Individual Information Retrieval Task for ICS 2205 Project **Document indexing** In [ ]: # Installing nltk library # pip install nltk 1. Parse the document to extract the data in the XML's < raw > tag In [ ]: import nltk import string from xml.etree.ElementTree import parse import os import re # Getting a list of all files in the folder folder = "docs" files = os.listdir(folder) # Looping over the files in the folder for file in files: # Parsing the .naf files using nltk's ElementTree parser file\_path = os.path.join(folder, file) tree = parse(file\_path) # Finding the <raw> tag raw\_tag = tree.find("raw") # Extracting contents from <raw> tag raw = raw\_tag.text # Remove all punctuation from the raw text using a regular expression raw\_text = re.sub(r'[{}]'.format(string.punctuation), '', raw) # Printing the extracted text without punctuation print(raw\_text) 2. Tokenise the documents' content In [ ]: # Imports from nltk.tokenize import word\_tokenize for file in files: # Parsing the .naf files using nltk's ElementTree parser file\_path = os.path.join(folder, file) tree = parse(file\_path) # Finding the <raw> tag raw\_tag = tree.find("raw") # Extracting contents from <raw> tag raw = raw\_tag.text # Remove all punctuation from the raw text using a regular expression raw\_text = re.sub(r'[{}]'.format(string.punctuation), '', raw) # Tokenizing the text into words or tokens tokens = word\_tokenize(raw\_text) # Printing the tokens print(tokens) 3. Perform case-folding, stop-word removal and stemming Case Folding: Lower Case: In [ ]: for file in files: # Parsing the .naf files using nltk's ElementTree parser file\_path = os.path.join(folder, file) tree = parse(file\_path) # Finding the <raw> tag raw\_tag = tree.find("raw") # Extracting contents from <raw> tag raw = raw\_tag.text # Remove all punctuation from the raw text using a regular expression raw\_text = re.sub(r'[{}]'.format(string.punctuation), '', raw) # Using the .lower() function to convert the string to lowercase lowercase = raw\_text.lower() # Printing lowercase text print(lowercase) • Upper Case: for file in files: # Parsing the .naf files using nltk's ElementTree parser file\_path = os.path.join(folder, file) tree = parse(file\_path) # Finding the <raw> tag raw\_tag = tree.find("raw") # Extracting contents from <raw> tag raw = raw\_tag.text # Remove all punctuation from the raw text using a regular expression raw\_text = re.sub(r'[{}]'.format(string.punctuation), '', raw) # Using the .upper() function to convert the string to uppercase uppercase = raw\_text.upper() print(uppercase) Stop-word removal: • In this example, the nltk.corpus.stopwords.words() function is used to get a list of English stop words. The stop words are removed from the list of tokens using a list comprehension. for file in files: # Parsing the .naf files using nltk's ElementTree parser file\_path = os.path.join(folder, file) tree = parse(file\_path) # Finding the <raw> tag raw\_tag = tree.find("raw") # Extracting contents from <raw> tag raw = raw\_tag.text # Remove all punctuation from the raw text using a regular expression raw\_text = re.sub(r'[{}]'.format(string.punctuation), '', raw) # Getting a list of English stop words stop\_words = nltk.corpus.stopwords.words("english") # Splitting the text into tokens tokens = nltk.tokenize.word\_tokenize(raw\_text) # Removing the stop words from the tokens filtered\_tokens = [token for token in tokens if token not in stop\_words] # Printing the resulting list of tokens print(filtered\_tokens) **Stemming:** • In this example, the Porter stemming algorithm was used The Porter stemming algorithm (or 'Porter stemmer') is a process for removing the commoner morphological and inflexional endings from words in English. Its main use is as part of a term normalisation process that is usually done when setting up Information Retrieval systems. [Martin, P. (1999). The Porter Stemming Algorithm. [Online]. Available: https://tartarus.org/martin/PorterStemmer] In [ ]: # Imports from nltk.stem import PorterStemmer for file in files: # Parsing the .naf files using nltk's ElementTree parser file\_path = os.path.join(folder, file) tree = parse(file\_path) # Finding the <raw> tag raw\_tag = tree.find("raw") # Extracting contents from <raw> tag raw = raw\_tag.text # Remove all punctuation from the raw text using a regular expression raw\_text = re.sub(r'[{}]'.format(string.punctuation), '', raw) # Creating a stemmer object stemmer = PorterStemmer() # Stemming the raw text stemmed\_text = [stemmer.stem(word) for word in raw\_text.split()] # Printing stemmed raw text print(stemmed\_text) 4. Build the term by document matrix containing the T F.IDF weight for each term within each document In [ ]: import math import numpy as np from nltk.tokenize import word\_tokenize def compute\_tfidf(documents): # Tokenize the documents tokens = [word\_tokenize(doc) for doc in documents] # Compute the term frequency (TF) for each term tf = [] for doc in tokens: term freq = {} for term in doc: if term in term\_freq: term\_freq[term] += 1 else: term\_freq[term] = 1 tf.append(term\_freq) # Compute the inverse document frequency (IDF) for each term idf = {} num\_docs = len(documents) for doc in tokens: for term in doc: if term in idf: continue df = sum(1 for d in tokens if term in d) idf[term] = math.log(num\_docs / df) # Compute the TF-IDF weight for each term tfidf = [] for doc in tf: doc tfidf = {} for term, freq in doc.items(): doc\_tfidf[term] = freq \* idf[term] tfidf.append(doc\_tfidf) # Create the term-by-document matrix terms = sorted(idf.keys()) term\_by\_doc\_matrix = np.zeros((len(terms), num\_docs)) for i, term in enumerate(terms): for j, doc in enumerate(tfidf): if term in doc: term\_by\_doc\_matrix[i, j] = doc[term] return terms, term\_by\_doc\_matrix # List of documents documents = [] for file in files: # Parsing the .naf files using nltk's ElementTree parser file\_path = os.path.join(folder, file) tree = parse(file\_path) # Finding the <raw> tag raw\_tag = tree.find("raw") # Extracting contents from <raw> tag raw\_text = raw\_tag.text # Adding the document to the list of documents documents.append(raw\_text) # Compute the term-by-document matrix terms, term\_by\_doc\_matrix = compute\_tfidf(documents) # Print the TF-IDF weights for each term in the matrix for i, row in enumerate(term\_by\_doc\_matrix): term = terms[i] print(f'TF-IDF weights for term "{term}": {row}') print('\n') **Querying Component** • ### Get a user query – note that it can be set within the notebook directly. Into a variable named query; Outputting all the quries: import xml.etree.ElementTree as ET # Get the current working directory cwd = os.getcwd() # Construct the path to the "queries" folder folder = os.path.join(cwd, "queries") # Get a list of all files in the "queries" folder files = os.listdir(folder) # Loop over the files in the "queries" folder for file in files: # Parse the .naf file using ElementTree file\_path = os.path.join(folder, file) tree = ET.parse(file\_path) # Find the <raw> tag raw\_tag = tree.find("raw") # Extract the contents of the <raw> tag query\_raw\_text = raw\_tag.text # Print the extracted text print(query\_raw\_text) Selecting one of the quries (choose from above): In [ ]: # Get the user's query query = input("Enter your query: ") # Print the user's query print("Your query:", query) • ## Preprocess the user query (tokenisation, case-folding, stop-word removal and stemming) In [ ]: # Imports import re # Converting the query to Lowercase query = query.lower() # Tokenizing the query into words or tokens tokens = word\_tokenize(query) # Getting a list of English stop words stop\_words = nltk.corpus.stopwords.words("english") # Removing the stop words from the tokens filtered\_tokens = [token for token in tokens if token not in stop\_words] # Creating a stemmer object stemmer = PorterStemmer() # Stemming the tokens query\_tokens = [stemmer.stem(token) for token in tokens] # Print the preprocessed tokens print("Preprocessed query:", query\_tokens) • ## Use cosine similarity to calculate the similarity between the query and each document Calculating cosine similarity between query and documents from collections import Counter # Get the current working directory cwd = os.getcwd() # Construct the path to the "docs" folder folder = os.path.join(cwd, "docs") # Get a list of all files in the "docs" folder files = os.listdir(folder) def cosine\_similarity(query\_tokens, doc\_tokens): # Create a Counter object for each list of tokens query\_counter = Counter(query\_tokens) doc\_counter = Counter(doc\_tokens) # Create a list of unique words present in either list of tokens unique\_words = set(query\_tokens).union(set(doc\_tokens)) # Calculate the dot product of the frequency vectors dot\_product = sum(query\_counter[word] \* doc\_counter[word] for word in unique\_words) # Calculate the Euclidean length of the frequency vectors query\_length = sum(query\_counter[word]\*\*2 for word in unique\_words) doc\_length = sum(doc\_counter[word]\*\*2 for word in unique\_words) # If either length is 0, return 0 as the cosine similarity if query\_length == 0 or doc\_length == 0: return 0 # Calculate the cosine similarity as the dot product of the frequency vectors divided by the product of their Euclidean Lengths cosine\_sim = dot\_product / (math.sqrt(query\_length) \* math.sqrt(doc\_length)) return cosine\_sim # Iterate over all the documents in the folder for file in files: # Parsing the .naf files using nltk's ElementTree parser file\_path = os.path.join(folder, file) tree = parse(file\_path) # Finding the <raw> tag raw\_tag = tree.find("raw") # Extracting contents from <raw> tag raw\_text\_docs = raw\_tag.text # Converting the raw\_text\_docs to lowercase raw\_text\_docs = raw\_text\_docs.lower() # Tokenizing the document into words or tokens tokens\_docs = word\_tokenize(raw\_text\_docs) # Removing stop words from the tokens filtered\_tokens\_docs = [token for token in tokens\_docs if token not in stop\_words] # Stemming the tokens stemmed\_tokens\_docs = [stemmer.stem(token) for token in filtered\_tokens\_docs] # Remove duplicates from the list no\_duplicates = [] for token in stemmed\_tokens\_docs: if token not in no\_duplicates: no\_duplicates.append(token) # Calculate the cosine similarity between the query and the current document cosine\_sim = cosine\_similarity(query\_tokens, no\_duplicates) # Print the cosine similarity between the query and the current document print("File: " f"{file} - {cosine\_sim}\n") • ## Output the list of documents as a ranked list results = [] for file in files: # Parsing the .naf files using nltk's ElementTree parser file\_path = os.path.join(folder, file) tree = parse(file path) # Finding the <raw> tag raw\_tag = tree.find("raw") # Extracting contents from <raw> tag raw\_text\_docs = raw\_tag.text # Converting the raw\_text\_docs to lowercase raw\_text\_docs = raw\_text\_docs.lower() # Tokenizing the document into words or tokens tokens\_docs = word\_tokenize(raw\_text\_docs) # Removing stop words from the tokens filtered\_tokens\_docs = [token for token in tokens\_docs if token not in stop\_words] # Stemming the tokens stemmed\_tokens\_docs = [stemmer.stem(token) for token in filtered\_tokens\_docs] # Remove duplicates from the list no\_duplicates = [] for token in stemmed\_tokens\_docs: if token not in no\_duplicates: no\_duplicates.append(token) # Calculate the cosine similarity between the query and the current document cosine\_sim = cosine\_similarity(query\_tokens, no\_duplicates) # Add the file name and cosine similarity to the results list results.append((file, cosine\_sim)) # Sort the results list by cosine similarity ranked\_results = sorted(results, key=lambda x: x[1], reverse=True) # Iterate over the ranked results and print the file name and cosine similarity for result in ranked\_results: file, cosine\_sim = result print(f"File: {file} - Similarity: {cosine\_sim}\n")