# LING575
# Summarization

2024.01

# #D2
# Process a docSet

# Overview

- Process XML -> Retrieve corresponding articles from DocSets
- Process Articles: Segment paragraph into sentences, tokenize sentences
- Summarization System Plan Overview

# extract.py

- Use XML.etree.ElementTree library to parse XML
- Extract topic ids and corresponding doc id from DocSetA
    - topic.docSetA.id
- Transform doc id into path
    - Case by case according to doc id format and year, but overall it is extracting dir_name, year and file_name, and generate the local path accordingly
    - e.g.
        - APW19990914.0234 -> /corpora/LDC/LDC02T31/apw/1999/19990914_APW_ENG
        - APW_ENG_20050609.0625 -> /corpora/LDC/LDC08T25/data/apw_eng/apw_eng_200506.xml
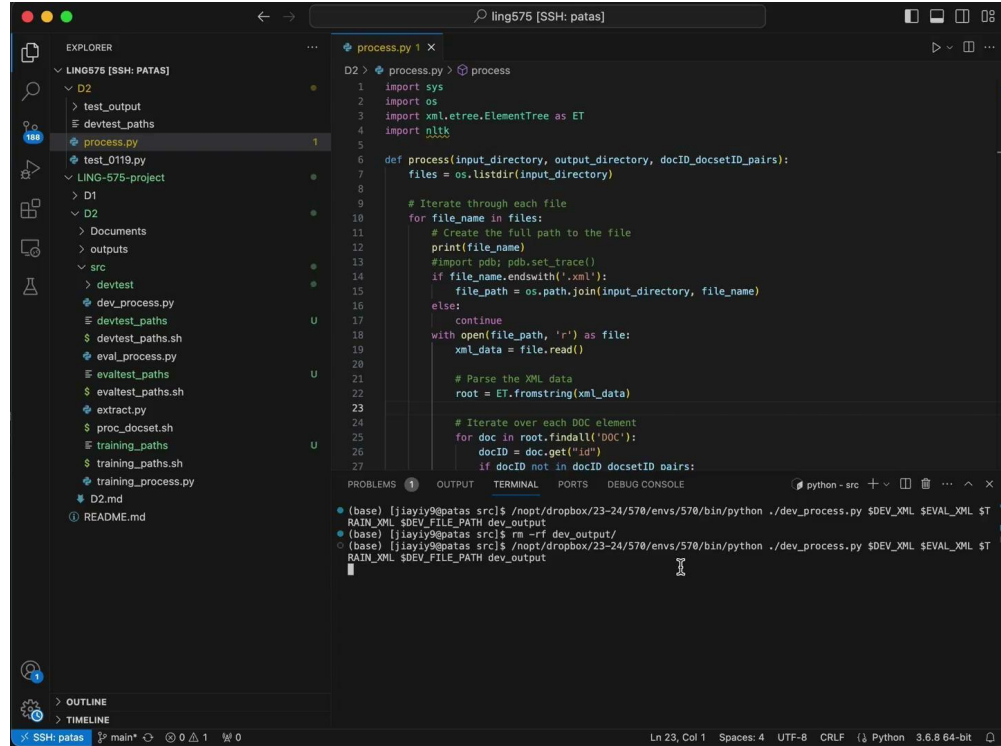
# Demo Video – 1

# process.py

- Use XML.etree.ElementTree library to parse XML
- If the file is not a standard XML (e.g. /corpora/LDC/LDC02T31/apw/1999/19990914_APW_ENG), then modify the file by adding <DOCSTREAM> tag
  - Find doc id by looking for DOCNO keyword
- If the file is standard XML, find the doc by looking for id keyword.
- Get headline and dateline from HEADLINE and DATELINE keyword.
- Get individual sentence by doc.TEXT.findall("P")
- Use nltk.word_tokenize to tokenize the sentences

# Demo Video - 2

# Methods we are considering for summarization

- TF-IDF
    - Create a Document-Term Matrix using TF-IDF
    - Score the sentences using TF-IDF score for each term
- LLR
    - Create a Document-Term Matrix using LLR
    - Score the sentences using LLR score for each term
- LSI & LDA topic modelling
    - Latent Semantic Indexing (LSI) and Latent Dirichlet Allocation (LDA)
    - Used to extract topics  from a collection of documents, and the topics can be used as feature
    - Can use key terms or phrases associated with dominant topics as content

# Methods we are considering for summarization

| Method | Pros | Cons |
|---|---|---|
| TF-IDF | - Simple and Intuitive<br>- Efficient for Extractive Summarization | - Rely heavily on word overlap<br>- Don't handle synonyms<br>- Sparse representation<br>- Sentence redundancy<br>- Lack of Context Understanding |
| LLR | - More nuanced analysis for context<br>- Handles Synonyms and Specific Contexts | - Data sensitivity<br>- Not Suitable for Abstractive Summarization involving rephrasing and paraphrasing. |
| LSI/LDA topic modeling | - Semantic Understanding<br>- Documents can be represented as a mixture of topics | - Require additional techniques to apply those extracted features<br>- Dimensionality Reduction, potentially loss of information<br>- Sensitivity to Hyperparameters |