# Homework #1

Introduction of Text Mining
NTU, Fall 2023

# Outline

➢ Programming Assignment 1
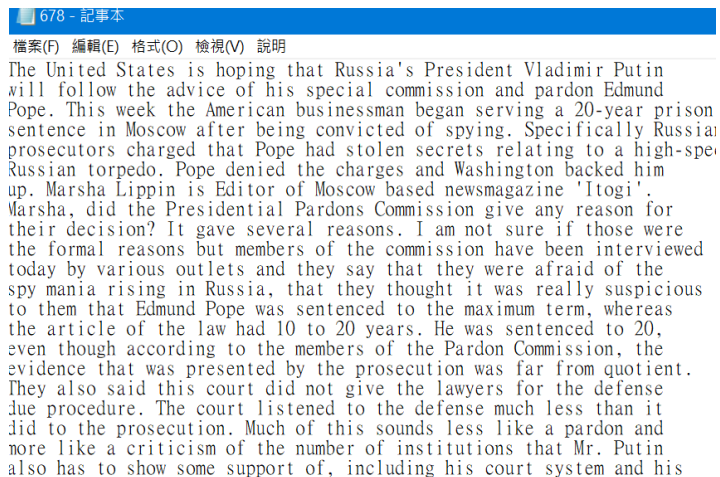
➢ Tools

➢ Submission

➢ Homework policy

# Outline

➢ **Programming Assignment 1**

➢ Tools

➢ Submission

➢ Homework policy

# Programming Assignment 1

**1.** Write a program to convert a set of documents into TF-IDF vectors

- ○ Text collection: 1095 news documents

- ○ 1.txt ~ 1095.txt

- ○ **Use `TfidfVectorizer`**

    - ■ Lowercase everything

    - ■ Filter out English stopwords

# Programming Assignment 1

**2.** Save each document's TF-IDF vector as a plain text file

- ○  1.vec, 2.vec, … 1095.vec

**3.** Load the TF-IDF vectors of documents 1 and 2, and calculate their cosine similarity.

- Please zip the vector files of documents 1 and 2, your code, and a report (indicating the similarity of doc 1 and doc 2) to TA.

- You have 2 weeks to do the work, deadline is **2023/10/8**.

# Outline

- ➢ Programming Assignment 1

- ➢ **Tools**

- ➢ Submission

- ➢ Homework policy

# Tool for Dataset

- **Download the dataset**

  - Manually download the dataset here

    https://cool.ntu.edu.tw/courses/28998/files/4425061?wrap=1

# Tool - Google Colab Tutorial

- **We suggest everyone use Google Colab to finish HW and provide some resources for someone who is not familiar with Google Colab.**
  - Introduction to Google Colab
  - Google Colab Tutorial
  - How to connect to your drive with google colab

# Outline

➢ Programming Assignment 1

➢ Tools

➢ **Submission**

➢ Homework policy

# Submission

- Deadline: **2023/10/8 (Sun.) 23:59 (GMT+8)**

- Your submission should include the following files:
  - D11725003
  - ├ report.pdf
  - ├ pa1.*(.py, .ipynb)
  - ├ output/
  -    ├ 1.vec
  -    ├ 2.vec

- Ensure your code can be executed successfully on **google colab** before your submission.

- You shall **NOT** hardcode any path in your python files, while the dataset given would be the absolute path to the directory.

# Outline

# How to find help

- Google!

- Post your question to NTU COOL

- Contact TAs by e-mail: d11725003@ntu.edu.tw