# Problem Set 3 - Natural Language Processing (UET)

https://github.com/CoderHung/Problem-set-3-Natural-language-processing-class-UET

## Overview

This project explores five topics:

- **Science**
- **Politics**
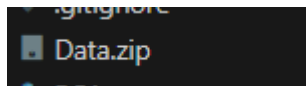- **Nature**
- **Sports**
- **Fashion**

This project uses the term-document matrix, TF-IDF, and t-SNE to visualize word vectors.

## Files

- **Term-Document Matrix**: `TD.csv`
- **TF-IDF**: `TF-IDF.csv`
- **t-SNE Visualization**: 150 random word vectors from TD matrix projected on a 2D plane

## Instructions

### Step 1

Unzip the dataset:



### Step 2

run term-document matrix.py to generate the term-document matrix:

```python
import os
import pandas as pd
import numpy as np
from sklearn.feature_extraction.text import CountVectorizer
import nltk
from nltk.corpus import words
import re

#nltk.download('words')
english_words = set(words.words())

def read_text_files(directory):
    documents = []
    filenames = []

    for filename in os.listdir(directory):
        if filename.endswith('.txt'):
            filepath = os.path.join(directory, filename)
            with open(filepath, 'r', encoding='utf-8') as file:
                text = file.read()
                words = [word for word in re.findall(r'\b\w+\b', text.lower()) if word in english_words]
                words = ' '.join(words)
                documents.append(words)
                filenames.append(filename)

    filenames = sorted(filenames, key=lambda x: int(x.split(' - ')[0]))
    filenames = [entry.split(' - ')[1] for entry in filenames]
    return documents, filenames


def create_term_document_matrix(documents):
    vectorizer = CountVectorizer(strip_accents='unicode',min_df = 0.03,stop_words='english')
    tdm = vectorizer.fit_transform(documents).toarray()
    tdm = np.transpose(tdm)
    return tdm, vectorizer.get_feature_names_out()

directory = 'Data'
documents, filenames = read_text_files(directory)
tdm, terms = create_term_document_matrix(documents)
tdm_df = pd.DataFrame(tdm ,columns=filenames, index=terms)
print(tdm_df.info)
tdm_df.to_csv('TD.csv')
```

## Step 3

run tf-idf.py to compute TF-IDF matrix:

```python
import os
import pandas as pd
import numpy as np
from sklearn.feature_extraction.text import TfidfVectorizer
import nltk
from nltk.corpus import words
import re

#nltk.download('words')
english_words = set(words.words())


def read_text_files(directory):
    documents = []
    filenames = []

    for filename in os.listdir(directory):
        if filename.endswith('.txt'):
            filepath = os.path.join(directory, filename)
            with open(filepath, 'r', encoding='utf-8') as file:
                text = file.read()

                words = [word for word in re.findall(r'\b\w+\b', text.lower()) if word in english_words]
                words = ' '.join(words)
                documents.append(words)
                filenames.append(filename)

    filenames = sorted(filenames, key=lambda x: int(x.split(' - ')[0]))
    filenames = [entry.split(' - ')[1] for entry in filenames]
    return documents, filenames


def create_term_document_matrix(documents):
    vectorizer = TfidfVectorizer(strip_accents='unicode',min_df = 0.03,stop_words='english')
    tdm = vectorizer.fit_transform(documents).toarray()
    tdm = np.transpose(tdm)
    return tdm, vectorizer.get_feature_names_out()

directory = 'Data'
documents, filenames = read_text_files(directory)
tfidf, terms = create_term_document_matrix(documents)
tfidf_df = pd.DataFrame(tfidf ,columns=filenames, index=terms)
print(tfidf_df.info)
tfidf_df.to_csv('TF-IDF.csv')
```

## Step 4

Visualize word vectors using t-SNE:

*(To get different results, modify the* `random_state` *variable in the t-SNE script.)*

```python
import pandas as pd
from sklearn.manifold import TSNE
import matplotlib.pyplot as plt
from mpl_toolkits.mplot3d import Axes3D
from adjustText import adjust_text
```

```python
random_state = 42
```

```python
def visualize_with_tsne(tfidf_matrix):
    tsne = TSNE(n_components=2, random_state= random_state)

    X = tfidf_matrix.sample(n=150, random_state= random_state)
    labels = X.index.values

    X_tsne = tsne.fit_transform(X.values)

    plt.figure(figsize=(20, 16))

    plt.scatter(X_tsne[:, 0], X_tsne[:, 1], s=5)

    texts = []
    for i, label in enumerate(labels):
        texts.append(plt.text(X_tsne[i, 0], X_tsne[i, 1], label, fontsize=8, ha='center', va='center'))

    adjust_text(texts, arrowprops=dict(arrowstyle="->", color='r', lw=0.5))

    plt.title('2D t-SNE with Word Labels')
    plt.xlabel('Dimension 1')
    plt.ylabel('Dimension 2')
    plt.show()
```

```python
file_name = 'TD.csv'
tfidf_matrix = pd.read_csv(file_name, index_col=0)
visualize_with_tsne(tfidf_matrix)
```