,591,593,595,597,599,601,603,605,607,609,611,613,615,617,619,621,623,625,627,629,631,633,635,637,639,641,643,645,647,649,651,653,655,657,659,661,663,665,667,669,671,673,675,677,679,681,683,685,687,689,691,693,695,697,699,701,703,705,707,709,711,713,715,717,719,721,723,725,727,729,731,733,735,737,739,741,743,745,747,749,751,753,755,757,759,761,763,765,767,769,771,773,775,777,779,781,783,785,787,789,791,793,795,797,799,801,803,805,807,809,811,813,815,817,819,821,823,825,827,829,831,833,835,837,839,841,843,845,847,849,851,853,855,857,859,861,863,865,867,869,871,873,875,877,879,881,883,885,887,889,891,893,895,897,899,901,903,905,907,909,911,913,915,917,919,921,923,925,927,929,931,933,935,937,939,941,943,945,947,949,951,953,955,957,959,961,963,965,967,969,971,973,975,977,979,981,983,985,987,989,991,993,995,997,999,1001,1003,1005,1007,1009,1011,1013,1015,1017,1019,1021,1023,1025,1027,1029,1031,1033,1035,1037,1039,1041,1043,1045,1047,1049,1051,1053,1055,1057,1059,1061,1063,1065,1067,1069,1071,1073,1075,1077,1079,1081,1083,1085,1087,1089,1091,1093,1095,1097,1099,1101,1103,1105,1107,1109,1111,1113,1115,1117,1119,1121,1123,1125,1127,1129,1131,1133,1135,1137,1139,1141,1143,1145,1147,1149,1151,1153,1155,1157,1159,1161,1163,1165,1167,1169,1171,1173,1175,1177,1179,1181,1183,1185,1187,1189,1191,1193,1195,1197,1199,1201,1203,1205,1207,1209,1211,1213,1215,1217,1219,1221,1223,1225,1227,1229,1231,1233,1235,1237,1239,1241,1243,1245,1247,1249,1251,1253,1255,1257,1259,1261,1263,1265,1267,1269,1271,1273,1275,1277,1279,1281,1283,1285,1287,1289,1291,1293,1295,1297,1299,1301,1303,1305,1307,1309,1311,1313,1315,1317,1319,1321,1323,1325,1327,1329,1331,1333,1335,1337,1339,1341,1343,1345,1347,1349,1351,1353,1355,1357,1359,1361,1363,1365,1367,1369,1371,1373,1375,1377,1379,1381,1383,1385,1387,1389,1391,1393,1395,1397,1399,1401,1403,1405,1407,1409,1411,1413,1415,1417,1419,1421,1423,1425,1427,1429,1431,1433,1435,1437,1439,1441,1443,1445,1447,1449,1451,1453,1455,1457,1459,1461,1463,1465,1467,1469,1471,1473,1475,1477,1479,1481,1483,1485,1487,1489,1491,1493,1495,1497,1499,1501,1503,1505,1507,1509,1511,1513,1515,1517,1519,1521,1523,1525,1527,1529,1531,1533,1535,1537,1539,1541,1543,1545,1547,1549,1551,1553,1555,1557,1559,1561,1563,1565,1567,1569,1571,1573,1575,1577,1579,1581,1583,1585,1587,1589,1591,1593,1595,1597,1599,1601,1603,1605,1607,1609,1611,1613,1615,1617,1619,1621,1623,1625,1627,1629,1631,1633,1635,1637,1639,1641,1643,1645,1647,1649,1651,1653,1655,1657,1659,1661,1663,1665,1667,1669,1671,1673,1675,1677,1679,1681,1683,1685,1687,1689,1691,1693,1695,1697,1699,1701,1703,1705,1707,1709,1711,1713,1715,1717,1719,1721,1723,1725,1727,1729,1731,1733,1735,1737,1739,1741,1743,1745,1747,1749,1751,1753,1755,1757,1759,1761,1763,1765,1767,1769,1771,1773,1775,1777,1779,1781,1783,1785,1787,1789,1791,1793,1795,1797,1799,1801,1803,1805,1807,1809,1811,1813,1815,1817,1819,1821,1823,1825,1827,1829,1831,1833,1835,1837,1839,1841,1843,1845,1847,1849,1851,1853,1855,1857,1859,1861,1863,1865,1867,1869,1871,1873,1875,1877,1879,1881,1883,1885,1887,1889,1891,1893,1895,1897,1899,1901,1903,1905,1907,1909,1911,1913,1915,1917,1919,1921,1923,1925,1927,1929,1931,1933,1935,1937,1939,1941,1943,1945,1947,1949,1951,1953,1955,1957,1959,1961,1963,1965,1967,1969,1971,1973,1975,1977,1979,1981,1983,1985,1987,1989,1991,1993,1995,1997,1999,2001,2003,2005,2007,2009,2011,2013,2015,2017,2019,2021,2023,2025,2027,2029,2031,2033,2035,2037,2039,2041,2043,2045,2047,2049,2051,2053,2055,2057,2059,2061,2063,2065,2067,2069,2071,2073,2075,2077,2079,2081,2083,2085,2087,2089,2091,2093,2095,2097,2099,2101,2103,2105,2107,2109,2111,2113,2115,2117,2119,2121,2123,2125,2127,2129,2131,2133,2135,2137,2139,2141,2143,2145,2147,2149,2151,2153,2155,2157,2159,2161,2163,2165,2167,2169,2171,2173,2175,2177,2179,2181,2183,2185,2187,2189,2191,2193,2195,2197,2199,2201,2203,2205,2207,2209,2211,2213,2215,2217,2219,2221,2223,2225,2227,2229,2231,2233,2235,2237,2239,2241,2243,2245,2247,2249,2251,2253,2255,2257,2259,2261,2263,2265,2267,2269,2271,2273,2275,2277,2279,2281,2283,2285,2287,2289,2291,2293,2295,2297,2299,2301,2303,2305,2307,2309,2311,2313,2315,2317,2319,2321,2323,2325,2327,2329,2331,2333,2335,2337,2339,2341,2343,2345,2347,2349,2351,2353,2355,2357,2359,2361,2363,2365,2367,2369,2371,2373,2375,2377,2379,2381,2383,2385,2387,2389,2391,2393,2395,2397,2399,2401,2403,2405,2407,2409,2411,2413,2415,2417,2419,2421,2423,2425,2427,2429,2431,2433,2435,2437,2439,2441,2443,2445,2447,2449,2451,2453,2455,2457,2459,2461,2463,2465,2467,2469,2471,2473,2475,2477,2479,2481,2483,2485,2487,2489,2491,2493,2495,2497,2499,2501,2503,2505,2507,2509,2511,2513,2515,2517,2519,2521,2523,2525,2527,2529,2531,2533,2535,2537,2539,2541,2543,2545,2547,2549,2551,2553,2555,2557,2559,2561,2563,2565,2567,2569,2571,2573,2575,2577,2579,2581,2583,2585,2587,2589,2591,2593,2595,2597,2599,2601,2603,2605,2607,2609,2611,2613,2615,2617,2619,2621,2623,2625,2627,2629,2631,2633,2635,2637,2639,2641,2643,2645,2647,2649,2651,2653,2655,2657,2659,2661,2663,2665,2667,2669,2671,2673,2675,2677,2679,2681,2683,2685,2687,2689,2691,2693,2695,2697,2699,2701,2703,2705,2707,2709,2711,2713,2715,2717,2719,2721,2723,2725,2727,2729,2731,2733,2735,2737,2739,2741,2743,2745,2747,2749,2751,2753,2755,2757,2759,2761,2763,2765,2767,2769,2771,2773,2775,2777,2779,2781,2783,2785,2787,2789,2791,2793,2795,2797,2799,2801,2803,2805,2807,2809,2811,2813,2815,2817,2819,2821,2823,2825,2827,2829,2831,2833,2835,2837,2839,2841,2843,2845,2847,2849,2851,2853,2855,2857,2859,2861,2863,2865,2867,2869,2871,2873,2875,2877,2879,2881,2883,2885,2887,2889,2891,2893,2895,2897,2899,2901,2903,2905,2907,2909,2911,2913,2915,2917,2919,2921,2923,2925,2927,2929,2931,2933,2935,2937,2939,2941,2943,2945,2947,2949,2951,2953,2955,2957,2959,2961,2963,2965,2967,2969,2971,2973,2975,2977,2979,2981,2983,2985,2987,2989,2991,2993,2995,2997,2999,3001,3003,3005,3007,3009,3011,3013,3015,3017,3019,3021,3023,3025,3027,3029,3031,3033,3035,3037,3039,3041,3043,3045,3047,3049,3051,3053,3055,3057,3059,3061,3063,3065,3067,3069,3071,3073,3075,3077,3079,3081,3083,3085,3087,3089,3091,3093,3095,3097,3099,3101,3103,3105,3107,3109,3111,3113,3115,3117,3119,3121,3123,3125,3127,3129,3131,3133,3135,3137,3139,3141,3143,3145,3147,3149,3151,3153,3155,3157,3159,3161,3163,3165,3167,3169,3171,3173,3175,3177,3179,3181,3183,3185,3187,3189,3191,3193,3195,3197,3199,3201,3203,3205,3207,3209,3211,3213,3215,3217,3219,3221,3223,3225,3227,3229,3231,3233,3235,3237,3239,3241,3243,3245,3247,3249,3251,3253,3255,3257,3259,3261,3263,3265,3267,3269,3271,3273,3275,3277,3279,3281,3283,3285,3287,3289,3291,3293,3295,3297,3299,3301,3303,3305,3307,3309,3311,3313,3315,3317,3319,3321,3323,3325,3327,3329,3331,3333,3335,3337,3339,3341,3343,3345,3347,3349,3351,3353,3355,3357,3359,3361,3363,3365,3367,3369,3371,3373,3375,3377,3379,3381,3383,3385,3387,3389,3391,3393,3395,3397,3399,3401,3403,3405,3407,3409,3411,3413,3415,3417,3419,3421,3423,3425,3427,3429,3431,3433,3435,3437,3439,3441,3443,3445,3447,3449,3451,3453,3455,3457,3459,3461,3463,3465,3467,3469,3471,3473,3475,3477,3479,3481,3483,3485,3487,3489,3491,3493,3495,3497,3499,3501,3503,3505,3507,3509,3511,3513,3515,3517,3519,3521,3523,3525,3527,3529,3531,3533,3535,3537,3539,3541,3543,3545,3547,3549,3551,3553,3555,3557,3559,3561,3563,3565,3567,3569,3571,3573,3575,3577,3579,3581,3583,3585,3587,3589,3591,3593,3595,3597,3599,3601,3603,3605,3607,3609,3611,3613,3615,3617,3619,3621,3623,3625,3627,3629,3631,3633,3635,3637,3639,3641,3643,3645,3647,3649,3651,3653,3655,3657,3659,3661,3663,3665,3667,3669,3671,3673,3675,3677,3679,3681,3683,3685,3687,3689,3691,3693,3695,3697,3699,3701,3703,3705,3707,3709,3711,3713,3715,3717,3719,3721,3723,3725,3727,3729,3731,3733,3735,3737,3739,3741,3743,3745,3747,3749,3751,3753,3755,3757,3759,3761,3763,3765,3767,3769,3771,3773,3775,3777,3779,3781,3783,3785,3787,3789,3791,3793,3795,3797,3799,3801,3803,3805,3807,3809,3811,3813,3815,3817,3819,3821,3823,3825,3827,3829,3831,3833,3835,3837,3839,3841,3843,3845,3847,3849,3851,3853,3855,3857,3859,3861,3863,3865,3867,3869,3871,3873,3875,3877,3879,3881,3883,3885,3887,3889,3891,3893,3895,3897,3899,3901,3903,3905,3907,3909,3911,3913,3915,3917,3919,3921,3923,3925,3927,3929,3931,3933,3935,3937,3939,3941,3943,3945,3947,3949,3951,3953,3955,3957,3959,3961,3963,3965,3967,3969,3971,3973,3975,3977,3979,3981,3983,3985,3987,3989,3991,3993,3995,3997,3999,4001,4003,4005,4007,4009,4011,4013,4015,4017,4019,4021,4023,4025,4027,4029,4031,4033,4035,4037,4039,4041,4043,4045,4047,4049,4051,4053,4055,4057,4059,4061,4063,4065,4067,4069,4071,4073,4075,4077,4079,4081,4083,4085,4087,4089,4091,4093,4095,4097,4099,4101,4103,4105,4107,4109,4111,4113,4115,4117,4119,4121,4123,4125,4127,4129,4131,4133,4135,4137,4139,4141,4143,4145,4147,4149,4151,4153,4155,4157,4159,4161,4163,4165,4167,4169,4171,4173,4175,4177,4179,4181,4183,4185,4187,4189,4191,4193,4195,4197,4199,4201,4203,4205,4207,4209,4211,4213,4215,4217,4219,4221,4223,4225,4227,4229,4231,4233,4235,4237,4239,4241,4243,4245,4247,4249,4251,4253,4255,4257,4259,4261,4263,4265,4267,4269,4271,4273,4275,4277,4279,4281,4283,4285,4287,4289,4291,4293,4295,4297,4299,4301,4303,4305,4307,4309,4311,4313,4315,4317,4319,4321,4323,4325,4327,4329,4331,4333,4335,4337,4339,4341,4343,4345,4347,4349,4351,4353,4355,4357,4359,4361,4363,4365,4367,4369,4371,4373,4375,4377,4379,4381,4383,4385,4387,4389,4391,4393,4395,4397,4399,4401,4403,4405,4407,4409,4411,4413,4415,4417,4419,4421,4423,4425,4427,4429,4431,4433,4435,4437,4439,4441,4443,4445,4447,4449,4451,4453,4455,4457,4459,4461,4463,4465,4467,4469,4471,4473,4475,4477,4479,4481,4483,4485,4487,4489,4491,4493,4495,4497,4499,4501,4503,4505,4507,4509,4511,4513,4515,4517,4519,4521,4523,4525,4527,4529,4531,4533,4535,4537,4539,4541,4543,4545,4547,4549,4551,4553,4555,4557,4559,4561,4563,4565,4567,4569,4571,4573,4575,4577,4579,4581,4583,4585,4587,4589,4591,4593,4595,4597,4599,4601,4603,4605,4607,4609,4611,4613,4615,4617,4619,4621,4623,4625,4627,4629,4631,4633,4635,4637,4639,4641,4643,4645,4647,4649,4651,4653,4655,4657,4659,4661,4663,4665,46
```

Elbow Method using ine



I find it bit confusing what to choose as best K as value of loss keeps decreasing but if we see closely decrease after but too it is decreasing but there slope is not that steeper . Now ofcourse how many clu "granular" cluster we need where each cluster is telling some diiferent "aspect" about data. But for sak

```

best_Kmean = KMeans(n_clusters=8, init='k-means++')
best_model= best_Kmean.fit(svd_train)

cluster1,cluster2,cluster3,cluster4,cluster5,cluster6,cluster7,cluster8=[],[],[],[],[],[],[],[]
for i in range(best_model.labels_.shape[0]):
    if best_model.labels_[i] == 0:
        cluster1.append(X_no_stop_50k_reviews[i])
    elif best_model.labels_[i] == 1:
        cluster2.append(X_no_stop_50k_reviews[i])
    elif best_model.labels_[i] == 2:
        cluster3.append(X_no_stop_50k_reviews[i])
    elif best_model.labels_[i] == 3:
        cluster4.append(X_no_stop_50k_reviews[i])
    elif best_model.labels_[i] == 4:
        cluster5.append(X_no_stop_50k_reviews[i])
    elif best_model.labels_[i] == 5:
        cluster6.append(X_no_stop_50k_reviews[i])
    elif best_model.labels_[i] == 6:
        cluster7.append(X_no_stop_50k_reviews[i])
    #elif best_model.labels_[i] == 7:
        #cluster8.append(X_no_stop_50k_reviews[i])
    #elif best_model.labels_[i] == 8:
        #cluster9.append(X_no_stop_50k_reviews[i])
    else:
        cluster8.append(X_no_stop_50k_reviews[i])

# Number of reviews in different clusters
print("No. of reviews in Cluster-1 : ",len(cluster1))
print("\nNo. of reviews in Cluster-2 : ",len(cluster2))
print("\nNo. of reviews in Cluster-3 : ",len(cluster3))
print("\nNo. of reviews in Cluster-4 : ",len(cluster4))
print("\nNo. of reviews in Cluster-5 : ",len(cluster5))
print("\nNo. of reviews in Cluster-6 : ",len(cluster6))
print("\nNo. of reviews in Cluster-7 : ",len(cluster7))
print("\nNo. of reviews in Cluster-8 : ",len(cluster8))

```

```
#print("\nNo. of reviews in Cluster-9 : ",len(cluster9))  
#print("\nNo. of reviews in Cluster-10 : ",len(cluster10))
```

⇨ No. of reviews in Cluster-1 : 1992

No. of reviews in Cluster-2 : 1

No. of reviews in Cluster-3 : 1

No. of reviews in Cluster-4 : 1

No. of reviews in Cluster-5 : 2

No. of reviews in Cluster-6 : 1

No. of reviews in Cluster-7 : 1

No. of reviews in Cluster-8 : 1

- ▼ [5.5] Wordclouds of clusters obtained in the above section

```
#for cluster 1
data=''
for i in cluster1:
    data+=str(i)
from wordcloud import WordCloud
wordcloud = WordCloud(background_color="white").generate(data)

# Display the generated image:
plt.imshow(wordcloud, interpolation='bilinear')
plt.axis("off")
plt.show()
```

2



```
#for cluster 2
data=''
for i in cluster2:
    data+=str(i)
from wordcloud import WordCloud
wordcloud = WordCloud(background_color="white").generate(data)

# Display the generated image:
plt.imshow(wordcloud, interpolation='bilinear')
plt.axis("off")
plt.show()
```



```
#for cluster 3
data=''
for i in cluster3:
    data+=str(i)
from wordcloud import WordCloud
wordcloud = WordCloud(background_color="white").generate(data)
```

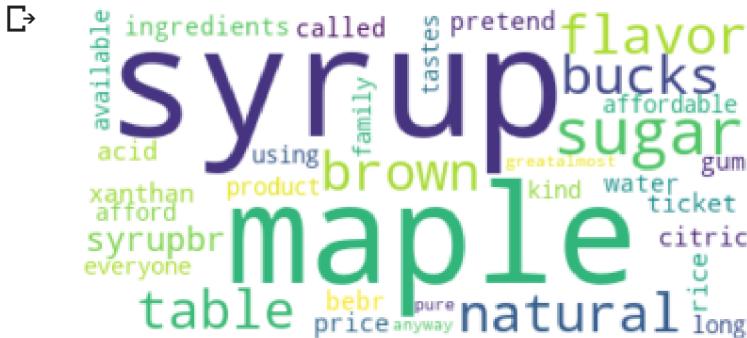
```
# Display the generated image:  
plt.imshow(wordcloud, interpolation='bilinear')  
plt.axis("off")  
plt.show()
```



```
#for cluster 4
data=''
for i in cluster4:
    data+=str(i)
```

```
data = sci(1)
from wordcloud import WordCloud
wordcloud = WordCloud(background_color="white").generate(data)

# Display the generated image:
plt.imshow(wordcloud, interpolation='bilinear')
plt.axis("off")
plt.show()
```



```
#for cluster 5
data=''
for i in cluster5:
    data+=str(i)
from wordcloud import WordCloud
wordcloud = WordCloud(background_color="white").generate(data)
```

```
# Display the generated image:  
plt.imshow(wordcloud, interpolation='bilinear')  
plt.axis("off")  
plt.show()
```



```
#for cluster 6
data=''
for i in cluster6:
    data+=str(i)
from wordcloud import WordCloud
wordcloud = WordCloud(background_color="white").generate(data)
```

```
# Display the generated image:
plt.imshow(wordcloud, interpolation='bilinear')
plt.axis("off")
plt.show()
```



```
#for cluster 7
data=''
for i in cluster7:
    data+=str(i)
from wordcloud import WordCloud
wordcloud = WordCloud(background_color="white").generate(data)
```

```
# Display the generated image:
plt.imshow(wordcloud, interpolation='bilinear')
plt.axis("off")
plt.show()
```



```
#for cluster 8
data=''
for i in cluster8:
    data+=str(i)
from wordcloud import WordCloud
wordcloud = WordCloud(background_color="white").generate(data)
```

```
# Display the generated image:
plt.imshow(wordcloud, interpolation='bilinear')
plt.axis("off")
plt.show()
```



- ▼ [5.6] Function that returns most similar words for a given word.

```
from sklearn.metrics.pairwise import cosine_similarity
def similar_words(word,n):
    top_words=[]
    cosine_sim=cosine_similarity(occ_matrix)
    val_word=cosine_sim[lst_2000.index(word)]
    idx=np.argsort(val_word)
    for i in range(n):
        top_words.append(lst_2000[idx[i]])
    return top_words

print('So the Top 10 words similar to word sucrose are ')
print(similar_words('sucrose',10))
print ("#####")
print('So the Top 10 words similar to word arsenic are ')
print(similar_words('arsenic',10))
print ("#####")
```

```
↳ So the Top 10 words similar to word sucrose are
['cocaine', 'naked', 'tea tried', 'cant taste', 'found great', 'meatloaf', 'hair proc
#####
So the Top 10 words similar to word arsenic are
['cocaine', 'hiding', 'product always', '5050', 'like ones', 'counta', 'fiber proteir
#####
```

- ▼ [6] Conclusions

We have done **k=8** for this assignment as mentioned earlier we could have tried many other number of requirements and how smaller chunks of cluster we want with how much granulaity.

In this assignment we took top 2000 features only.

As far as finding of similar word is concerned i have used cosine similarity and we see to some extent when i tried similar word for arsenic i got 'cocaine' , 'hiding' , 'cant taste' which kind of makes sense.