

# Group Assignment 4

## Page Rank

Team 2: Thanh Le, Antoine Si

PAGE RANK

# BUILD ADJACENCY MATRIX

```
# function to build the adjacency matrix
def build_adjacency_matrix(self, input_file):
    with open(input_file, 'r') as file:
        lines = file.readlines()

        # record the number of pages
        self.page_count = (int)(lines[0].strip())

        # record the number of links
        self.link_count = (int)(lines[1].strip())

        # initialize a matrix of 0s with size of page_count x page_count
        adjacency_matrix = np.zeros((self.page_count, self.page_count))

        # iterate over each line in the input file to record the link to the matrix
        for line in lines[2:]:
            link = line.split()
            src = int(link[0])
            dst = int(link[1])
            adjacency_matrix[src][dst] = 1

    return adjacency_matrix
```

# BUILD TRANSITION PROBABILITY MATRIX

```
# function to build the transition probability matrix
def build_transition_probability_matrix(self, adjacency_matrix):
    transition_probability_matrix = np.nan_to_num(np.divide(adjacency_matrix, adjacency_matrix.sum(axis=1)[: , None]), nan=0)
    return transition_probability_matrix
```

- Given an adjacency matrix A
- If a row of A has no 1's, then replace each element by  $1/N$ .
- Otherwise, divide each 1 in A by the number of 1's in its row
  - Ex: if there is a row with three 1's, then each of them is replaced by  $1/3$ .

# BUILD TRANSITION PROBABILITY MATRIX

```
# function to build the transition probability matrix with teleporting
def build_transition_probability_matrix_with_teleporting(self, transition_probability_matrix, teleportation_rate):
    transition_probability_matrix_with_teleporting = transition_probability_matrix * (1 - teleportation_rate)
    transition_probability_matrix_with_teleporting += (teleportation_rate / self.page_count)
    return transition_probability_matrix_with_teleporting
```

- Given a transition probability matrix  $B$  and teleportation rate  $\alpha$
- Multiply  $B$  by  $1 - \alpha$
- Add  $\alpha/N$  to every entry of the resulting matrix.

# PAGE RANK

```
# function to implement pagerank algorithm
# input_file - input file that follows the format provided in the assignment description
def pagerank(self, input_file):
    # record the start time
    start_time = time.time()

    # build the adjacency matrix
    adjacency_matrix = self.build_adjacency_matrix(input_file)

    # build the transition probability matrix
    transition_probability_matrix = self.build_transition_probability_matrix(adjacency_matrix)

    # build the transition probability matrix with teleporting
    teleportation_rate = 0.15
    transition_probability_matrix_with_teleporting = self.build_transition_probability_matrix_with_teleporting(transition_probability_matrix, teleportation_rate)
```

# PAGE RANK

```
# initialize initial probability distribution vector
probability_vector = np.ones(self.page_count)/self.page_count

# initialize epsilon value as threshold for indicating steady state
eps = 1e-5

# initialize count of iterations
iterations = 0

# calculate the page rank
while True:
    last_probability_vector = probability_vector
    probability_vector = np.dot(last_probability_vector, transition_probability_matrix_with_teleporting)
    iterations += 1
    if np.sum(np.abs(probability_vector - last_probability_vector))/self.page_count < eps:
        break

# initialize a dictionary to store the result where key is the page_id and value is the corresponding pagerank value
result = {k: v for k, v in enumerate(probability_vector)}

# sort the page rank
result = dict(sorted(result.items(), key=lambda item: item[1], reverse=True))

# record the end time
end_time = time.time()

# print the output
print("Number of pages: ", self.page_count)
print("Number of links: ", self.link_count)
print("\nPageRank calculated in {} iterations ({} seconds)".format(iterations, end_time - start_time))
print("Page ID\t\tPageRank Value")
for page_id in result:
    print("{}\t\t{}".format(page_id, result[page_id]))
```

The initial vector in power iteration loop is initialized with  $1/\text{num\_of\_pages}$  in all the elements

The vector is updated over each iteration and is calculated by the dot product of the last vector and the transition probability matrix with teleporting

# PAGE RANK

```
# initialize initial probability distribution vector
probability_vector = np.ones(self.page_count)/self.page_count

# initialize epsilon value as threshold for indicating steady state
eps = 1e-5

# initialize count of iterations
iterations = 0

# calculate the page rank
while True:
    last_probability_vector = probability_vector
    probability_vector = np.dot(last_probability_vector, transition_probability_matrix_with_teleporting)
    iterations += 1
    if np.sum(np.abs(probability_vector - last_probability_vector))/self.page_count < eps:
        break

# initialize a dictionary to store the result where key is the page_id and value is the corresponding page rank value
result = {k: v for k, v in enumerate(probability_vector)}

# sort the page rank
result = dict(sorted(result.items(), key=lambda item: item[1], reverse=True))

# record the end time
end_time = time.time()

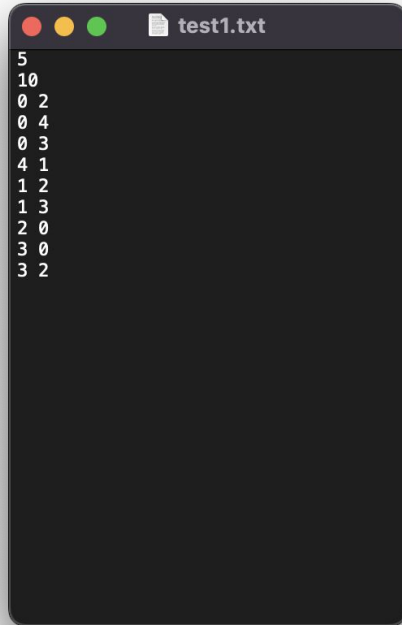
# print the output
print("Number of pages: ", self.page_count)
print("Number of links: ", self.link_count)
print("\nPageRank calculated in {} iterations ({} seconds)".format(iterations, end_time - start_time))
print("Page ID\t\tPageRank Value")
for page_id in result:
    print("{}\t\t{}".format(page_id, result[page_id]))
```

Epsilon value is set to be  $1e-5$  in this assignment. Smaller value yields more accurate result but takes more time (more iteration).

Determine the steady state by checking if the difference between the current iteration and the last one falls below the specified epsilon.



# EXPERIMENTAL RESULT (TEST1.TXT)



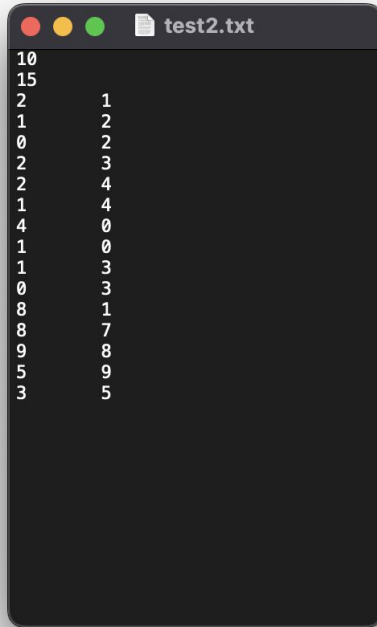
```
5
10
0 2
0 4
0 3
4 1
1 2
1 3
2 0
3 0
3 2
```

Number of pages: 5  
Number of links: 10

PageRank calculated in 19 iterations (0.00019788742065429688 seconds)

Page ID	PageRank Value
0	0.31899680476330367
2	0.25168131582784947
3	0.17661591062273685
1	0.13232630994072855
4	0.12037965884538143

# EXPERIMENTAL RESULT (TEST2.TXT)



Page ID	Number of links
10	
15	
2	1
1	2
0	2
2	3
2	4
1	4
4	0
1	0
1	3
0	3
8	1
8	7
9	8
5	9
3	5

```
Number of pages: 10
Number of links: 15

PageRank calculated in 78 iterations (0.0006330013275146484 seconds)
Page ID      PageRank Value
8            0.000146371827523615
9            0.00013953238693068264
5            0.0001321575836596206
3            0.00012420551176896475
1            0.00011417120956767715
0            0.00011092571318972676
2            9.487597832823341e-05
7            8.484167612694584e-05
4            7.276889727873644e-05
6            1.696858297238537e-05
```

# EXPERIMENTAL RESULT

Number of pages: 5  
Number of links: 10

PageRank calculated in 5 iterations (0.00014090538024902344 seconds)

Page ID	PageRank Value
0	0.31854055729166664
2	0.2509200321180555
3	0.17461386458333333
1	0.13776766319444445
4	0.11815788281250002

Number of pages: 10  
Number of links: 15

PageRank calculated in 2 iterations (9.918212890625e-05 seconds)

Page ID	PageRank Value
3	0.10453333333333337
5	0.10347083333333336
8	0.09745000000000001
9	0.09745000000000001
0	0.08558541666666668
2	0.08222083333333335
1	0.07726250000000001
7	0.054950000000000006
4	0.05300208333333335
6	0.01245

Epsilon = 1e-2

Number of pages: 5  
Number of links: 10

PageRank calculated in 38 iterations (0.0002722740173339844 seconds)

Page ID	PageRank Value
0	0.31899306233476216
2	0.25168227014986366
3	0.17661913720967856
1	0.1323241623655669
4	0.12038136794012859

Number of pages: 10  
Number of links: 15

PageRank calculated in 184 iterations (0.0006999969482421875 seconds)

Page ID	PageRank Value
8	1.4230775553975328e-08
9	1.3565821473402226e-08
5	1.2848817580776256e-08
3	1.2075689635463531e-08
1	1.1100120054327817e-08
0	1.0784581666256505e-08
2	9.224170907035266e-09
7	8.248601325899551e-09
4	7.0748439915250015e-09
6	1.649744352000145e-09

Epsilon = 1e-10

# CONCLUSION

- Being able to detect the steady state is the metric for determine the PageRank's accuracy.
- The value of epsilon should be picked wisely to be able to detect the steady state while not keeping the program running in much time → PageRank's efficiency.

EXTRA CREDIT

[illegible]

- id (column A = index 1)
- text (column K = index 11)

for building the  
index

# INDEX FORMAT

Each term entry in index will be of the form:

term: [(doc\_id1, term\_frequency1), (doc\_id2, term\_frequency2), ...]

The positional data and weights are discarded in this assignment since the goal is focus on performance comparison.

# TRADITIONAL INDEX

## Using the one from Group Assignment 3

```
# function to build an index using PyLucene
def build_lucene_index(self):
    directory_path = os.path.dirname(self.path) + '/index'

    # construct the directory to store the index on local file system
    self.directory = FSDirectory.open(File(directory_path).toPath())

    # configure an index write
    config = IndexWriterConfig()

    # always overwrite existing index to avoid duplicate files
    config.setOpenMode(IndexWriterConfig.OpenMode.CREATE)

    # construct writer, reader, and searcher
    self.writer = IndexWriter(self.directory, config)
    self.reader = DirectoryReader.open(self.directory)
    self.searcher = IndexSearcher(self.reader)

    # iterate through each line in the dataset file
    workbook = openpyxl.load_workbook(self.path).active
    for i in range(3, workbook.max_row + 1, 2):

        # get the tweet id in column 1
        id = workbook.cell(row=i, column=1).value

        # get the tweet text in column 11
        text = workbook.cell(row=i, column=11).value

        # construct doc and add to the IndexWriter
        self.add_doc(id, text)

    # close the writer
    self.writer.close()
```

```
# function to add a new document with specific id and text to the IndexWriter
def add_doc(self, id, text):
    # create a new document
    doc = Document()

    # configure how metadata is stored in the index
    metaType = FieldType()
    metaType.setStored(True)
    metaType.setTokenized(True)

    # configure how content data is stored in the index
    # store the doc id, term frequency, and positions
    contentType = FieldType()
    contentType.setIndexOptions(IndexOptions.DOCS_AND_FREQS_AND_POSITIONS)
    contentType.setStoreTermVectors(True)
    contentType.setStoreTermVectorPositions(True)
    contentType.setStored(True)
    contentType.setTokenized(True)

    # add the title and content field to the document
    doc.add(Field("id", id, metaType))
    doc.add(Field("text", text, contentType))

    # add the document to the IndexWriter
    self.writer.addDocument(doc)
```

```
def build_traditional_index(self):
    print("----- Traditional Indexing -----")
    self.index = collections.defaultdict(lambda: collections.defaultdict(list))

    # record the start time
    start_time = time.time()

    # build index using PyLucene
    self.build_lucene_index()

    # iterate through each document in the lucene index
    self.doc_list = self.get_doc_list()
    for doc in self.doc_list:

        # get the doc_id of the current document
        doc_id = doc.doc

        # get the term vector of the current document
        term_vector = self.reader.getTermVector(doc_id, "text")

        # check if there is term in the document
        if term_vector is not None:
            term_iter = BytesRefIterator.cast_(term_vector.iterator())

            # iterate through each term in the document
            while term_iter.next():
                terms_enum = TermsEnum.cast_(term_iter)

                # convert the term from UTF-8 byte format to string
                term = terms_enum.term().utf8ToString()

                # get the postings list of the term
                postings = terms_enum.postings(None, PostingsEnum.ALL)

                # iterate through each posting
                if postings.nextDoc() is not DocIdSetIterator.NO_MORE_DOCS:

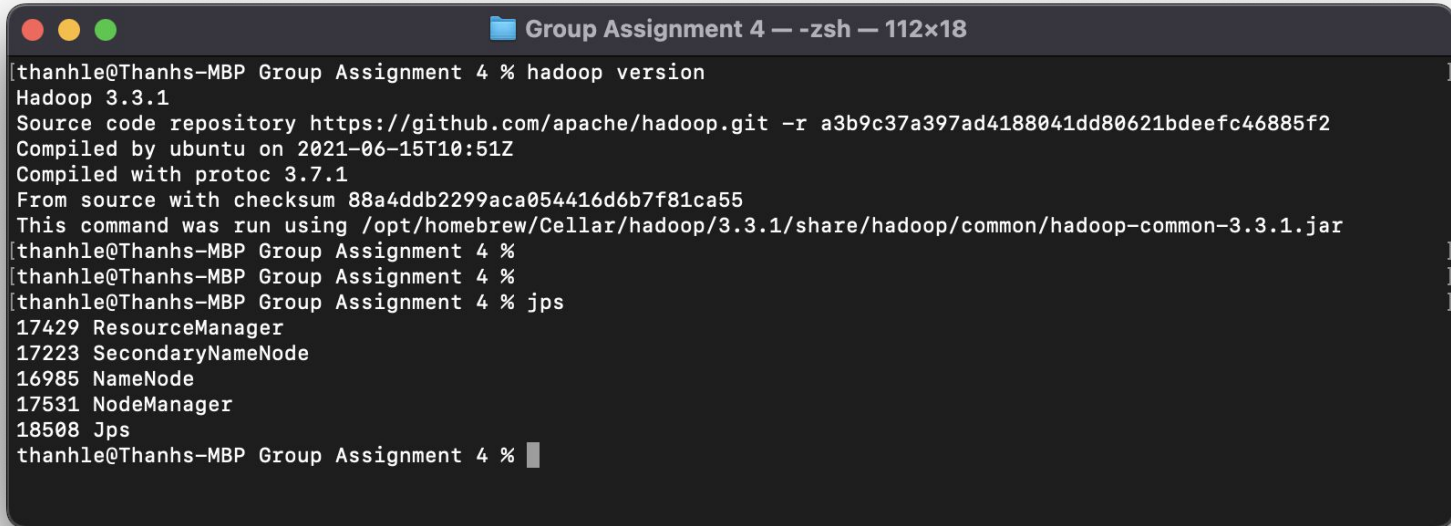
                    # get the term frequency
                    freq = postings.freq()

                    # record the term frequency in the index
                    self.index[term][doc_id] = freq

    # record the end time
    end_time = time.time()
```

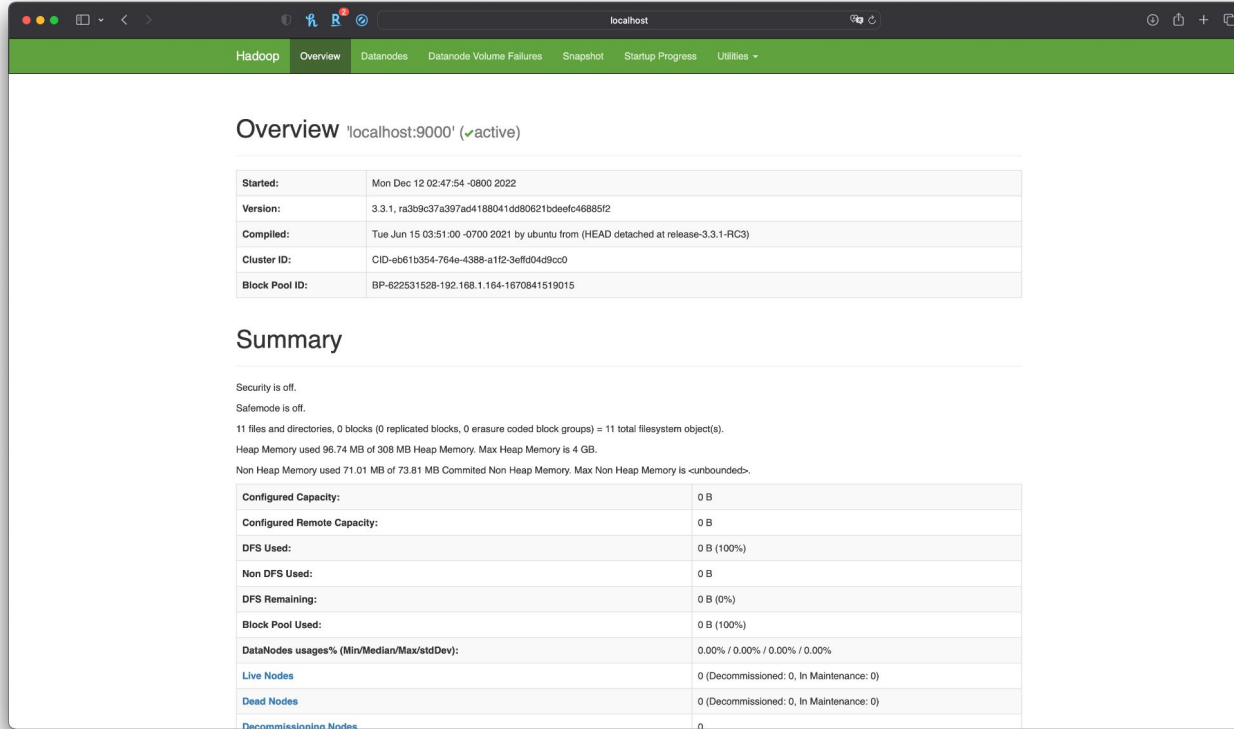


# INSTALLING HADOOP

A terminal window titled "Group Assignment 4 — -zsh — 112x18" with standard macOS window controls (red, yellow, green buttons). The terminal shows the command "hadoop version" being executed. The output provides details about the Hadoop installation, including the version (3.3.1), source code repository, compilation date, and checksum. It also lists the processes running in the background: ResourceManager, SecondaryNameNode, NameNode, NodeManager, and Jps.

```
[thanhle@Thanks-MBP Group Assignment 4 % hadoop version]
Hadoop 3.3.1
Source code repository https://github.com/apache/hadoop.git -r a3b9c37a397ad4188041dd80621bdeefc46885f2
Compiled by ubuntu on 2021-06-15T10:51Z
Compiled with protoc 3.7.1
From source with checksum 88a4ddb2299aca054416d6b7f81ca55
This command was run using /opt/homebrew/Cellar/hadoop/3.3.1/share/hadoop/common/hadoop-common-3.3.1.jar
[thanhle@Thanks-MBP Group Assignment 4 %]
[thanhle@Thanks-MBP Group Assignment 4 %]
[thanhle@Thanks-MBP Group Assignment 4 % jps]
17429 ResourceManager
17223 SecondaryNameNode
16985 NameNode
17531 NodeManager
18508 Jps
[thanhle@Thanks-MBP Group Assignment 4 %]
```

# SETTING HADOOP CLUSTER LOCALLY



The screenshot displays the Hadoop web interface in a browser window titled 'localhost'. The interface has a green navigation bar with tabs: 'Hadoop', 'Overview', 'Datanodes', 'Datanode Volume Failures', 'Snapshot', 'Startup Progress', and 'Utilities'. The 'Overview' tab is selected, showing the status of the 'localhost:9000' cluster, which is marked as '(active)'.

**Overview 'localhost:9000' (✓active)**

Started:	Mon Dec 12 02:47:54 -0800 2022
Version:	3.3.1, ra3b9c37a397ad4188041dd90621bdeefc46885f2
Compiled:	Tue Jun 15 03:51:00 -0700 2021 by ubuntu from (HEAD detached at release-3.3.1-RC3)
Cluster ID:	CID-eb61b354-764e-4388-a1f2-3eff0d49cc0
Block Pool ID:	BP-622531528-192.168.1.164-1670841519015

**Summary**

Security is off.  
Safemode is off.

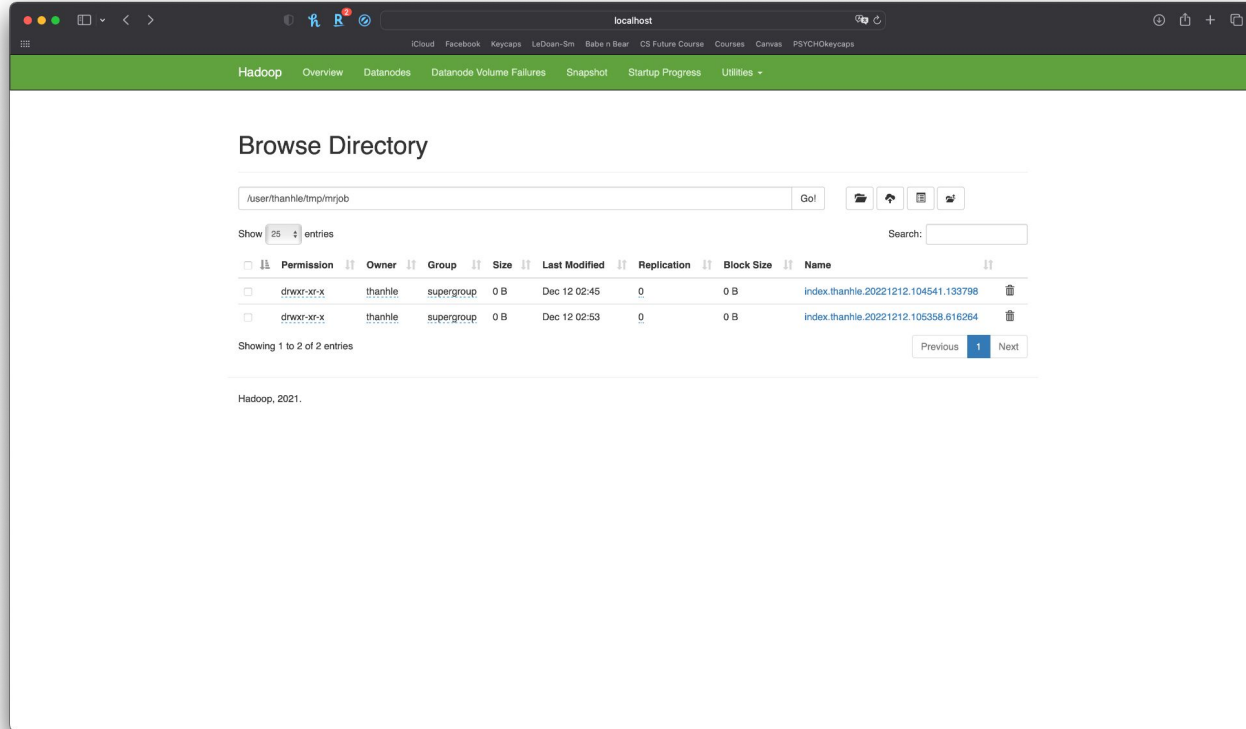
11 files and directories, 0 blocks (0 replicated blocks, 0 erasure coded block groups) = 11 total filesystem object(s).

Heap Memory used 96.74 MB of 308 MB Heap Memory. Max Heap Memory is 4 GB.

Non Heap Memory used 71.01 MB of 73.81 MB Committed Non Heap Memory. Max Non Heap Memory is <unbounded>.

Configured Capacity:	0 B
Configured Remote Capacity:	0 B
DFS Used:	0 B (100%)
Non DFS Used:	0 B
DFS Remaining:	0 B (0%)
Block Pool Used:	0 B (100%)
DataNodes usages% (Min/Median/Max/stdDev):	0.00% / 0.00% / 0.00% / 0.00%
<a href="#">Live Nodes</a>	0 (Decommissioned: 0, In Maintenance: 0)
<a href="#">Dead Nodes</a>	0 (Decommissioned: 0, In Maintenance: 0)
<a href="#">Decommissioning Nodes</a>	0

# UPLOADING FILE TO HADOOP



# LIBRARIES

```
# Python 3.0
import re
import sys
import os
import collections
import time

import lucene
from java.io import File
from org.apache.lucene.document import Document, Field, FieldType
from org.apache.lucene.util import BytesRefIterator

from org.apache.lucene.index import IndexWriter, IndexWriterConfig, PostingsEnum, IndexOptions, TermsEnum
from org.apache.lucene.store import FSDirectory

from org.apache.lucene.search import IndexSearcher, MatchAllDocsQuery, DocIdSetIterator
from org.apache.lucene.index import DirectoryReader

import openpyxl
from mrjob.job import MRJob
from mrjob.step import MRStep
```

- Openpyxl: python library to open and read excel files for tweets dataset
- MRjob: python library for MapReduce developed by YELP

# MAP REDUCE

Using MapReduce operation  
to count the word frequency  
of each tweets

```
# regex to find non-alphabetical words
WORD_RE = re.compile(r"[\w']+")
```

```
# map-reduce job to count word frequency
class word_count(MRJob):
```

```
    def mapper_raw(self, input_path, input_uri):
        # iterate through each line in the dataset file
        workbook = openpyxl.load_workbook(input_path).active
        for i in range(3, workbook.max_row + 1, 2):

            # get the tweet id in column 1
            id = workbook.cell(row=i, column=1).value

            # get the tweet text in column 11
            text = workbook.cell(row=i, column=11).value

            # iterate through each term in the text and remove non-alphabetically words
            for term in WORD_RE.findall(text):

                # yield the ((id, term), term_frequency) pair
                yield ((id, term.lower()), 1)

    def reducer(self, key, value):
        # sum all term frequency with the same key (id, term)
        yield (key, sum(value))

    def steps(self):
        return [
            MRStep(
                mapper_raw=self.mapper_raw,
                reducer=self.reducer
            )
        ]
```

# BUILDING DISTRIBUTED INDEX

```
# function to build the distributed index using map reduce
def build_map_reduce_index(self):
    print("----- Distributed Indexing using Map Reduce -----")
    self.index = collections.defaultdict(lambda: collections.defaultdict(list))

    # record the start time
    start_time = time.time()

    # instantiate the map-reduce job to count word frequency running on local hadoop cluster
    job = word_count(args=['-r', 'hadoop', self.path])

    # run the map-reduce job
    with job.make_runner() as runner:
        runner.run()

        # parse the output of the job and record the data to the index
        for (doc_id, term), freq in job.parse_output(runner.cat_output()):
            self.index[term][doc_id] = freq

    # record the end time
    end_time = time.time()

    # print the indexing time
    print("TF-IDF Index built in", end_time - start_time, "seconds.\n")
```

# EXPERIMENTAL RESULT (WORD\_COUNT OUTPUT)

["1595256559616990000", "colleagues"]	1
["1595256559616990000", "continue"]	1
["1595256559616990000", "covid"]	1
["1595256559616990000", "family"]	1
["1595256559616990000", "for"]	1
["1595256559616990000", "get"]	1
["1595256559616990000", "huh"]	1
["1595256559616990000", "in"]	1
["1595256559616990000", "message"]	1
["1595256559616990000", "millisecond"]	1
["1595256559616990000", "never"]	1
["1595256559616990000", "one"]	1
["1595256559616990000", "over"]	1
["1595256559616990000", "real"]	1
["1595256559616990000", "said"]	1
["1595256559616990000", "should"]	1
["1595256559616990000", "so"]	1
["1595256559616990000", "that"]	1
["1595256559616990000", "the"]	2
["1595256559616990000", "thug"]	1
["1595256559616990000", "to"]	2
["1595256559616990000", "uncledanny58"]	1
["1595256559616990000", "ve"]	1
["1595256559616990000", "we"]	1
["1595256559616990000", "worked"]	1
["1595256559616990000", "would"]	1
["1595256559616990000", "you"]	1

["1595256560220960000", "comparable"]	1
["1595256560220960000", "covid"]	1
["1595256560220960000", "dear"]	1
["1595256560220960000", "document"]	1
["1595256560220960000", "donald"]	1
["1595256560220960000", "election"]	1
["1595256560220960000", "insurrection"]	1
["1595256560220960000", "is"]	1
["1595256560220960000", "lies"]	4
["1595256560220960000", "miles"]	1
["1595256560220960000", "million"]	1
["1595256560220960000", "ncomparing"]	1
["1595256560220960000", "ndo"]	2
["1595256560220960000", "nnot"]	1
["1595256560220960000", "not"]	1
["1595256560220960000", "nstop"]	1
["1595256560220960000", "nwe"]	1
["1595256560220960000", "press"]	1
["1595256560220960000", "related"]	1
["1595256560220960000", "remotely"]	1
["1595256560220960000", "sides\u00e2"]	1
["1595256560220960000", "stolen"]	1
["1595256560220960000", "to"]	2
["1595256560220960000", "trillion"]	1
["1595256560220960000", "trump"]	1
["1595256560220960000", "trying"]	1
["1595256560220960000", "wedding"]	1

["1595256561395730000", "daily"]	1
["1595256561395730000", "died suddenly"]	1
["1595256561395730000", "end"]	1
["1595256561395730000", "enough"]	1
["1595256561395730000", "ever"]	1
["1595256561395730000", "everyone"]	1
["1595256561395730000", "got"]	1
["1595256561395730000", "happening"]	1
["1595256561395730000", "have"]	1
["1595256561395730000", "hear"]	1
["1595256561395730000", "heart"]	1
["1595256561395730000", "i"]	2
["1595256561395730000", "know"]	1
["1595256561395730000", "lucky"]	1
["1595256561395730000", "me"]	1
["1595256561395730000", "multiple"]	1
["1595256561395730000", "none"]	1
["1595256561395730000", "oddly"]	1
["1595256561395730000", "of"]	1
["1595256561395730000", "or"]	1
["1595256561395730000", "over"]	1
["1595256561395730000", "shot"]	1
["1595256561395730000", "sicker"]	1
["1595256561395730000", "stars"]	1
["1595256561395730000", "stroke"]	1
["1595256561395730000", "than"]	1
["1595256561395730000", "the"]	1

# EXPERIMENTAL RESULT (RUNTIME)

```
----- Traditional Indexing -----  
TF-IDF Index built in 16.255134105682373 seconds.
```

```
----- Distributed Indexing using Map Reduce -----  
No configs specified for inline runner  
TF-IDF Index built in 12.771775245666504 seconds.
```

```
----- Traditional Indexing -----  
TF-IDF Index built in 16.14408802986145 seconds.
```

```
----- Distributed Indexing using Map Reduce -----  
No configs specified for inline runner  
TF-IDF Index built in 12.734126091003418 seconds.
```

The distributed index results in a faster runtime on multiple tries compared to the traditional index



# TRADITIONAL VS. DISTRIBUTED INDEXING

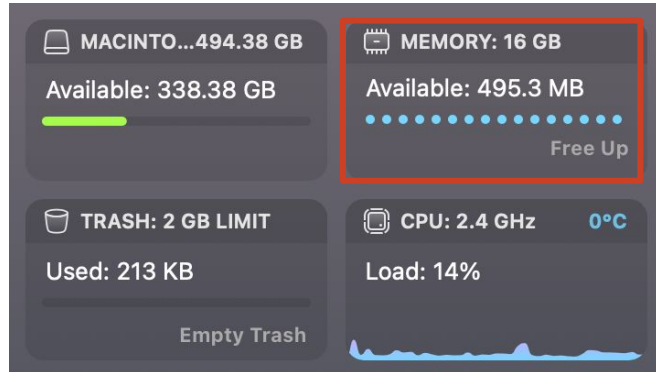
The data set given was around 15.3 MB or about 25,013 tweets.

The data we have shown with both traditional and distributed processing showed runtimes of about 16 seconds and 13 seconds respectively. This makes the speedup approximately 1.23 times faster for distributed compared to traditional.

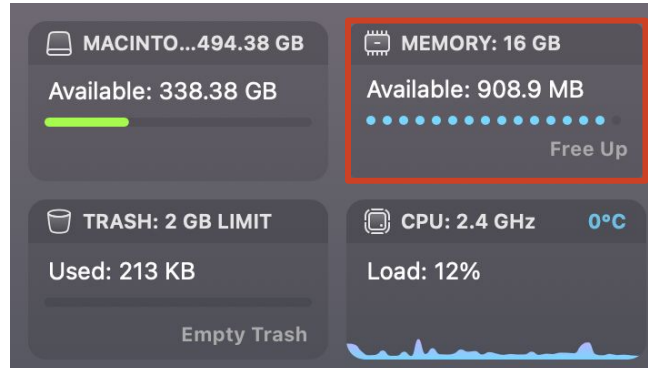
If the dataset were multitudes larger, distributed processing would be needed as distributed processing was designed for scalability, since it runs the data in parallel on different nodes/computers in a system rather than on a single computer in traditional processing.

# EXPERIMENTAL RESULT (MEMORY USAGE)

- Free up the RAM memory before each run.
- Keep track of RAM memory usage during indexing with Map Reduce and Hadoop.



During indexing



After indexing

Much memory is used during indexing. The memory is released after usage.

# CONCLUSION

- The distributed index built using Map Reducing running on Hadoop cluster shows better performance in run time.
- Analysis of speed up:
  - Hadoop optimizes memory and disk spill use. It aims to use as much memory as possible without triggering swapping → memory tuning to maximize the performance
  - Each MapReduce job in Hadoop collects the information about the various input records read, number of reducer records, number of records pipelined for further execution, swap memory, heap size set, etc.
  - Hadoop controls the amount of mapper and the size of each job. When dealing with large files, Hadoop splits the file into smaller chunks so that mapper can run it in parallel → speeding up the process.

# TASK DIVISION

- Both of us first started working on implementing the index.py file by ourselves and discussed with each other when any problem happened
- We compared our query results with each other to make sure we got the same outputs
- We combined our works into a final version with detailed comments in each step which helps explaining the codes
- Antoine created the output.txt file
- Thanh commented the codes and submitted the zip file
- We worked on the presentation slides together