

## OverAll Procedure-

- Step 1:** Crawling over the Internet, and saving information about retail companies into training.txt file.
- Step 2:** Cleaning the data-The data is in form of paragraphs and sentences. It is tokenized and from that only those tokens are chosen which potentially contains information about the vendor. For example, only the Nouns are kept and words which are tagged under Name, Organization, Location or Entity are kept.
- Step 3:** At this step, the cleaned words are changed into vectors via Word2Vec. Word2Vec changes each word into a vector containing an array of numbers such that mathematical computations like clustering can be done on textual data. Now these words are used for forming Clusters using K Means Clustering.
- Step 4:** Now we input a vendor (whose details we need to find out), We crawl over the training.txt file and extract only those paragraphs which are pertaining to the input vendor. We do it via NER Tagger. We check the Organization for each paragraph in our training set and in which paragraphs the (Organization==Input Vendor Name) and if this has happened more than thrice it is considered as a suitable paragraph for the Vendor. This is our Test Data
- Step 5:** Now we use our Kmeans model created in Step 3 and determine to which cluster the words in our Test Data belongs to. Thus we create a Cluster frequency matrix like Cluster 1:5 words, Cluster 2:1 words, Cluster 3:0 words. Now we select the top 7 words (which has occurred maximum number of times in the Test Data) from the cluster containing maximum number of words and top 3 words from the 2nd max cluster.

## Code Explained in Details with Screenshot of Outputs displayed, ran on small sample of Data and Output predicted for Vendor Dell:-

### 1. **crawling.py** (Used for Web Crawling to create the Train Data)

```
from bs4 import BeautifulSoup #BeautifulSoup is a library for data crawling
import urllib

#Paste the URL here which you want to crawl the data from
thisurl = "https://en.wikipedia.org/wiki/Hewlett-Packard"
soup = BeautifulSoup(urllib.urlopen(thisurl).read(), 'html.parser')

#Find the relevant data in the document
data="".join([x.text for x in soup.find(id='mw-content-text').find_all('p')]).encode('utf-8');

#File is opened in Append mode so that it is appended everytime a URI is crawled
with open("trainingdata.txt", "a") as myfile:
    myfile.write(data)
myfile.close()
```

## Output: (training.txt)

```
acturers. The company was well known for its innovations in supply chain management and electro manufacturing\delivering individual PCs configured to customer specifications.[6][7] Dell entered the market for IT services. The company has since made additional acquisitions in serving complete solutions for enterprise customers.[8][9]', u'Dell was listed at number 51 in total information prevents the company from being ranked by Fortune. In 2015, it was the third largest in the world.[12] Dell is the sixth largest company in Texas by total revenue, according to Fortune magazine in the Greater Austin area.[14] It was a publicly traded company (NASDAQ: DELL), as well as until October 30, 2013.', u'On September 7, 2016, Dell Inc. acquired EMC Corporation to form Dell Technologies its origins to 1984, when Michael Dell created Dell Computer Corporation, which at the time was a quartered company sold IBM PC-compatible computers built from stock components.[18] Dell dropped out of his family. In 1985, the company produced the first computer of its own design, the Turbo PC, which was sold to consumers and custom assembled each ordered unit according to a selection of options. The Lee Walker, a 51-year-old venture capitalist, as president and chief operating officer, to serve as recruiting members to the board of directors when the company went public in 1988. Walker retired Dell Systems to transform the company from a fast-growing medium-sized firm into a billion-dollar company.
```

## 2. **cleaning.py** (Used for cleaning the paragraphs and changing them into useful tokens)

```
import pandas as pd
from nltk.stem.porter import PorterStemmer
from nltk.tokenize import word_tokenize,sent_tokenize
from nltk.tag import pos_tag
from nltk.tag import StanfordNERTagger
import nltk.data
import re
```

### **#For reading sentences from the training module and saving it to a list Query1**

```
query1=[]
with open("Desktop/trainingdata.txt", "r") as myfile:
    query1.append(myfile.read())
myfile.close()
```

### **#Tokenizing sentences into words and removing punctuations**

```
def remove_punctuations(text):
    wordlist=[]
    sentences=sent_tokenize(text)
    for sentence in sentences:
        words=[word for word in sentence.split()]
        words=[re.sub(r'[^\w\s]','',word) for word in words]
        wordlist.append(words)
    return wordlist
```

### **#Keeping only Nouns from the sentence**

```
def descriptive_words(words):
```

```

tagged_word=[]
meaningful_words=[]
tags=['NN','NNS','NNP']

for word in words:
    tagged_word.append(pos_tag(word))
    for words in tagged_word:
        for word in words:
            if word[1] in tags:
                meaningful_words.append(word[0])

return meaningful_words

```

### **#Using NER Tagger to find Name, Organisation, Entity**

```

def remove_names(text):
    meaningful_words=[]
    tagged_word=[]
    tags=['LOCATION', 'ORGANIZATION', 'PERSON']
    st = StanfordNERTagger('english.all.3class.distsim.crf.ser.gz')

    for sentence in text:
        tagged_word.append(st.tag(sentence))

    for words in tagged_word:
        for word in words:
            if word[1] in tags:
                meaningful_words.append(word[0])
    return meaningful_words

```

### **#Using the above functions for Cleaning of Sentences in Query1**

```
words_req=[]
```

```

for query in query1:
    #Removing Punctuations and changing to lower space
    letters_only = remove_punctuations(query)

```

```

#5 Read only NOuns,Pronouns,interjections (descriptive words)
meaningful_words=descriptive_words(letters_only)

```

```

#6 Keeping Time, Location, Organization, Person, Money, Percent, Date using NER
removed_words=remove_names(letters_only)

```

```

#6.We would not Stem our data because we would convert our data into vectors in the next
    step, Stemming them would loose the information saved in them
    stemmed_words=[words for words in meaningful_words+removed_words]
    words_req.append(stemmed_words)

```

**Output** - After Cleaning every paragraph changes into the form of

```
[['HewlettPackard', 'Company', 'HP', 'information', 'technology', 'company', 'Palo', 'Alto', 'enterprises', 'customers', 'government', 'health', 'education', 'sectors', 'company', 'onecar', 'st', 'equipment', 'HP', 'worlds', 'PC', 'manufacturer', 'Lenovo', 'HP', 'data', 'storage', 'ha', 'stry', 'standard', 'servers', 'storage', 'devices', 'products', 'software', 'range', 'printers', 'officesupply', 'retailers', 'software', 'partners', 'technology', 'vendors', 'HP', 'service', 'ent', 'instruments', 'part', 'business', 'Agilent', 'merger', 'Compaq', 'acquisition', 'EDS', 'buyout', 'Palm', 'Inc', 'September', 'HP', 'bidding', 'var', 'share', 'offer', 'Dell', 'Octob', 'split', 'November', 'companies', 'HP', 'Inc', 'Hewlett', 'Packard', 'Enterprise', u'Hewlett', 'Hewlett', u'David', u'Dave', u'Packard', u'HP', u'Lenovo', u'HP', u'HP', u'HewlettPackard', u', u'Packard', u'Enterprise']]
```

#### 4. **clustering.py** (Changing words into vectors and then applying K Means Clustering on them)

```
from gensim import models
from sklearn.cluster import Means
```

##### **#Loading Googles PreTrained Vector Model**

```
wordvec = models.Word2Vec.load_word2vec_format('GoogleNews-vectors-      nega
tive300.bin', binary=True)
vecs=[]
```

##### **#Changing the words into Vectors so that Clusters can be formed out of them**

```
for sentence in words_req:
    for words in sentence:
        word=wordvec[words]
        vecs.append(word)
```

##### **#Applying Clustering to it**

```
kmeans = KMeans(init='k-means++', n_clusters=10, n_iter=100)
kmeans.fit(vecs)
```

#### **Output-**

Word converted to a vector ( this is a vector for a single word)

```
[ 0.00524902 -0.14355469 -0.06933594  0.12353516  0.13183594 -0.08886719
-0.07128906 -0.21679608 -0.19726562  0.05566406 -0.07568359 -0.30805938
 0.10400391 -0.00081635  0.1328125   0.11279297  0.07275391 -0.046875
 0.06591797  0.09423028  0.19042969  0.13671075 -0.23632012 -0.11065234
 0.06542969 -0.05322266 -0.30859375  0.09179688  0.18847656 -0.16699219
-0.15625    -0.13085938 -0.08251953  0.21289062 -0.35546875 -0.13183594
 0.09619141  0.26367108 -0.09472656  0.10359375  0.10693359 -0.41601562
 0.26953125 -0.02770996  0.17578125 -0.11279297 -0.00411987  0.14550781
 0.15625    0.26757812 -0.01794434  0.09863281  0.05297852 -0.03125
-0.16308594 -0.05810547 -0.34375    -0.17285156  0.11425781 -0.09033203
 0.13476562  0.27929688 -0.04980469  0.12988281  0.17578125 -0.22167969
-0.01190186  0.140625   -0.18164062  0.11065234  0.16113281  0.21404375
-0.21191406  0.12695312 -0.10009766  0.13671075  0.12695312  0.01531982
 0.10449219 -0.02783203 -0.06030273  0.0222168  0.10164062 -0.06730281
 0.04907227  0.15429688 -0.25        0.13964844  0.29492188  0.10644531
 0.3359375   -0.22265625 -0.125      -0.05297852  0.19238281  0.06835938
 0.06902422 -0.05200195  0.14453125  0.00440608 -0.01013104 -0.1404375
 0.21777344 -0.1953125   -0.390625   0.07763672 -0.57421875 -0.07910156
-0.04052734 -0.1075     0.25390625  0.15722656  0.125      0.140625]
```

Clusters formed on these vectors

```
Clustering Done
KMeans(copy_x=True, init='k-means++', max_iter=300, n_clusters=3, n_init=10,
       n_jobs=1, precompute_distances='auto', random_state=None, tol=0.0001,
       verbose=0)
```

5. **predict.py** (First find out the paragraphs which points out to your input vendor as your train data, then identify to which clusters your train data belongs to.)

```
from gensim import models
from sklearn.cluster import KMeans
import cleaning.py
```

**#Searching for only those Sentences from the training module which belongs to the vendor we are searching for DELL in this case and save it in Query3**

```
query2=[]
with open("Desktop/trainingdata.txt", "r") as myfile:
    query2.append(myfile.read())
myfile.close()
```

```
query3=[]
for query in query2:
    letters_only = remove_punctuations(query)
    removed_words=remove_names(letters_only)
    if "DELL" in removed_words:
        query3.append(query)
```

**#Now Predicting to which Cluster the words describing the Vendor belongs to**

```
cluster_data={'0':0,'1':0,'2':0,'3':0,'4':0,'5':0,'6':0,'7':0,'8':0,'9':0}
cluster_words={'0':[],'1':[],'2':[],'3':[],'4':[],'5':[],'6':[],'7':[],'8':[],'9':[]}
```

```
for pos,vec in enumerate(vecs):
    cluster_no=kmeans.predict([vec])
    cluster_data[str(cluster_no[0])]+=1
    cluster_words[str(cluster_no[0])].append(words[pos])
```

**#find cluster which contains maximum number of words, print the top 7 (which occurred the most in the training set for that vendor) from that cluster**

```
max_cluster=max(cluster_data, key=cluster_data.get)
imp_words_1=set(cluster_words[max_cluster])
```

**#Findind the 2nd max cluster,print the top 3 (which occurred the most in the training set for that vendor) from that cluster**

```
cluster_data[max_cluster]=0
max_cluster=max(cluster_data, key=cluster_data.get)
imp_words_2=set(cluster_words[max_cluster])
```

## Output-

Cluster Matrix formed

```
Cluster matrix formed
{'1': 5, '0': 0, '3': 0, '2': 5, '5': 0, '4': 0, '7': 0, '6': 0, '9': 0, '8': 0}
{'1': [u'deal', u'Hardwar', u'Softwar', u'Hardwar', u'Softwar'], '0': [], '3': [], '2': [u'Dell', u'Laptop', u'
set([u'Hardwar', u'Softwar', u'deal'])
```

**Final Predicted Output for Vendor name== “DELL”**

```
set([u'Hardwar', u'Softwar', u'deal'])
set([u'Laptop', u'Dell'])
```