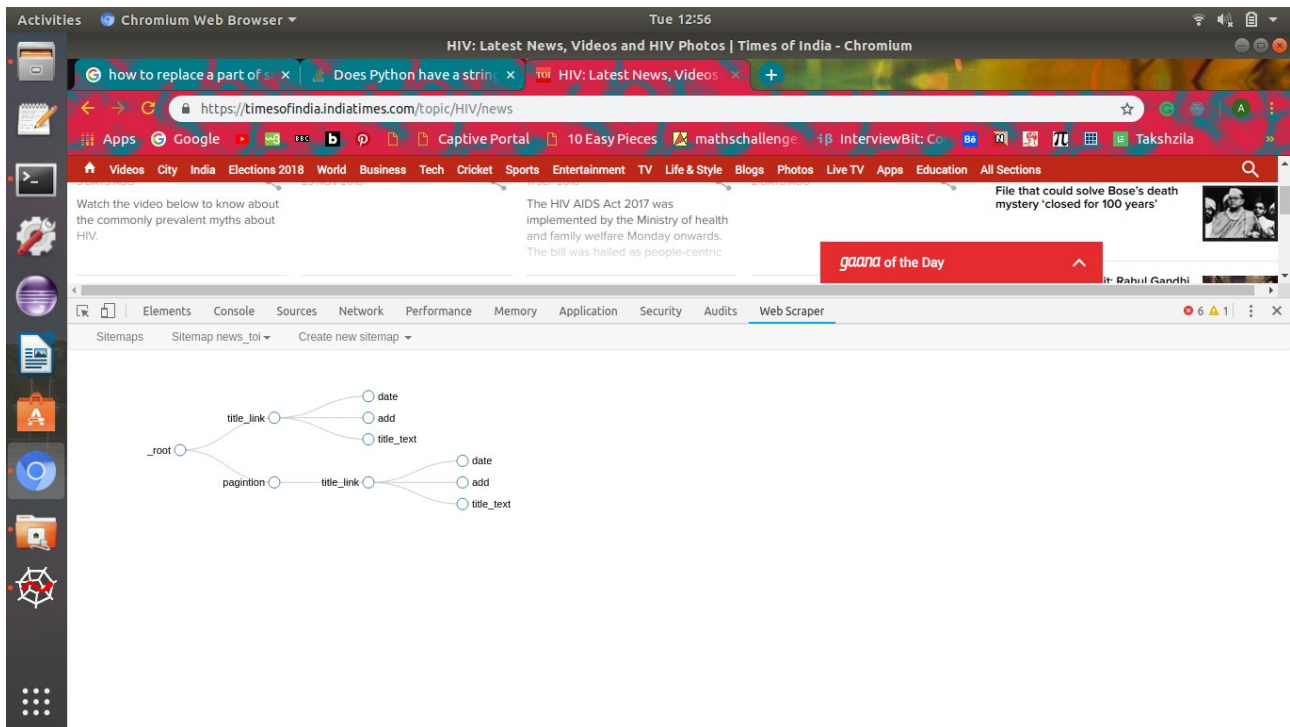# Assignment Report

## 1. Web Scraping using Chrome Extension

I have attached the screenshot of the sitemap selector graph for web scraping.



After this I have exported the file in CSV (Comma Separated Values) format for further analysis.

The extracted CSV file is shown below.



After this step I extracted the articles using BeautifulSoup Web scraping module in Python.

## 2. Analysing data

The average word lenght of articles is 383 words with a maximum of 1782 words. The following graph shows an analysis on word length of articles.



The NLTK tree for an article is shown below.

# Frequency distribution for city names mentioned in articles with their legends.



| column | | | | | | | | | ▼ |
|---|---|---|---|---|---|---|---|---|---|
| Data | | | | | | | | | |

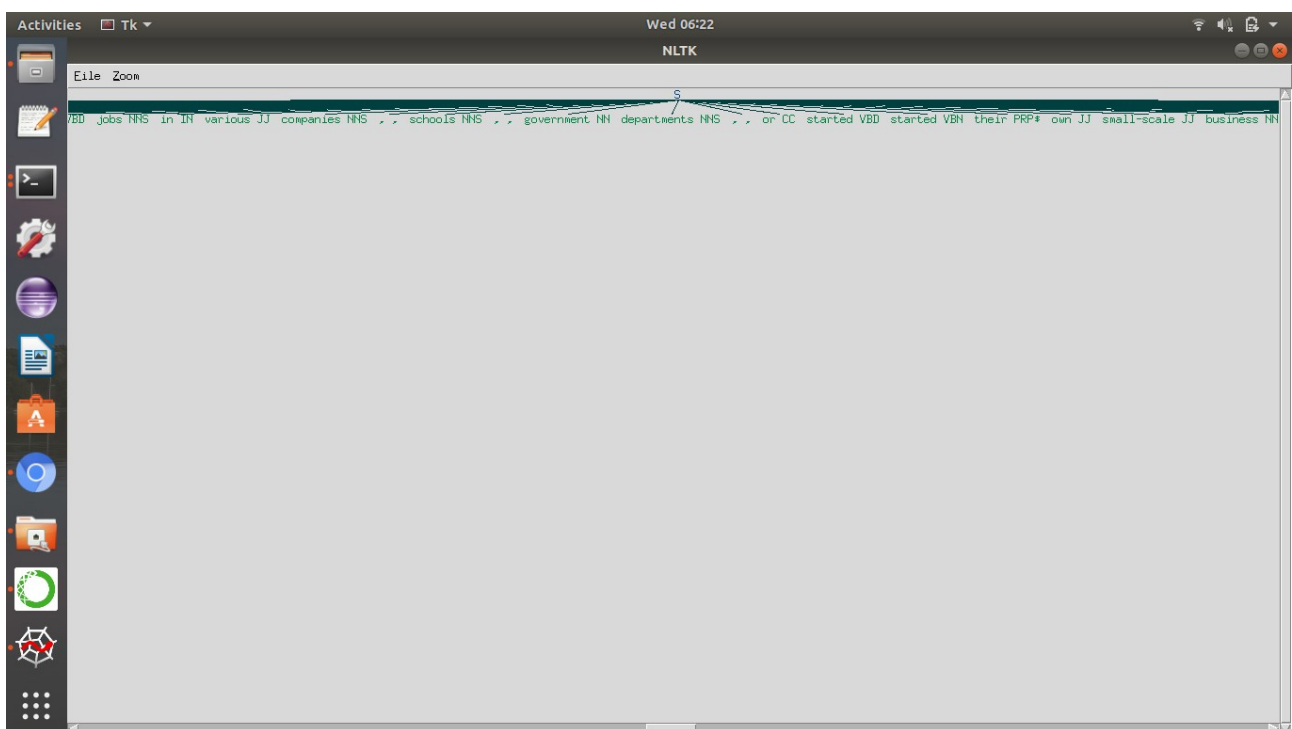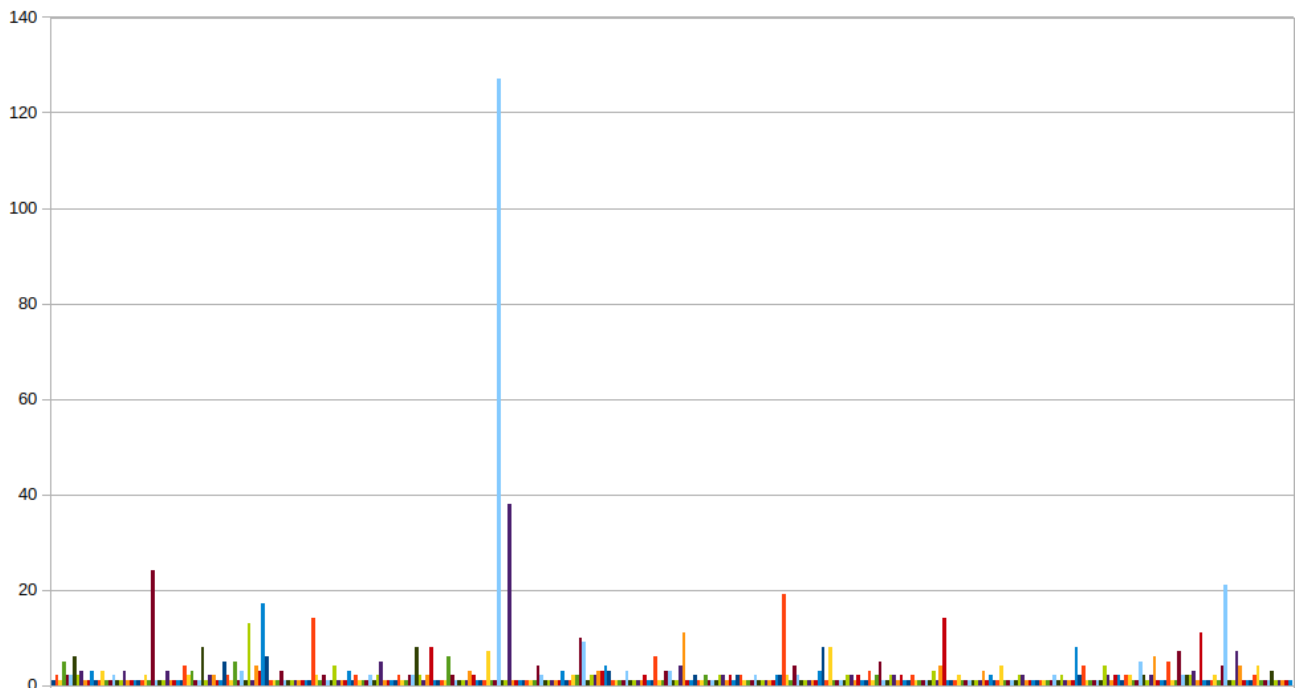| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| ■ Aadhaar | ■ Abbupura | ■ Achhnera | ■ Africa | ■ Agasaim | ■ Agra | ■ Ahmedabad | ■ x2 | ■ American | ■ Amritsar |
| ■ Anand | ■ Anandaghar | ■ Anita | ■ Ankit | ■ Anthiyur | ■ Anupurba | ■ Architecture | ■ Asia | ■ Asian | ■ Asothar |
| ■ Assam | ■ Atrocities | ■ August | ■ Bagalkot | ■ Bagh | ■ Balasubramanian | ■ Baltimore | ■ Bandra | ■ Bangarmau | ■ Bannur |
| ■ Bansal | ■ Banyan | ■ Bareilly | ■ Basti | ■ Beijing | ■ Being | ■ Belagavi | ■ Bengal | ■ Bengali | ■ Bengaluru |
| ■ Bhadohi | ■ Bharuch | ■ Bihar | ■ x1 | ■ Breach | ■ British | ■ Byculla | ■ Calcutta | ■ Canada | ■ Canadian |
| ■ Care | ■ Cauvery | ■ Centre | ■ Chandigarh | ■ Chandrapur | ■ Chennai | ■ Chikkannan | ■ Chikkaran | ■ Child | ■ China |
| ■ Chinese | ■ Chinmay | ■ Chinna | ■ Chitorgarh | ■ Christian | ■ Corporation | ■ Cortalim | ■ Crime | ■ Cuddalore | ■ Darunavir |
| ■ Davanagere | ■ Deepak | ■ Deepu | ■ Delhi | ■ Diwali | ■ Dumbledore | ■ Education | ■ Egmore | ■ Egra | ■ ELISA |
| ■ Eminence | ■ England | ■ Epidemiology | ■ Erode | ■ Europe | ■ Faizabad | ■ Farrukhabad | ■ Farukhabad | ■ Fatehgarh | ■ Fatehpur |
| ■ Fatima | ■ Florida | ■ Food | ■ France | ■ French | ■ FY19 | ■ Galla | ■ Ganiga | ■ Gautam | ■ Geneva |
| ■ Ghaziabad | ■ Gilead | ■ Goa | ■ Gorakhpur | ■ Ground | ■ Grover | ■ Gujarat | ■ Gujjar | ■ Guru | ■ HAART |
| ■ Halappa | ■ Haryana | ■ Hassan | ■ Health | ■ Hence | ■ Himachal | ■ Himahcal | ■ Hindi | ■ Hindustan | ■ Hollywood |
| ■ Howrah | ■ Hunsur | ■ Hyderabad | ■ Ichalkaranji | ■ Independence | ■ India | ■ India.Robotic | ■ India.The | ■ Indian | ■ Indonesia |
| ■ Industry | ■ Iraq | ■ Isanpur | ■ Jaipur | ■ Jalandhar | ■ Jalna | ■ Jamakhandi | ■ Jandiala | ■ Jharkhand | ■ Jodhewal |
| ■ July | ■ Junagadh | ■ Jurists | ■ Kaithal | ■ Kalina | ■ Kalyani | ■ Kannauj | ■ Kannur | ■ Kanpur | ■ Karnataka |
| ■ Karol | ■ Karur | ■ Kasaragod | ■ Kerala | ■ Kharar | ■ Kirmidiyapur | ■ Kolhapur | ■ Kolkata | ■ Kollidam | ■ Kopargaon |
| ■ Kota | ■ Kozhikode | ■ Krishnarajasagar | ■ Kukatpally | ■ Kumar | ■ Lancet | ■ Latur | ■ x | ■ Lohara | ■ Lok |
| ■ London | ■ Lucknow | ■ Ludhiana | ■ Madarpura | ■ Madhapur | ■ Madiaon | ■ Madurai | ■ Maharashtra | ■ Mahindra | ■ Malad |
| ■ Malappuram | ■ Malaria | ■ Malleswaram | ■ Malsian | ■ Management | ■ Mandi | ■ Manipur | ■ Manitoba | ■ Manmad | ■ Marathwada |
| ■ Maryland | ■ Mathura | ■ Maurya | ■ Mayyil | ■ Median | ■ Medicine | ■ Meena | ■ Meerut | ■ Meghalaya | ■ Men |
| ■ Microbiology | ■ Midnapore | ■ Modi | ■ Mombasa | ■ Mukkombu | ■ Mumbai | ■ Muneeshwar | ■ Murshidababd | ■ Mysuru | ■ Nadu |
| ■ Nagaland | ■ Nagapattinam | ■ Nair | ■ Namakkal | ■ Nangloi | ■ Narsinghpur | ■ Nashik | ■ Neelam | ■ New | ■ NHAI |
| ■ Nipah | ■ Niphad | ■ Nobel | ■ Noida | ■ North | ■ north-eastern | ■ Nursing | ■ Oceanography | ■ Orissa | ■ Osmania |
| ■ Oxford | ■ Padikattugal | ■ Palawi | ■ Palayamkottai | ■ Palghar | ■ Parmanand | ■ Parra | ■ Parvathy | ■ Pathanapuram | ■ Patna |
| ■ Periyapatna | ■ Persons | ■ PGIMER | ■ Phagwara | ■ Pharmacopoeia | ■ Pollachi | ■ Potato | ■ Pradesh | ■ Prayas | ■ Premganj |
| ■ Punjab | ■ Puri | ■ Q3 | ■ Queer | ■ Rajasthan | ■ Raltegravir | ■ Rihan | ■ Ritonavir | ■ Rivona | ■ Rohini |
| ■ Rohtak | ■ Roorkee | ■ Russia | ■ Russian | ■ Russians | ■ Sabang | ■ Sabha | ■ Sachin | ■ Sahnewal | ■ Sahodaran |
| ■ Sainikpuri | ■ Salem | ■ Salempur | ■ Sanskrit | ■ Sarai | ■ Science | ■ Secunderabad | ■ Sengulam | ■ Sharma | ■ Sheen |
| ■ Shipurkar | ■ Shiv | ■ Shivajinagar | ■ Sikkaran | ■ Sikkim | ■ Sion | ■ Soundaravel | ■ South | ■ Sriram | ■ States |
| ■ Std | ■ Subramanian | ■ Sudarshan | ■ Sumana | ■ Sumant | ■ Surat | ■ Swaddi | ■ Swami | ■ Swamy | ■ Sweden |
| ■ tailor-made | ■ Tambaram | ■ Tamil | ■ Tapentadol | ■ Technical | ■ Telangana | ■ Telangana.State | ■ Texas | ■ Thai | ■ Thailand |
| ■ Thakur | ■ Thakurdwara | ■ Thana | ■ Thane | ■ Thanjavur | ■ Theni | ■ Thiruvananthapuram | ■ Tibba | ■ Tirwa | ■ Tohana |
| ■ Toronto | ■ Treg | ■ Trichy | ■ Tripura | ■ Tuberculosis | ■ Udai | ■ Udaipur | ■ Udupi | ■ United | ■ Unnao |
| ■ UP | ■ Upperveda | ■ Uttar | ■ Uttarakhand | ■ Vadodara | ■ Varanasi | ■ Vellore | ■ Vihaan | ■ Vijayapura | ■ Vinay |
| ■ Vivekananda | ■ Warora | ■ West | ■ Witwatersrand | ■ World | ■ Xinhua | ■ Yeshwantpur | ■ Yousufguda | | |

These named entities are extracted using the IO_tag and parts of speech(POS) tag in NLTK.

## 3. Modules Used:

1. BeautifulSoup4 Module(bs4)
2. Pandas
3. Requests
4. JSON Module
5. Natural Language Toolkit (NLTK)
6. os Module
7. Matplotlib

## 4. Python Scripts

## For Data Analysis

```
import requests
import pandas as pd
import bs4
import json

url_list=[]

field=["title_link-href"]

data=pd.read_csv(r"/home/arushi/
news_toi.csv",skipinitialspace=True,usecols=field,index_col=None)

list1=data.iloc[:,0].tolist()

#print(dic[:30])


title=[]
d_time=[]
content=[]

print(len(list1))
print("\n\n")



for i in range(len(list1)):
    res=requests.get(list1[i])

    print(i)
    print("------****res.text length****----\n")
    print(len(res.text))
    print("\n")
    soup=bs4.BeautifulSoup(res.text)
```

```
    ele_text=soup.select("div .Normal")
    ele_date=soup.select("time")
    ele_title=soup.select("arttitle")
    #print(ele_date)
    #print(ele_title)
  # print(type(ele_text))
    print(ele_date[0].attrs)
    print(ele_title[0].getText())
    print(ele_text[0].getText())
    print("--------new article------\n\n")
    #append lists
    title.append(ele_title[0].getText())
    d_time.append(ele_date[0].attrs)
    content.append(ele_text[0].getText())
    x=title[i-1152]
    if "HIV/AIDS" in x:
        y=x.replace("HIV/AIDS","HIV-AIDS")
    toi2=open(r"/home/arushi/toi_news_articles/%s"%y,"w")
    toi2.write(json.dumps(d_time[i])+"\n")
    toi2.write(content[i])



print(len(title))
print(title)
print("\n")

print(len(d_time))
print(d_time)
print("\n")

print(len(content))
print(content)
print("\n")

new_file=[title,d_time,content]
zip(*new_file)
print(zip(*new_file))

for t,d,c in zip(*new_file):
    print(t)
    print(d)
    print(c)


x=title[0]
print(content[0])
toi2=open(r"/home/arushi/toi_news_articles/%s"%x,"w")
toi2.write(json.dumps(d_time[0])+"\n")
toi2.write(content[0])
```

```python
#str1=data[1:1]
#print(str1+"hello")




res=requests.get(r"https://timesofindia.indiatimes.com/india/hiv-positive-
women-marries-twice-kills-first-husband-with-the-help-of-second/articleshow/
63043159.cms")

print(type(res))

soup=bs4.BeautifulSoup(res.text)
print(type(soup))


ele_text=soup.select("div .Normal")
ele_date=soup.select("time")
print(ele_date)
ele_title=soup.select("arttitle")
print(ele_title)

print(type(ele_text))

print(ele_text[0].getText())


import nltk
from nltk.tokenize import sent_tokenize, word_tokenize
from nltk import ne_chunk, pos_tag
import os
import matplotlib.pyplot as plt

plt.style.use('ggplot')
from nltk import conlltags2tree, tree2conlltags


#from spacy.en import English



#new fuction definition
def entities(text):
    return ne_chunk(
        pos_tag(
            word_tokenize(text)))

#nltk.download('words')
```

```python
#nltk.download('punkt')
#nltk.download('maxent_ne_chunker')
#nltk.download('averaged_perceptron_tagger')

path=r"/home/arushi/toi_news_articles"

len_art=[]
city_name=[]

for filename in os.listdir(path):
    print(filename)
    toi2=open(r"/home/arushi/toi_news_articles/"+filename,"r")
    data=toi2.read().replace('\n', '')
    #len_art.append(len(data.split()))
    words=word_tokenize(data)

    #print(nltk.pos_tag(words))

    tree=entities(data)
    iob_tags = tree2conlltags(tree)
    #print(iob_tags)

    for tup in iob_tags:
        if(tup[2]=="B-GPE" or tup[2]=="O_GPE" or tup[2]=="I-GPE"):
            city_name.append(tup[0])

    #print(tree)
    #tree.draw()



print(city_name)

import pandas as pd
df = pd.DataFrame(city_name, columns=["commmn"])
df.to_csv('city_list.csv', index=False)


city_set=set(city_name)
word_tag_fd=nltk.FreqDist(words)


plt.hist(city_name)

fig_size = plt.rcParams["figure.figsize"]
fig_size[0] = 12
fig_size[1] = 9
plt.rcParams["figure.figsize"] = fig_size
plt.xticks(city_set)
plt.show()
```

```
'''
print(len_art)
plt.plot(len_art)
plt.ylabel('article length')
plt.show()

print(max(len_art))
print(min(len_art))
print(sum(len_art)/len(len_art))
'''
```

**Neural Network Model**

```
import pandas as pd
import numpy as np
import pickle
from keras.preprocessing.text import Tokenizer
from keras.models import Sequential
from keras.layers import Activation, Dense, Dropout
from sklearn.preprocessing import LabelBinarizer
import sklearn.datasets as skds
from pathlib import Path

# For reproducibility
np.random.seed(1237)

# Source file directory
path_train = "/home/arushi/toi_news_articles/'85% Raj rural women don't know about HIV-AIDS'"

files_train = skds.load_files(path_train,load_content=False)

label_index = files_train.target
label_names = files_train.target_names
labelled_files = files_train.filenames

data_tags = ["filename","category","news"]
data_list = []

# Read and add data from file to a list
i=0
for f in labelled_fmodel = Sequential()
model.add(Dense(512, input_shape=(vocab_size,)))
model.add(Activation('relu'))
model.add(Dropout(0.3))
model.add(Dense(512))
model.add(Activation('relu'))
model.add(Dropout(0.3))
model.add(Dense(num_labels))
model.add(Activation('softmax'))
model.summary()
```

```
model.compile(loss='categorical_crossentropy',
        optimizer='adam',
        metrics=['accuracy'])

history = model.fit(x_train, y_train,
            batch_size=batch_size,
            epochs=30,
            verbose=1,
            validation_split=0.1)iles:
    data_list.append((f,label_names[label_index[i]],Path(f).read_text()))
    i += 1

# We have training data available as dictionary filename, category, data
data = pd.DataFrame.from_records(data_list, columns=data_tags)

encoder = LabelBinarizer()
encoder.fit(train_tags)
y_train = encoder.transform(train_tags)
y_test = encoder.transform(test_tags)
```