Exercise 1: Data cleaning and text tokenization

- and output as a comma separated document.
 - tokenize the file are return back to master. Master concatenate these list from workers and write in a file in which each line give
- tokens of one document.
- First I read the filepaths into a list and then splitted this list among workers. Each worker open files in its list and then clean and

- 1. Cleaning: remove all punctuations, numbers and The list of common English stopwords used in this exercise sheet can be found in the reference [4]. 2. Tokenize: Tokenize your documents so it is easy to process in the next task i.e. Tokenize words

- that are cleaned and tokenized.
- oped program should take a set of documents in raw format as input, and outputs a set of documents Please explain your solution and how you distribute work among workers.

- the above mentioned words/numbers. You can take help of some python libraries i.e. NLTK. The devel-
- In a text document you encounter words that are not helpful for your final model. For example, punctuations and numbers, meaningless words, common English stopwords etc. Your first task is to preprocess your data so you can remove these words from the documents. Your solution should be based on MPI framework. You have to develop (write code) a distributed algorithm that cleans the data by removing

- comm = MPI.COMM WORLD # setting up the MPI communicator total worker = comm.Get size() # getting number of workers rank = comm.Get rank() # storing rank of each worker start time = MPI.Wtime() # time variable to find time at the end

Each row of return list represent on document'''

if rank == master: # Condition of checking master to distribute data

filenames in os.walk(dir):

file paths.append(abs path train)

cleanData = [(lambda x: [word for word in

in docListl

if len(filenames) > 0:

fileList = getDocumentList('testDir')

chunk size = len(fileList) // total worker

Splitting file path list between all workers

comm.send(filesInChunk, dest=i, tag=1)

tokenData = comm.recv(source=i, tag=3)

Each Document is separated by new line''' with open('TokenizeData.txt', 'w') as f: for tokendoc in tokenDocList:

filesInChunk = comm.recv(source=master, tag=1)

', 'computer', 'means', 'never', 'say', 'sorry']]

Exercise 2: Calculate Term Frequency (TF)

Execution time : 287.82005559999925 PS C:\Users\wasif\DDA Lab\Exercise02>

a document d can be calculated as.

verify if you get the same result at the end.

comm = MPI.COMM WORLD # setting up the MPI communicator total worker = comm.Get size() # getting number of workers

start time = MPI.Wtime() # time variable to find time at the end

Normalizing by dividing total number of terms in the document

:param tokenizeDocChunk: Chunk of Tokenized Data list

:return: Term frequency of each word in the document

Calculating Term Frequency for each word in the document list by using counter

words = dict(Counter(lines)) # counting term frequency for each word

words.update({key: words[key] / numOftokens}) # updating normalize term frequency

removing new line and commas and splitting each line with comma(,) into a list

docIds = range(1, len(tokenizeDoc) + 1) # for document ids stored as key in dictionary chunk size = len(tokenizeDoc) // total worker # each chunk size processed by worker

comm.send(tokenizeDoc[(i - 1) * chunk size:i * chunk size], dest=i, tag=1)

docFreqDict = calculate Tf(tokenizeDoc[i * chunk size:], docIds[i * chunk size:])

receiving chunks of term frequencies from other workers and concatenating them

print("Execution Time :", MPI.Wtime() - start time) # total time takes to task complition

Term Frequencies : document1 {'path': 0.037037037037037035, 'newsgroups': 0.037037037037037035, 'subject': 0.037037037 037037035, 'argic': 0.037037037037037035, 'date': 0.037037037037037035, 'apr': 0.07407407407407407, 'organization': 0. 037037037037035, 'dimension': 0.037037037037037035, 'bbs': 0.037037037037035, 'lines': 0.037037037037035, 'se nder': 0.0370370370370377035, 'daemon': 0.037037037037037035, 'aswer': 0.037037037037037035, 'one': 0.03703703703703703 5, 'question': 0.037037037037037035, 'get': 0.037037037037035, 'retarded': 0.037037037037037035, 'system': 0.037037 037037037035, 'phone': 0.037037037037037035, 'cute': 0.037037037037035, 'quote': 0.037037037037037035, 'computer': 0.037037037037035, 'means': 0.037037037037037035, 'never': 0.037037037037035, 'say': 0.037037037037035, 'sorr

Term Frequencies : document1 {'path': 0.037037037037037035, 'newsgroups': 0.037037037037037035, 'subject': 0.037037037 037037035, 'argic': 0.037037037037037035, 'date': 0.037037037037035, 'apr': 0.07407407407407407, 'organization': 0. 037037037037037035, 'dimension': 0.037037037037037035, 'bbs': 0.037037037037035, 'lines': 0.037037037037035, 'se nder': 0.037037037037037035, 'daemon': 0.037037037037037035, 'aswer': 0.037037037037037035, 'one': 0.03703703703703703 5, 'question': 0.037037037037037035, 'get': 0.037037037037035, 'retarded': 0.037037037037037035, 'system': 0.037037 037037037035, 'phone': 0.037037037037037035, 'cute': 0.037037037037035, 'quote': 0.037037037037037035, 'computer': 0.037037037037035, 'means': 0.037037037037037035, 'never': 0.037037037037035, 'say': 0.037037037037035, 'sorr

Term Frequencies : document1 {'path': 0.037037037037037035, 'newsgroups': 0.037037037037037035, 'subject': 0.037037037 037037035, 'argic': 0.037037037037037035, 'date': 0.037037037037035, 'apr': 0.07407407407407407, 'organization': 0. 037037037037037035, 'dimension': 0.037037037037037035, 'bbs': 0.037037037037037035, 'lines': 0.037037037037037 nder': 0.037037037037037705, 'daemon': 0.037037037037037035, 'aswer': 0.037037037037037035, 'one': 0.03703703703703703 5, 'question': 0.037037037037037035, 'get': 0.037037037037035, 'retarded': 0.037037037037037035, 'system': 0.037037 037037037035, 'phone': 0.037037037037037035, 'cute': 0.037037037037035, 'quote': 0.037037037037037035, 'computer': 0.037037037037035, 'means': 0.037037037037037035, 'never': 0.037037037037035, 'say': 0.037037037037037035, 'sorr

Execution Time

1.7063034999955562

1.3614088999966043

1.4460126999983913

(2)

comm.send(docIds[(i - 1) * chunk size:i * chunk size], dest=i, tag=2)

tokenizeDoc = [(line.rstrip('\n')).replace("'", "").replace(' ', '').split(",") for line in f]

numOftokens = sum(words.values()) # total terms in a document

tfDictChunk.update({'document' + str(docIds[i]): words})

if rank == master: # Condition of checking master to distribute data

reading tokenized data text file from last exercise

for i in range(1, total worker): # for loop to send data # sending the chunk of tokenized data list to workers

sending the document id for keys in dictionary

docFreqDictChunk = comm.recv(source=i, tag=3)

writing the TF dictionary to a file for later exercises

receiving the chunk of tokenized data list from master

PS C:\Users\wasif\DDA Lab\Exercise02> mpiexec.exe -n 2 python .\exer2.py

PS C:\Users\wasif\DDA Lab\Exercise02> mpiexec.exe -n 4 python .\exer2.py

PS C:\Users\wasif\DDA Lab\Exercise02> mpiexec.exe -n 8 python .\exer2.py

Results are same when experimenting with different number of workers.

Workers P

Exercise 3: Calculate Inverse Document Frequency ((IDF)

the corpus that contain a particular token. The IDF formula is given as

The Inverse Document Frequency ((IDF) is counting the number of documents in the corpus and counting the number of documents that contain a token. While the TF is calculated on a per-document basis, the IDF is computed on the basis of the entire corpus. The final IDF score of a token t in the corpus C is obtained by taking logarithm of document count in the corpus divided by the number of documents in

 $IDF(t) = \log \frac{|C|}{\sum_{d \in C} \mathbb{1}(t, d)} ,$

where $\mathbb{1}(t,d)$ is an indicator function which returns 1 if a token t exists in document d. |C| is the

Develop an solution using MPI framework and write cod. Please explain how you parallelize (or distribute) IDF(t,d) calculation. Also explain your strategy from the data division and calculation division point of view as well. Perform a small experiment by varying number workers i.e. $P = \{2, 4, 8\}$.

Reading file containg tf from last exercise as dictionary and then slicing and sending this dictionary to other workers. Each worker

Calculates sum of all term frequency for all term in a document and then master joins these dictionary and divide C with each

if (idfDictChunk.get(wordKey) != None): # if token already exist, add 1 in token

idfDictChunk.update({wordKey: idfDictChunk.get(wordKey) + 1})

chunk size = len(tfDict) // total worker # each chunk size processed by worker

tfDictChunk = {'document' + str(k): tfDict['document' + str(k)] for k in

if 'document' + str(k) in tfDict}

tfDictChunk = {'document' + str(k): tfDict['document' + str(k)] for k in range((i * chunk size) + 1, len(tfDict) + 1) if

'''receiving chunks from workers and concatinating them to make one dictionary

diving C with each term in this dictionary and taking log to give final Idf Dictionary

print("Execution Time :", MPI.Wtime() - start time) # total time takes to task completion

Idf for terms starts with abilit: {'ability': 3.9925408972444756, 'abilities': 6.074696144795908}

Idf for terms starts with abilit: {'ability': 3.9925408972444756, 'abilities': 6.074696144795908}

Idf for terms starts with abilit: {'ability': 3.9925408972444756, 'abilities': 6.074696144795908}

In this exercise you have to think about how you can combine the complete pipeline in a single parallel/distributed program that can run worker $P = \{2,4,8\}$. Develop an solution using MPI framework and write code. Please explain how you parallelize (or distribute) the complete pipeline. Also explain your strategy from the data division and calculation division point of view as well. Perform a small experiment by varying number workers i.e. $P = \{2, 4, 8\}$. Also verify if you get the same result at the

Reading TF and IDF from files and storing in dictionaries. Splitting the TF dictionary among all workers while broadcasting IDf

Workers P

2

print("Idf for terms starts with abilit: ", {k:v for k,v in finalIdfDict.items() if k.startswith('abilit')

Execution Time

1.9553215999985696

1.8562213000041083

1.91034939999372

(3)

range(((i - 1) * chunk size) + 1, (i * chunk size) + 1)

tokenizeDocChunk = comm.recv(source=master, tag=1)

docIds = comm.recv(source=master, tag=2)

Showing term frequencies for first and last document.

y': 0.037037037037037035}

y': 0.037037037037037035}

y': 0.037037037037037035}

Execution Time : 1.7063034999955562

Execution Time : 1.3614088999966043

Execution Time : 1.4460126999983913

total number of documents in the corpus.

from collections import Counter

master = 0 # master with id=0

def calculateSumTF(tfDictChunk):

idfDictChunk = {}

for key in tfDictChunk:

return idfDictChunk

docDic = tfDictChunk[key] for wordKey in docDic:

tfDict = json.load(open("TF.txt"))

sending to other workers

master calculation of own chunk

for i in range(1, total worker):

finalIdfDict = {}

else:

for key in dict(idfDict):

Execution Time : 1.9553215999985696

Execution Time : 1.8562213000041083

Execution Time : 1.91034939999372

Same results when changing number of workers.

idfDict += Counter(idfDictChunk)

C = len(tfDict) # total number of documents

comm.send(tfDictChunk, dest=i, tag=1)

idfDict = Counter(calculateSumTF(tfDictChunk))

idfDictChunk = comm.recv(source=i, tag=3)

json.dump(finalIdfDict, open("IDF.txt", 'w'))

tfDictChunk = comm.recv(source=master, tag=1)

Showing Idf for words starting with ast.Remaining Idf are stored in IDF.txt.

printing term starting with abilit

from mpi4py import MPI

In []: import json

import math

Also verify if you get the same result at the end.

term in this dictionary and taking log to give final Idf Dictionary

comm = MPI.COMM WORLD # setting up the MPI communicator total worker = comm.Get size() # getting number of workers

start time = MPI.Wtime() # time variable to find time at the end

else: # if token does not exist put 1 idfDictChunk.update({wordKey: 1})

reading term frequency dictionary from last exercise

for i in range(1, total worker): # for loop to send data

'document' + str(k) in tfDict}

in which each term tell in how many documents it appears'''

finalIdfDict.update({key: math.log(C / idfDict[key])})

receiving the chunk of term frequency dictionary list from master

writing the IDF dictionary to a file for later exercises

comm.send(calculateSumTF(tfDictChunk), dest=master, tag=3)

PS C:\Users\wasif\DDA Lab\Exercise02> mpiexec.exe -n 2 python .\Exer3.py

PS C:\Users\wasif\DDA Lab\Exercise02> mpiexec.exe -n 4 python .\Exer3.py

PS C:\Users\wasif\DDA Lab\Exercise02> mpiexec.exe -n 8 python .\Exer3.py

You have already calculated TF(t,d) in Exercise 2 and IDF(t) is Exercise 3.

dictionary. Calculating tf-idf at each process and then merging the results at master process.

comm = MPI.COMM WORLD # setting up the MPI communicator total worker = comm.Get size() # getting number of workers

start_time = MPI.Wtime() # time variable to find time at the end

if rank == master: # Condition of checking master to distribute data#

tfDictChunk = {'document' + str(k): tfDict['document' + str(k)] for k in range((i * chunk_size) + 1, len(tfDict) + 1) if

receiving the chunk of term frequency dictionary list from master

0.27475669968507416, 'system': 0.06466544815178625, 'phone': 0.10796016145501523, 'cute': 0.22190053836506876, 'quote' : 0.14847906322296858, 'computer': 0.0756437472743931, 'means': 0.10307011736638945, 'never': 0.08134351429313073, 'sa

Tf-Idf: document19997: {'path': 0.0, 'newsgroups': 0.0, 'subject': 0.0, 'argic': 0.17836078540935624, 'date': 0.0, 'a pr': 0.005880024274495046, 'organization': 0.0015179456336085145, 'dimension': 0.23110281093167986, 'bbs': 0.138597415

54036764, 'lines': 0.00010757961030013146, 'sender': 0.02251219099773859, 'daemon': 0.17905308663202857, 'aswer': 0.36 679027930685193, 'one': 0.03591382824385434, 'question': 0.08065059628507597, 'get': 0.04888129516026783, 'retarded': 0.27475669968507416, 'system': 0.06466544815178625, 'phone': 0.10796016145501523, 'cute': 0.22190053836506876, 'quote' : 0.14847906322296858, 'computer': 0.0756437472743931, 'means': 0.10307011736638945, 'never': 0.08134351429313073, 'sa

Execution Time

2.139510100001644

2.0363853000017116

2.084294799999043

'document' + str(k) in tfDict}

reading term frequency dictionary from last exercise

reading term frequency dictionary from last exercise

rank = comm.Get rank() # storing rank of each worker

slicing dictionary for sending to workers

if rank == master: # Condition of checking master to distribute data

Calculating sum of all term frequency for all term in a document

:return: dictionary of every term frequency across all documents

:param tfDictChunk: List of Term frequency dictionary for each document

rank = comm.Get_rank() # storing rank of each worker

receiving the document ids for keys in dictionary

ordering the term frequency dictionary with key as document ids

print("Term Frequencies : document19997", orderDocFreq['document1321'])

comm.send(calculate Tf(tokenizeDocChunk, docIds), dest=master, tag=3)

docFreqDict.update(docFreqDictChunk)

json.dump(orderDocFreq, open("TF.txt", 'w'))

orderDocFreq = dict(sorted(docFreqDict.items()))

with open('TokenizeData.txt', 'r') as f:

master calculation of own chunk

for i in range(1, total worker):

else:

rank = comm.Get rank() # storing rank of each worker

def calculate Tf(tokenizeDocChunk, docIds):

i = 0 # for iterating document ids

for lines in tokenizeDocChunk:

for key in words:

:param docIds: Chunk of document ids list

tokens in the document d.

from collections import Counter

master = 0 # master with id=0

tfDictChunk = {}

i += 1

return tfDictChunk

from mpi4py import MPI

In []: import json

comm.send(cleanTokenizeData(filesInChunk), dest=master, tag=3)

print('Tokenized and clean Data List:', tokenDocList) print('Execution time :', MPI.Wtime() - start time)

master calculation cleaning the chunk

for i in range(1, total_worker):

tokenDocList += tokenData

else: # workers calculations

less number of workers.

for i in range(1, total worker): # for loop to send data

tokenDocList = cleanTokenizeData(fileList[i * chunk size:])

'''Writing the tokenize data list to a file to be used in next exercise.

This list of tokenized and clean data is written in 'TokenizeData.txt' which I have placed it in code folder.

The Term Frequency (TF) is calculated by counting the number of times a token occurs in the document. This TF score is relative to a specific document, therefore you need to normalized it by dividing with the total number of tokens appearing in the document. A normalized TF score for a specific token t in

 $TF(t,d) = \frac{n^d(t)}{\sum_{t' \in d} n^d(t')} ,$

Develop an solution using MPI framework and write code. Please explain how you parallelize (or distribute) TF(t,d) calculation. Also explain your strategy from the data division and calculation division point of view as well. Perform a small experiment by varying number workers i.e. $P = \{2, 4, 8\}$. Also

where $n^{d}(t)$ is the number of times a token t appears in a document d and |d| is the total number of

Reading tokenize file from last exercise and storing in a list. Splitted this list into workers. Eeach worker calculates Term Frequency for each word in the document list by using counter and Normalize by dividing total number of terms in the document.At master gathered these chunks into dictionary where every key is a document Id and value is Term frequency of tokens in that documents

(1)

f.write(str(tokendoc)[1: -1] + '\n') # removing brackets before writing

I ran the code on 8 workers. More time is taken because of IO operation of reading files. Thats why I could not experiment with

fessor', 'turkish', 'studies', 'university', 'michigan', 'list', 'goes', 'serdar', 'argic', 'closed', 'university', 'minnesota', 'tom', 'goodrich', 'professor', 'history', 'indiana', 'university', 'penns ylvania', 'tibor', 'professor', 'emeritus', 'turkish', 'studies', 'columbia', 'university', 'justin', 'mccarthy', 'professor', 'history', 'university', 'louisville', 'jon', 'mandaville', 'professor', 'h istory', 'portland', 'state', 'university', 'oregon', 'robert', 'olson', 'professor', 'history', 'uni versity', 'kentucky', 'madeline', 'zilfi', 'professor', 'history', 'university', 'maryland', 'james', 'professor', 'turkish', 'studies', 'university', 'michigan', 'list', 'goes', 'serdar', 'argic', 'clo sed', 'roads', 'mountain', 'passes', 'might', 'serve', 'ways', 'escape', 'turks', 'proceeded', 'work' , 'extermination', 'ohanus', 'appressian', 'soviet', 'armenia', 'today', 'longer', 'exists', 'single' 'turkish', 'soul', 'sahak', 'melkonian'], ['path', 'newsgroups', 'subject', 'ajerk', 'date', 'apr', 'organization', 'dimension', 'bbs', 'lines', 'sender', 'daemon', 'apr', 'good', 'case', 'rights', 'a bortion', 'system', 'phone', 'cute', 'quote', 'computer', 'means', 'never', 'say', 'sorry'], ['path', 'newsgroups', 'subject', 'argic', 'date', 'apr', 'organization', 'dimension', 'bbs', 'lines', 'sende r', 'daemon', 'apr', 'aswer', 'one', 'question', 'get', 'retarded', 'system', 'phone', 'cute', 'quote

for file in filenames:

from nltk.corpus import stopwords

master = 0 # master with id=0

def cleanTokenizeData(filesInChunk):

for file in filesInChunk: f = open(file, "r") docList.append(f.read())

def getDocumentList(dir):

file paths = []

return file paths

for dir path,

In []: import os

import nltk import re

import string

from mpi4py import MPI

docList = []

return cleanData

'''Opening files in the Chunk and removing number ,english stopwords,punctuations and

'''Reading files absoulute path from dataset directory and returning it in a list '''

nltk.word tokenize(re.sub(r'[0-9]', '', x).lower().strip(string.punctuation)) if word.isalnum() and word not in stopwords.words('english') and len(word) > 2]) (doc)

abs path train = os.path.abspath(os.path.join(dir path, file)) # for complete refrence of

words smaller than 2 characters from text and then tokenizing the documents

Lab Course: Distributed Data Analytics **Exercise Sheet 2**

Syed Wasif Murtaza Jafri- 311226

- 8 Exercise 4: Calculate Term Frequency Inverse Document Frequency (TF-IDF) scores In this exercise you will find the TF-IDF(t,d) for a given token t and a document d in the corpus C is the product of TF(t,d) and IDF(t), which is represented as, TF- $IDF(t,d) = TF(t,d) \times IDF(t)$,
- def calculateTfIdf(tfDictChunk, idfDic): ''' Calculating tf-idf by multiplying tf of each term in document with idf of that term''' tf idfDictChunk = {} for key in tfDictChunk: docDic = tfDictChunk[key] tdfIdfDoc = {} for wordKey in docDic: if (idfDic.get(wordKey) != None): tdfIdfDoc.update({wordKey: docDic.get(wordKey) * idfDic.get(wordKey)}) tf idfDictChunk.update({key: tdfIdfDoc}) return tf idfDictChunk

master = 0 # master with id=0

chunk size = len(tfDict) // total worker # each chunk size processed by worker for i in range(1, total_worker): # for loop to send data # slicing tf dictionary for sending to workers tfDictChunk = {'document' + str(k): tfDict['document' + str(k)] for k in $range(((i - 1) * chunk_size) + 1, (i * chunk_size) + 1)$ if 'document' + str(k) in tfDict} comm.send(tfDictChunk, dest=i, tag=1) # sending Idf dictionary to all workers comm.send(idfDic, dest=i, tag=2)

tfDict = json.load(open("TF.txt"))

idfDic = json.load(open("IDF.txt"))

master calculation of own chunk

C = len(tfDict) # total number of documents

tf IdfDict = calculateTfIdf(tfDictChunk, idfDic) ""receiving chunks from workers and concatinating them to make one dictionary of documents where each value is also a dictionary of Terms idf in that document''' for i in range(1, total_worker): tf_idfDictChunk = comm.recv(source=i, tag=3) tf IdfDict.update(tf idfDictChunk) print("Tf-Idf : document19997:", tf_IdfDict['document19997']) print("length of Tf-Idf Dictionary : ", len(tf_IdfDict)) print("Execution Time :", MPI.Wtime() - start_time) # total time takes to task completion

else: # receiving array parts to sum

receiving whole IDF dictionary list from master idfDic = comm.recv(source=master, tag=2) comm.send(calculateTfIdf(tfDictChunk, idfDic), dest=master, tag=3) PS C:\Users\wasif\DDA Lab\Exercise02> mpiexec.exe -n 2 python .\exer4.py Tf-Idf: document19997: {'path': 0.0, 'newsgroups': 0.0, 'subject': 0.0, 'argic': 0.17836078540935624, 'date': 0.0, 'a pr': 0.005880024274495046, 'organization': 0.0015179456336085145, 'dimension': 0.23110281093167986, 'bbs': 0.138597415 54036764, 'lines': 0.00010757961030013146, 'sender': 0.02251219099773859, 'daemon': 0.17905308663202857, 'aswer': 0.36 679027930685193, 'one': 0.03591382824385434, 'question': 0.08065059628507597, 'get': 0.04888129516026783, 'retarded':

: 0.06747638893140558, 'sorry': 0.11653352414953516}

PS C:\Users\wasif\DDA Lab\Exercise02> mpiexec.exe -n 4 python .\exer4.py

Workers P

2

4

8

tfDictChunk = comm.recv(source=master, tag=1)

writing the tf-IDF dictionary to a file json.dump(tf_IdfDict, open("tf_Idf.txt", 'w'))

54036764, 'lines': 0.00010757961030013146, 'sender': 0.02251219099773859, 'daemon': 0.17905308663202857, 'aswer': 0.36 679027930685193, 'one': 0.03591382824385434, 'question': 0.08065059628507597, 'get': 0.04888129516026783, 'retarded': 0.27475669968507416, 'system': 0.06466544815178625, 'phone': 0.10796016145501523, 'cute': 0.22190053836506876, 'quote' : 0.14847906322296858, 'computer': 0.0756437472743931, 'means': 0.10307011736638945, 'never': 0.08134351429313073, 'sa y': 0.06747638893140558, 'sorry': 0.11653352414953516} length of Tf-Idf Dictionary : 19997 Execution Time : 2.0363853000017116 PS C:\Users\wasif\DDA Lab\Exercise02> mpiexec.exe -n **8** python .\exer4.py Tf-Idf : document19997: {'path': 0.0, 'newsgroups': 0.0, 'subject': 0.0, 'argic': 0.17836078540935624, 'date': 0.0, 'a pr': 0.005880024274495046, 'organization': 0.0015179456336085145, 'dimension': 0.23110281093167986, 'bbs': 0.138597415

length of Tf-Idf Dictionary : 19997

y': 0.06747638893140558, 'sorry': 0.11653352414953516}

length of Tf-Idf Dictionary : 19997 Execution Time : 2.139510100001644

- Execution Time : 2.084294799999043 Getting same results for document19997 for changing number of workers.

end.

import json

from mpi4py import MPI