

VIETNAM NATIONAL UNIVERSITY
HO CHI MINH UNIVERSITY OF SCIENCE



Course: CSC14118 - INTRODUCTION TO BIG DATA

Lab 02: Document Clustering with Hadoop MapReduce

Group: 21KHDL08

Members: Hồ Đinh Duy Lực - 21127351

Tô Khánh Linh - 21127638

Đoàn Thị Yên Nhi - 21127660

Lê Mỹ Khánh Quỳnh - 21127681

Class: 21KHDL

Instructors: Mr. Lê Ngọc Thành

Ms. Nguyễn Ngọc Thảo

Mr. Bùi Huỳnh Trung Nam

Mr. Đỗ Trọng Lê

August 2024, Ho Chi Minh City

TABLE OF CONTENT

1. Group information.....	3
2. Tasks Evaluation	3
3. Part 1: Data Preprocessing	4
3.1. Task 1.1 Text Cleaning and Term Frequency	4
3.2. Task 1.2 Low-Frequency Term Elimination	13
3.3. Task 1.3 Top 10 Most Frequent Words	20
3.4. Task 1.4 TF-IDF	24
3.5. Task 1.5 Highest average TF-IDF	32
4. Part 2: K-Means Algorithm	39
4.1. Task 2.1 K-Means on 2D Data	39
4.2. Task 2.2 K-Means on Preprocessed Data.....	47
4.3. Task 2.3 Scalable K-Means++ Initialization	57
5. Conclusion.....	61
6. References	61

1. Group information

Student ID	Fullname	Group ID
21127638	Tô Khánh Linh	KHDL- 08
21127660	Đoàn Thị Yên Nhi	
21127351	Hồ Đinh Duy Lực	
21127681	Lê Mỹ Khánh Quỳnh	

2. Tasks Evaluation

	Task	Completed
Part 1: Data Preprocessing	1.1. Text Cleaning and Term Frequency	100%
	1.2. Low-Frequency Term Elimination	100%
	1.3. Top 10 Most Frequent Words	100%
	1.4. TF-IDF	100%
	1.5. Highest average TF-IDF	100%
Part 2: K-Means Algorithm	2.1. K-Means on 2D Data	100%
	2.2. K-Means on Preprocessed Data	100%
	2.3. Scalable K-Means++ Initialization	100%

3. Part 1: Data Preprocessing

3.1. Task 1.1 Text Cleaning and Term Frequency

a. **Data Description:** In this section, we use the following data files:

- The input consists of files organized in five folders within the unprocessed `bbc_dataset` directory. These folders are named `business`, `entertainment`, `politics`, `sport`, and `tech`. Each folder contains multiple `.txt` files, sequentially labeled from `001.txt` to `N.txt`.
- There is also a stopwords file (`stopwords.txt`) used to filter out stop words from the text files.
- Additionally, the processing code requires reading two files: `bbc.terms` and `bbc.docs`, for the following purposes:
 - `bbc.terms`: This file is used to assign an ID to each term. After removing stop words, each term is given an ID based on its line number.
 - `bbc.docs`: This file is used to assign an ID to each text file within a specific folder, with each file's ID corresponding to its line number.

b. **MapReduce Job Implementation**

In this task, we use **1 Mapper** phase and **1 Reducer** phase. Specifically, the execution steps are as follows:

- **In Mapper phase:**

- This mapper processes each line of a text file, cleans and filters the words (remove stop words and regex), and then outputs key-value pairs where the key is a combination of the term ID and the document ID, and the value is always 1.

Step 1: Before map (set up):

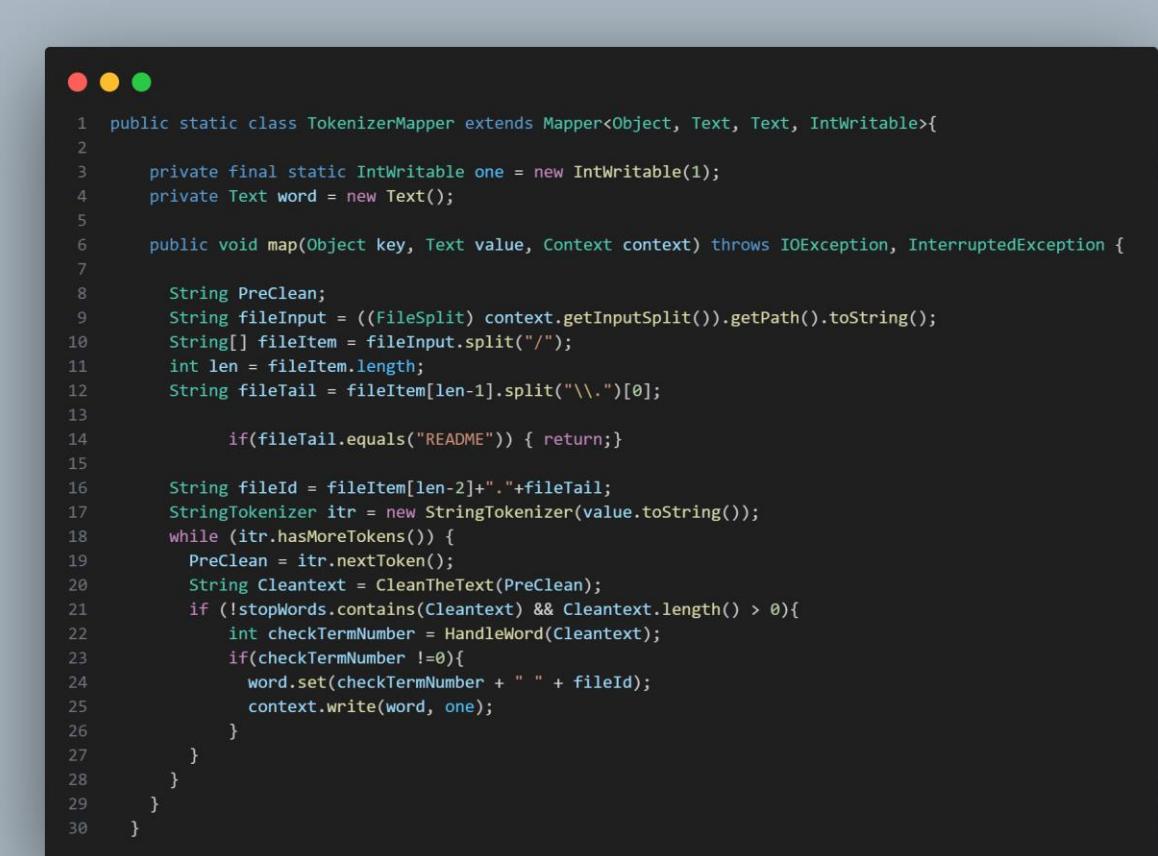
- During the map phase, before mapping the data, the code will read and process auxiliary files as follows:

- For the stopwords file, it will extract words from the stopwords file and add them to a string set.
- For the bbc.terms, it will extract the words from the file and read it into a `HashMap<String, Integer>` which `String` indicates the word in `bbc.term`, `Integer` indicates the line of that word.
- For the bbc.docs, it will extract the documents id from the file into to into a `HashMap<String, Integer>` `String` is the name of the file and `Integer` indicates the line of that file.

Step 2: Mapping:

- **File Identification:** The mapper retrieves the file path of the current input split using `context.getInputSplit()`. Then it extracts the file name (`fileTail`) and its parent directory name (`fileItem[len-2]`), combining them to form a `fileId` (e.g `folder.filename`).
- **Skipping README Files:** If the file name is "README", the mapper skips by returning early.
- **Tokenizing and Processing the Text:** Each word is cleaned by removing special characters and converting it to lowercase using the `CleanTheText` method. The cleaned word is then checked against a list of stopwords and its length, if it's not a stopword and has a length greater than zero, it's can be process further
- **Handling the Word:** The cleaned word is passed to `HandleWord`, which checks if it exists in the `termId` dictionary (a mapping of words to term IDs). If the word exists in `termId`, its corresponding term ID (line of that word) is returned. If it doesn't exist, 0 is returned.
- **Emitting Key-Value Pairs:** If the word has a valid term ID (`checkTermNumber != 0`), the mapper creates a key by combining the term ID and `fileId`.
 - **Key:** `checkTermNumber` (the ID of the word) + " " + `fileId` (the ID of the file which is `folder.filename` (e.g `bussiness.001`))

- **Value:** 1
- Example output: [K,V] [1 business.001, 1]



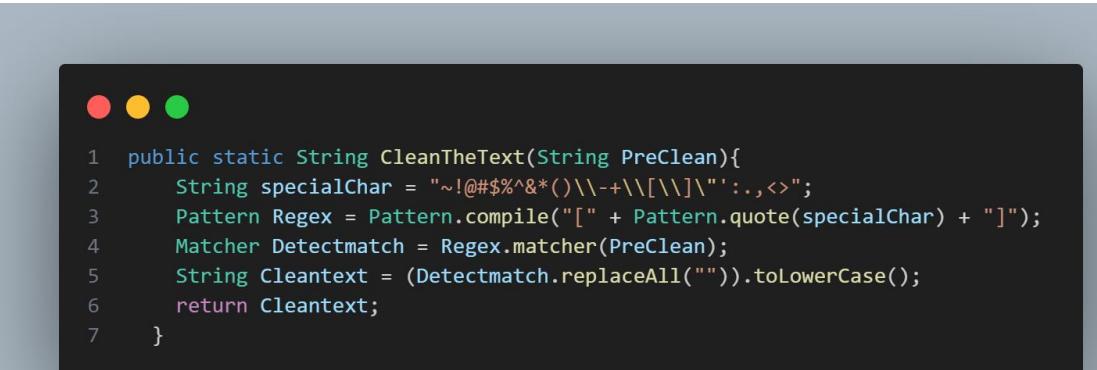
```

1  public static class TokenizerMapper extends Mapper<Object, Text, Text, IntWritable>{
2
3      private final static IntWritable one = new IntWritable(1);
4      private Text word = new Text();
5
6      public void map(Object key, Text value, Context context) throws IOException, InterruptedException {
7
8          String PreClean;
9          String fileInput = ((FileSplit) context.getInputSplit()).getPath().toString();
10         String[] fileItem = fileInput.split("/");
11         int len = fileItem.length;
12         String fileTail = fileItem[len-1].split("\\.")[0];
13
14         if(fileTail.equals("README")) { return;}
15
16         String fileId = fileItem[len-2]+". "+fileTail;
17         StringTokenizer itr = new StringTokenizer(value.toString());
18         while (itr.hasMoreTokens()) {
19             PreClean = itr.nextToken();
20             String Cleantext = CleanTheText(PreClean);
21             if (!stopWords.contains(Cleantext) && Cleantext.length() > 0){
22                 int checkTermNumber = HandleWord(Cleantext);
23                 if(checkTermNumber !=0){
24                     word.set(checkTermNumber + " " + fileId);
25                     context.write(word, one);
26                 }
27             }
28         }
29     }
30 }

```

Figure 1: Map code

- CleanTheText:



```

1  public static String CleanTheText(String PreClean){
2      String specialChar = "~!@#$%^&*()'\\-+\\[\\]\\\"':,<>";
3      Pattern Regex = Pattern.compile("[ " + Pattern.quote(specialChar) + "]");
4      Matcher Detectmatch = Regex.matcher(PreClean);
5      String Cleantext = (Detectmatch.replaceAll("")).toLowerCase();
6      return Cleantext;
7  }

```

Figure 2: CleanTheText method

- HandleWord:

```

1 public static int HandleWord(String Cleantext){
2     Integer posTerm = termId.get(Cleantext);
3     if(posTerm != null){ return posTerm; }
4     return 0;
5 }
```

Figure 3: HandleWord method

- **In Reducer phase:**

Step 3: Reducing:

- During the reduce phase, the process is straightforward. When key-value pairs with the same key are encountered, their values are summed up to generate the final value. This step is essential because certain words may appear multiple times within a single file, making aggregation necessary.
- Once the reduce phase is completed, the result will be a key-value pair where the key remains the same as in the map phase, and the value reflects the total count of a word's occurrences within a file. Thus, the key-value format will be:
 - **Key:** checkTermNumber (the ID of the word) + " " + fileId (the ID of the file which is folder.filename (e.g business.001))
 - **Value:** total number of occurrences in the file

```

1 public static class IntSumReducer extends Reducer<Text,IntWritable,Text,IntWritable> {
2     private IntWritable result = new IntWritable();
3     public void reduce(Text key, Iterable<IntWritable> values,Context context) throws IOException, InterruptedException {
4         int sum = 0;
5         for (IntWritable val : values) {
6             sum += val.get();
7         }
8         result.set(sum);
9         context.write(key, result);
10    }
11 }
```

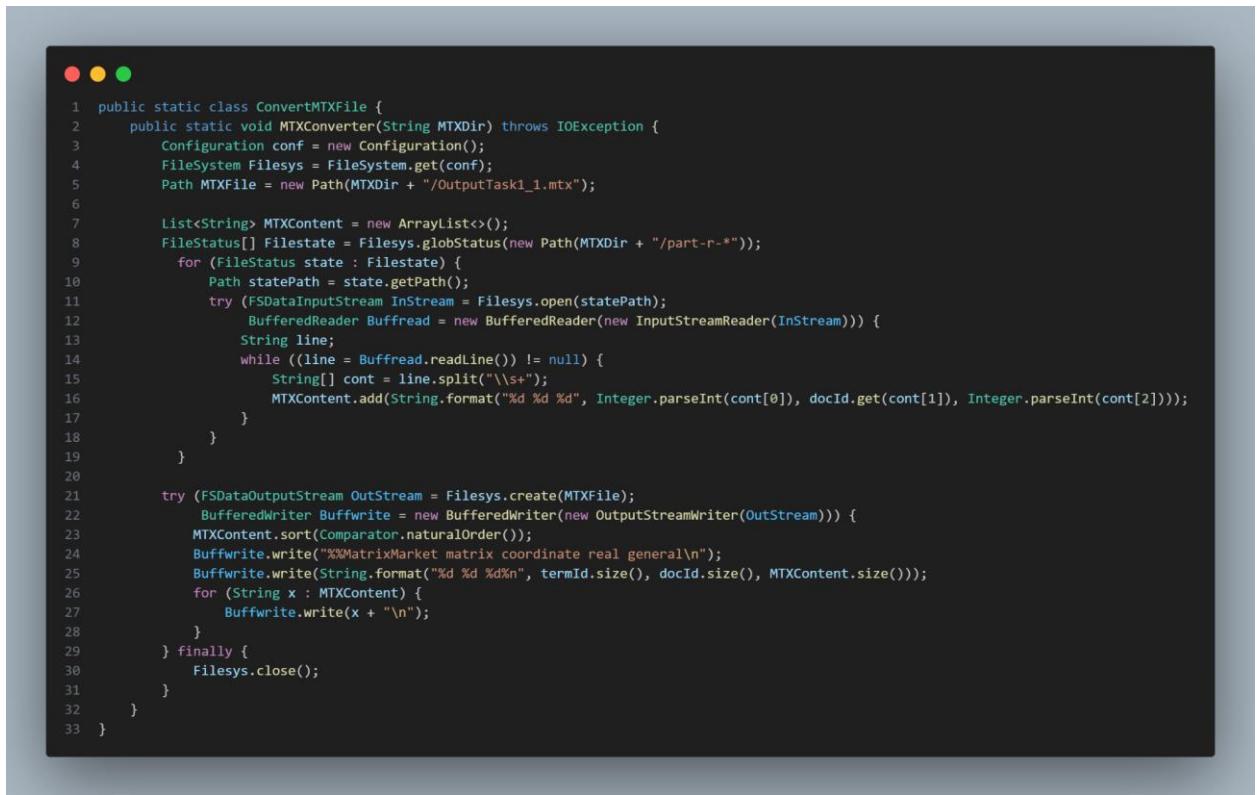
Figure 4: Reduce code

- **MTX converter phase:** This method essentially converts and aggregates the output of a MapReduce job into a single, structured Matrix Market file

Step 4: Output formatting:

- After reducing, we will continue converting and reformatting the output file to MTX file by calling the MTXConverter method in the ConvertMTXFile by passing in the output directory location.
- **Configuration Setup:** The method starts by creating a Hadoop Configuration object and obtaining a reference to the file system (FileSystem Filesys) using that configuration.
- **File Path Initialization:** It then defines the path for the output .mtx file, MTXFile, which will store the final converted data. The file will be named OutputTask1_1.mtx and stored in the specified directory (MTXDir).
- **Reading Input Files:** The method searches for all files in the MTXDir directory that match the pattern part-r-*. These files typically represent the output of the reduce tasks in a MapReduce job. For each matching file, it opens the file using FSDataInputStream and reads its content line by line using a BufferedReader.
- **Processing and Formatting Data:** Each line of the input files is split into components which is term ID, document ID, and frequency. The document ID is mapped using a docId mapping (business.001 have docId = 1), and the components are then formatted into a string of three integers which are termId and docId, frequency. These formatted strings are added to the MTXContent list.
- **Writing the .mtx File:** After reading all the input files and storing the formatted data in MTXContent, the list is sorted. The method then creates the output .mtx file and writes the header information. Finally, it writes each formatted line from MTXContent to the .mtx file.

- When finishing the Reducer phase, the program will then run a function to reformat the structure of the output and return an output file with 3 columns of data(TermId DocId Frequency), namely:
 - TermId: the ID of the term
 - DocId: the ID of the file
 - Frequency: the number of occurrences of the term in a file



```

1  public static class ConvertMTXFile {
2      public static void MTXConverter(String MTXDir) throws IOException {
3          Configuration conf = new Configuration();
4          FileSystem Filesys = FileSystem.get(conf);
5          Path MTXFile = new Path(MTXDir + "/OutputTask1_1 mtx");
6
7          List<String> MTXContent = new ArrayList<>();
8          FileStatus[] Filestate = Filesys.globStatus(new Path(MTXDir + "/part-r-*"));
9          for (FileStatus state : Filestate) {
10              Path statePath = state.getPath();
11              try (FSDataInputStream InStream = Filesys.open(statePath)) {
12                  BufferedReader Buffread = new BufferedReader(new InputStreamReader(InStream));
13                  String line;
14                  while ((line = Buffread.readLine()) != null) {
15                      String[] cont = line.split("\\s+");
16                      MTXContent.add(String.format("%d %d %d", Integer.parseInt(cont[0]), docId.get(cont[1]), Integer.parseInt(cont[2])));
17                  }
18              }
19          }
20
21          try (FSDataOutputStream OutStream = Filesys.create(MTXFile);
22               BufferedWriter Buffwrite = new BufferedWriter(new OutputStreamWriter(OutStream))) {
23              MTXContent.sort(Comparator.naturalOrder());
24              Buffwrite.write("%XMMatrixMarket matrix coordinate real general\n");
25              Buffwrite.write(String.format("%d %d %d", termId.size(), docId.size(), MTXContent.size()));
26              for (String x : MTXContent) {
27                  Buffwrite.write(x + "\n");
28              }
29          } finally {
30              Filesys.close();
31          }
32      }
33  }

```

Figure 5: Output converter code

c. Task Flow:

- The task run with this flow:

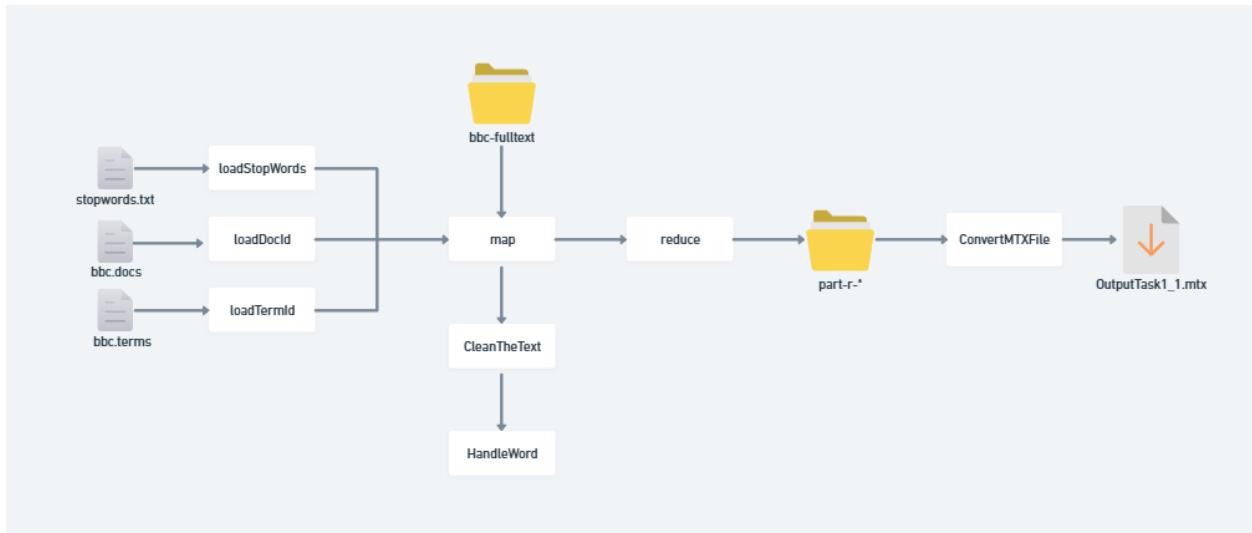


Figure 6: Figures describe task 1.1 flow

- **An example for this task:**

- This output file will be named "task_1_2.mtx"
- For example, we have a file name 001.txt, and it put in folder Business, and the content in file is: “The Ad Sales boost Time ad”.
 - We can see that the word “The ” is stop word and if in file terms, we have that:
 - ad in line 1
 - sales in line 2
 - boost in line 3
 - time in line 4
- In file docs, we have that Business.001 is in the line 1
- We will describe an example for this task as the below figure:



Figure 7: Figures describe map-reduce job

- From the input file, it's clear that the word "the" is a stopword and will be removed. The remaining words will be paired with the file-folder ID to create key-value pairs. During the reduction process, these pairs will be aggregated into a key in the format (termId foldername.fileid), where the value represents the total number of occurrences.

d. How to run:

- Configure this section:

- We need to update the file paths according to your browser. Passing the stopwords file path to the loadStopWords, bbc.docs file path to the loadDocId , and bbc.terms file path to the loadTermId.

```

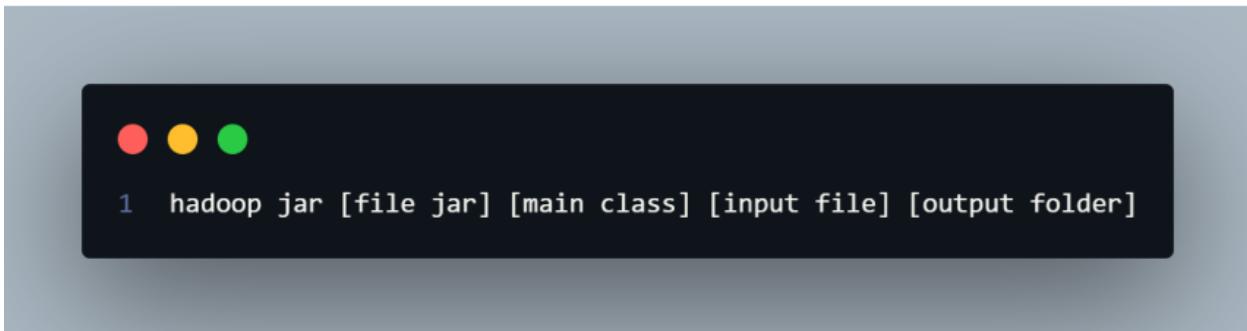
● ● ●

1 loadStopWords("./input/stopwords.txt");
2 loadDocId("./input/bbc.docs");
3 loadTermId("./input/bbc.terms");

```

Figure 8 : Configure Path

- **Command to run the program:** After completing the program, the results will return to the task_1_1.mtx file:

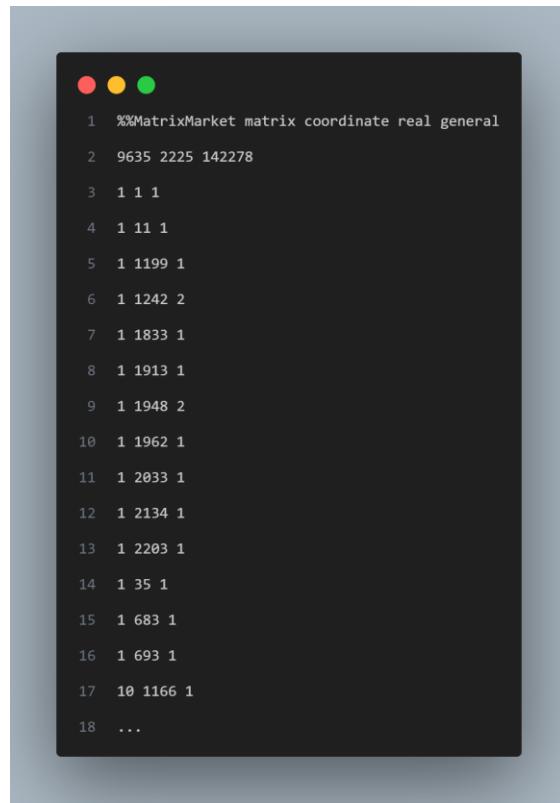


```
1 hadoop jar [file jar] [main class] [input file] [output folder]
```

Figure 9: Command

e. Result:

- Sample of file task_1_1.mtx:



```
1 %%MatrixMarket matrix coordinate real general
2 9635 2225 142278
3 1 1 1
4 1 11 1
5 1 1199 1
6 1 1242 2
7 1 1833 1
8 1 1913 1
9 1 1948 2
10 1 1962 1
11 1 2033 1
12 1 2134 1
13 1 2203 1
14 1 35 1
15 1 683 1
16 1 693 1
17 10 1166 1
18 ...
```

Figure 10: Sample of output task 1.1

- From the result file, we see that a file that meets the requirements is returned, including 3 columns: terms, docs and frequency.

3.2. Task 1.2 Low-Frequency Term Elimination

a. Data Description:

- MTX file generated from Task 1.1 named.
- Additionally, the processing code requires reading two files: bbc.terms and bbc.docs, for the following purposes:
 - bbc.terms: This file is used to assign an ID to each term. After removing stop words, each term is given an ID based on its line number.
 - bbc.docs: This file is used to assign an ID to each text file within a specific folder, with each file's ID corresponding to its line number.

b. MapReduce Job Implementation:

- In Mapper phase:

This mapper reads lines of text input, skips the first and second row which is the header of the MTX file, it splits the line into tokens. It then emits the first token as the key (termId) and the combination of the second and third tokens as the value (docId and frequency).

Step 1: Before map (set up):

- During the map phase, before mapping the data, the code will read and process auxiliary files as follows:
 - For the bbc.terms, it will extract the words from the file and read it into a `HashMap<String, Integer>` which String indicates the word in `bbc.terms`, Integer indicates the line of the word.
 - For the bbc.docs, it will extract the documents id from the file into to into a `HashMap<String, Integer>` String is the name of the file and Integer indicates the line of the file.

Step 2: Mapping:

- **Skipping the header of the MTX file:** The method first checks if row is less than 2. If it is, it increments the row counter and skips processing the current line. This effectively skips the first row of the input data.
- **Processing the Line:** It split into tokens using whitespace as the delimiter (`split("\\s+")`). The first token (`cont[0]`) is set as the termId. The second and

third tokens (cont[1] and cont[2]) which are docId and frequency are combined into a single string separated by a space and set as pairVal.

- **Emitting Key-Value Pairs:** The method then emits a key-value pair using context.write(termId, pairVal). Here, termId is the key, and pairVal is the value:
 - **Key:** termId (the ID of the word)
 - **Value:** pairVal (the ID of the doc +” “+ frequency)

```

1  public static class TokenizerMapper extends Mapper<Object, Text, Text, Text>{
2      private Text termId = new Text();
3      private Text pairVal = new Text();
4      private int row=1;
5
6      public void map(Object key, Text value, Context context) throws IOException, InterruptedException {
7          if(row<2) {
8              row++;
9          }
10         else {
11             String [] cont = value.toString().split("\\s+");
12             termId.set(cont[0]);
13             pairVal.set(cont[1]+ " "+cont[2]);
14             context.write(termId,pairVal);
15         }
16     }
17 }

```

Figure 11: Map code

- **In Reducer phase:** This reducer takes a group of values for each key and adds up the numbers of the frequency of the termId from those values. If the total is 3 or more, it writes out each value again with the same key.
 - **Reduce Method:** Iteration Over Values:
 - The method iterates over each Text value in the values iterable which is the docId and frequency. Each value is converted to a string and added to the PairVal list.
 - The string value is then split into its components using whitespace as the delimiter (split("\\s+")). So it can take the second (index 1) component which is the frequency of the termId to add it up.

- The second component (index 1 which is the **frequency** of the **termId** in the **docId**) of the split string is parsed into an integer (Valadd) and added to sum.
- **Conditional:** After iterating through all the values, the reducer checks if the add up sum (**frequency** of the **termId**) is greater than or equal to 3. If the condition is met, the reducer iterates over the PairVal list and writes each value back out as a key-value pair, with the original key and the value from PairVal.



```

1  public static class IntSumReducer
2      extends Reducer<Text,Text,Text,Text> {
3          private Text result = new Text();
4
5          public void reduce(Text key, Iterable<Text> values, Context context )
6              throws IOException, InterruptedException {
7
8              ArrayList<String> PairVal = new ArrayList<>();
9              Integer sum = 0;
10             for (Text val : values) {
11                 String strVal = val.toString();
12                 PairVal.add(strVal);
13                 String ValCon[] = strVal.split("\\s+");
14                 Integer Valadd = Integer.parseInt(ValCon[1]);
15                 sum = sum + Valadd;
16             }
17             if (sum>=3){
18                 for (String val : PairVal) {
19                     result.set(val);
20                     context.write(key, result);
21                 }
22             }
23         }
24     }

```

Figure 12: Reduce code

- **MTX converter phase:** This method essentially converts and aggregates the output of a MapReduce job into a single, structured Matrix Market file
- Step 3: Output formatting:
- It has exactly the same code flow with task 1.1 but the output file would be OutputTask1_2.mtx.

c. Task Flow:

- The task run with this flow:

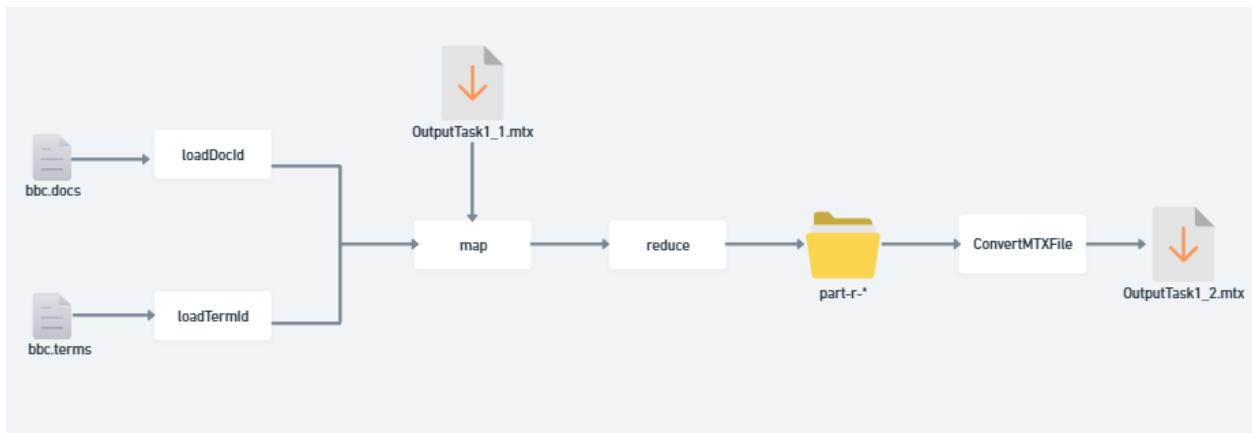


Figure 13: Figures describe task 1.2 flow

- An example for this task: We will describe an example for this task as the below figure:

- For example, we have an input file as follows:

1	1	1	3
2	1	2	3
3	2	1	1
4	3	1	2
5	3	2	1

Figure 14: Example input

- We can see that in the above file contains 3 columns show:
 - Word has termId = 1 appear in file has docId is 1 and 2 with frequency is 3 and 3
 - Word has termId = 2 appear in file has docId is 1 and frequency is 1
 - Word has termId = 3 appear in file has docId is 1 and 2 with frequency is 2 and 1

- The description of the map reduce will be show in the figures below:

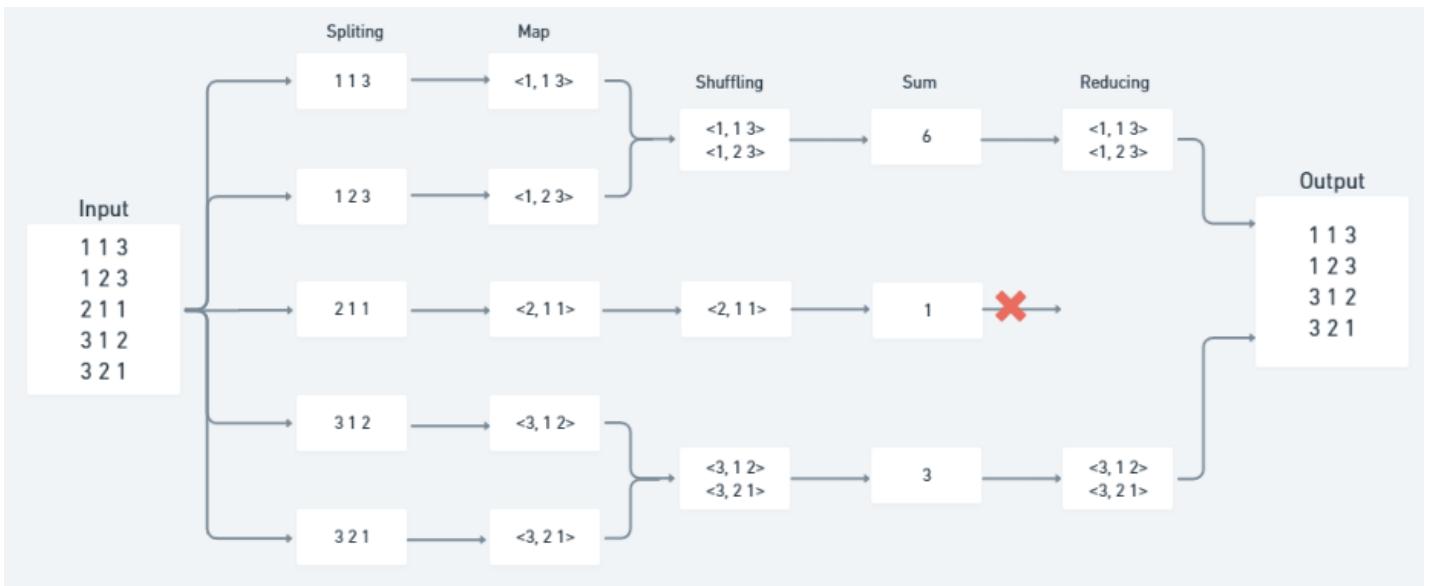


Figure 15: Figures describe map-reduce job

- Given the input, the word with termId = 2 does not meet the requirement of having a total occurrence count of at least 3. Consequently, after the map-reduce process, only the words with termId = 1 and termId = 3 satisfy the condition and remain.

d. How to run:

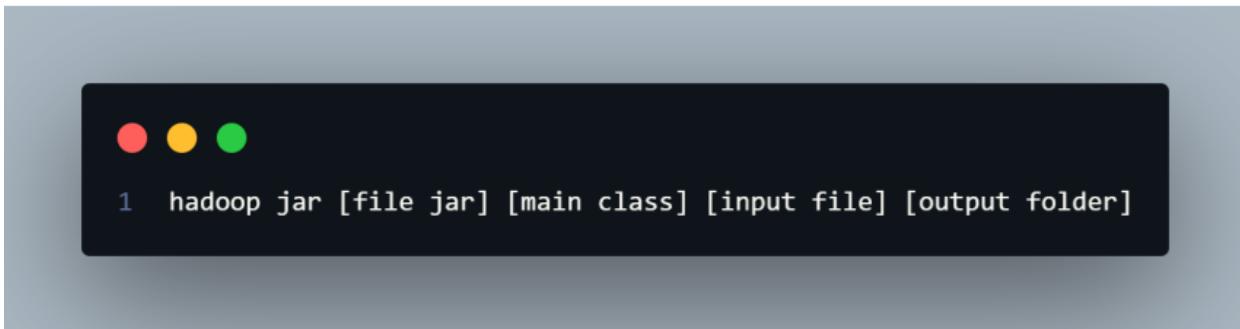
- Configure this section:

- We need to update the file paths according to your browser. Passing bbc.docs file path to the loadDocId , and bbc.terms file path to the loadTermId.



Figure 16: Configure Path

- **Command to run the program:** After completing the program, the results will return to the OutputTask1_1.mtx file:

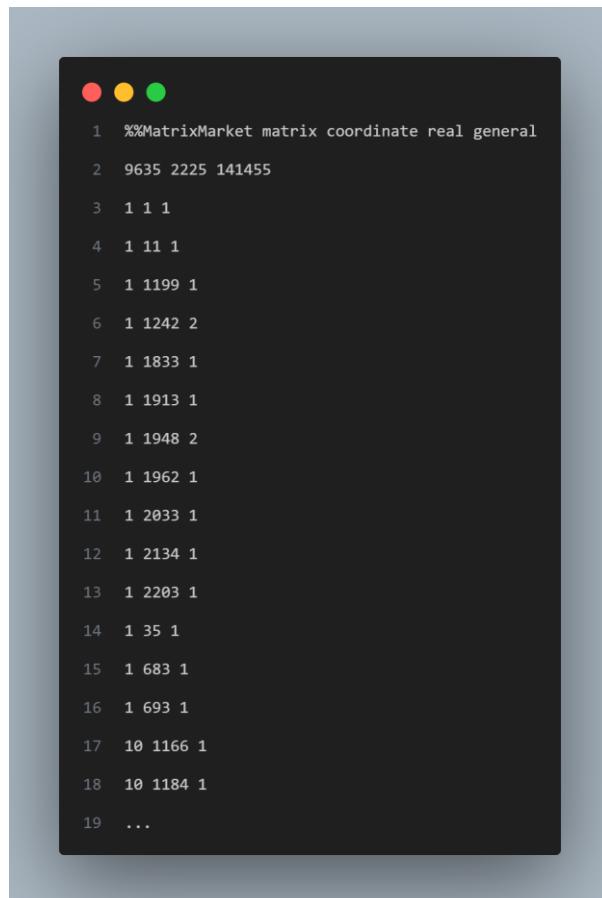


```
1 hadoop jar [file jar] [main class] [input file] [output folder]
```

Figure 17: Command

e. Results:

- Sample of file task_1_2.mtx:



```
1 %%MatrixMarket matrix coordinate real general
2 9635 2225 141455
3 1 1 1
4 1 11 1
5 1 1199 1
6 1 1242 2
7 1 1833 1
8 1 1913 1
9 1 1948 2
10 1 1962 1
11 1 2033 1
12 1 2134 1
13 1 2203 1
14 1 35 1
15 1 683 1
16 1 693 1
17 10 1166 1
18 10 1184 1
19 ...
```

Figure 18: Sample of output task 1.2

- From the result file, the MTX file size has been subtracted 828 lines.

3.3. Task 1.3 Top 10 Most Frequent Words

- a. **Data Description:** Output mtx file from task 1.2, this is a Matrix Market data file containing termID, docID and frequency. Each column represents the id of the term, the id of the document containing that term and the frequency of the term appearing in that document. The two first rows contain the header and the matrix size.

b. **MapReduce Job Implementation:**

The MapReduce job for this task is structured into 2 phases: the Map phase and the Reduce phase.

- **Map Phase:**



```

1  public static class TokenizerMapper extends Mapper<Object, Text, Text, IntWritable> {
2      public void map(Object key, Text value, Context context) throws IOException, InterruptedException {
3          String[] tokens = value.toString().split("\\s+");
4          if (tokens.length == 3) {
5              String termId = tokens[0];
6              int freq = Integer.parseInt(tokens[2]);
7              context.write(new Text(termId), new IntWritable(freq));
8          }
9      }
10 }

```

Figure 19: Map code

- Input: Each line of the input MTX file represents a term - document pair and its corresponding frequency.
- Process: The Mapper Processes each line by splitting it into tokens to extract the termid and frequency. It then emits a key-value pair where the key is the termid, and the value is the frequency. This allows us to group all frequencies of a term across different documents.
- Output: The output of the Mapper is a list of key-value pairs, where each key is a termid, and each value is the frequency of that term in a specific document.

- Reduce Phase:



```

1  public static class TopWordsReducer extends Reducer<Text, IntWritable, Text, IntWritable> {
2      private Map<String, Integer> termFrequencyMap = new HashMap<>();
3      private static final int TOP_N = 10;
4
5      @Override
6      public void reduce(Text key, Iterable<IntWritable> values, Context context) throws IOException, InterruptedException {
7          int totalFrequency = 0;
8          for (IntWritable value : values) {
9              totalFrequency += value.get();
10         }
11         termFrequencyMap.put(key.toString(), totalFrequency);
12     }
13
14     @Override
15     protected void cleanup(Context context) throws IOException, InterruptedException {
16         // Sắp xếp và ghi kết quả
17         List<Map.Entry<String, Integer>> sortedEntries = new ArrayList<>(termFrequencyMap.entrySet());
18         sortedEntries.sort(Map.Entry.comparingByValue(Comparator.reverseOrder()));
19
20         int count = 0;
21         for (Map.Entry<String, Integer> entry : sortedEntries) {
22             if (count < TOP_N) {
23                 context.write(new Text(entry.getKey()), new IntWritable(entry.getValue()));
24                 count++;
25             } else {
26                 break;
27             }
28         }
29     }

```

Figure 20: Reduce code

- Input: The Reducer receives key-value pairs from the Mapper, where each key is a termid, and the corresponding values are the list of frequencies of that term across all documents.
- Process: The Reducer sums these frequencies to obtain the total frequency for each term across all documents. It stores these totals in a map for further processing.
- Sorting and Selection: After processing all the terms, the Reducer sorts the terms by their total frequency in descending order and selects the top 10 most frequent terms.
- Output: The output of the Reducer is the top 10 most frequent terms, written to the file **task_1_3.txt** in the format [termid] [frequency].

- The description of the map-reduce will be show in the figures below (example finding top 2):

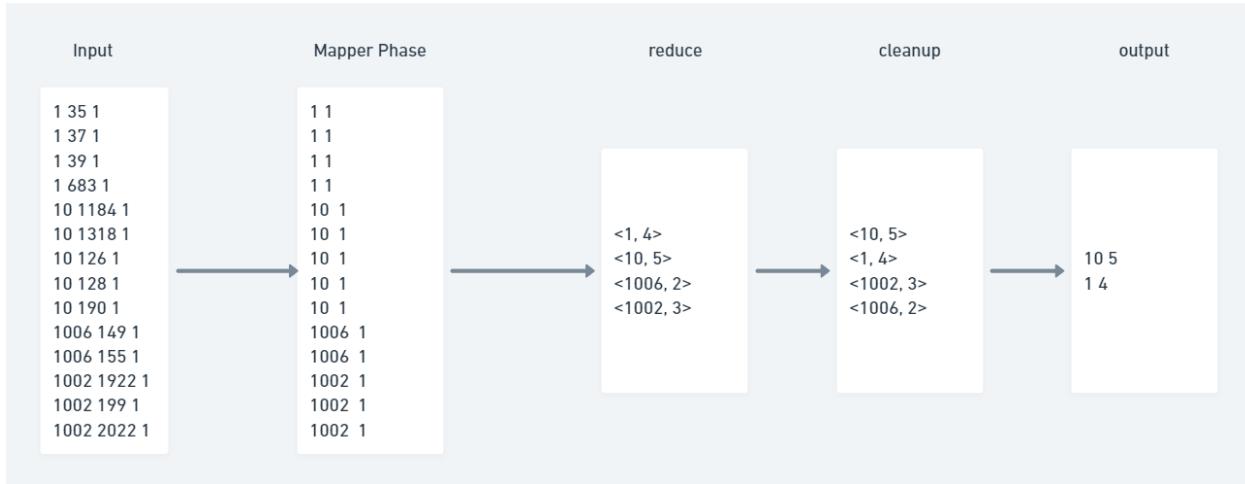


Figure 21: describe map reduce flow

c. Task Flow:

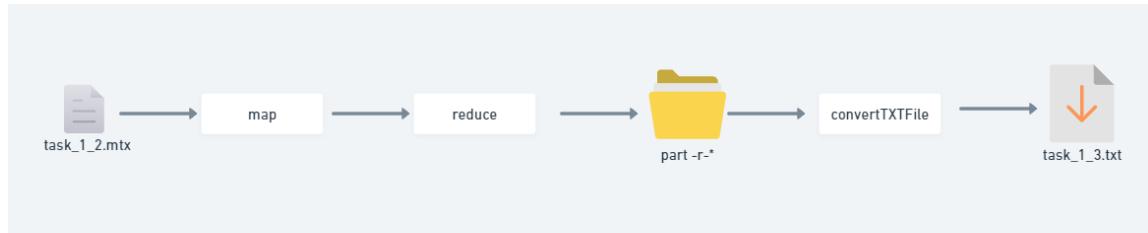
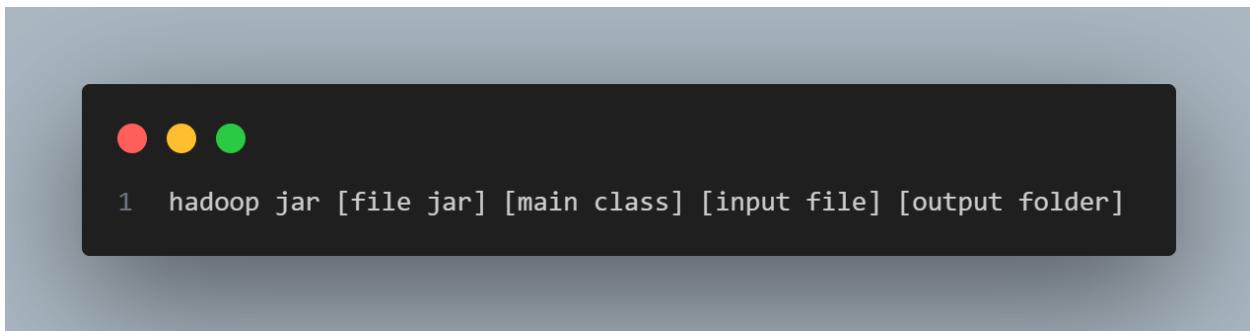


Figure 22: describe task 1.3 flow

- Step 1: Data Preparation
 - The MTX file generated from Task 1.2 is used as input.
- Step 2: Map
 - Each term-document pair is processed to emit term frequency key-value pairs.
- Step 3: Reduce
 - The term frequencies are summed, and the top 10 terms are selected based on their total frequency.
- Step 4: Sorting and Output
 - The terms are sorted in descending order, and the top 10 are written to the output file task_1_3.txt.

d. How to run:

- Command to run the program:

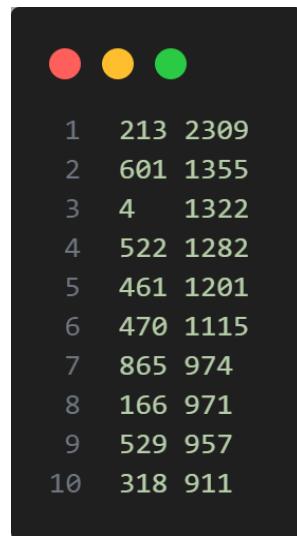


```
hadoop jar [file jar] [main class] [input file] [output folder]
```

Figure 23: Command

e. Result:

- File named task_1_3.txt, which contains the top 10 most frequent terms across all documents. Each line of this file is formatted as [termid] [frequency], where termid is the unique identifier of the term, and frequency is the total number of occurrences of that term across all documents.
- Sample of task_1_3.txt



TermID	Frequency
1	213 2309
2	601 1355
3	4 1322
4	522 1282
5	461 1201
6	470 1115
7	865 974
8	166 971
9	529 957
10	318 911

Figure 24: Output (task_1_3.txt)

3.4. Task 1.4 TF-IDF

- a. Data Description:** Output mtx file from task 1.2, this is a Matrix Market data file containing termID, docID and frequency. Each column represents the id of the

term, the id of the document containing that term and the frequency of the term appearing in that document. The two first rows contain the header and the matrix size.

b. MapReduce Job Implementation:

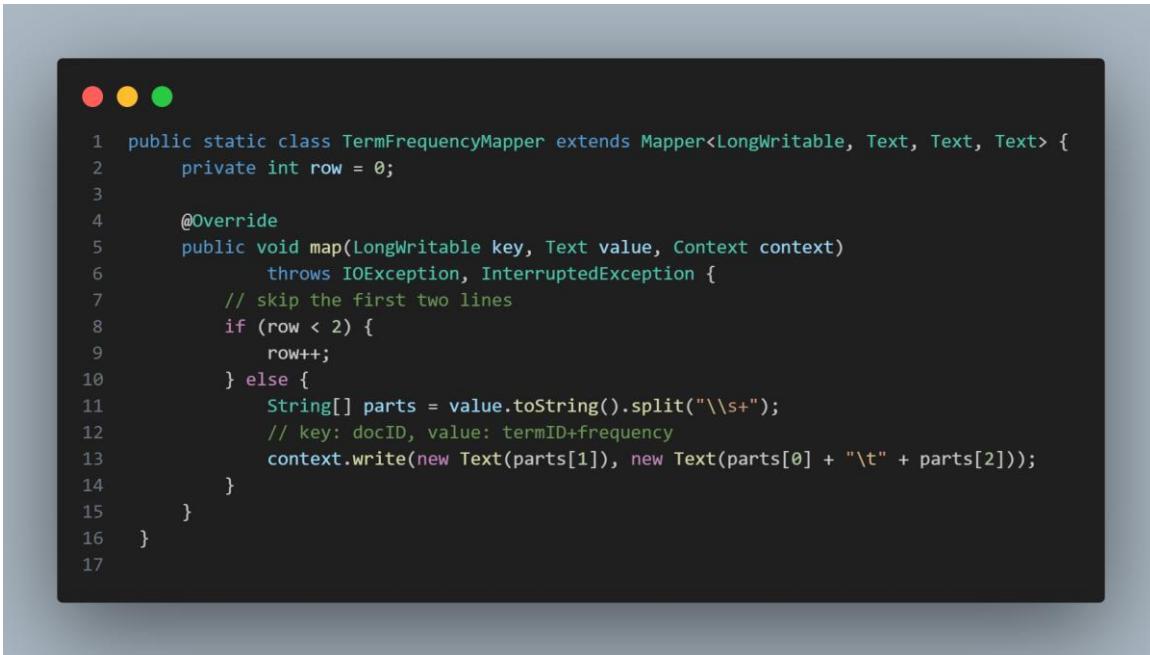
In this task we divide into **2 Map phases** and **2 Reduce phases** for the MapReduce process into **2 distinct Jobs**:

- **Job 1** calculates the Term Frequency (TF) value for each term in each document class.
- **Job 2** calculates and outputs the TF-IDF score for each term in each document class.

Before going into the map reduce work, we will have to save the **total document count value (totalDocs)** in the dataset to use for calculating the TF-IDF score. We used the [**saveToTotalDocumentConfig**](#) function to save the total document value into the configuration. The **totalDocs** value is read from the input file in the second line, where the size values are contained.

1. Job 1: calculates the Term Frequency (TF)

- **Map Phase**
 - Input data: An input data file in data description (after skipping the first 2 lines) containing 3 parts: *termID*, *docID*, and *frequency*.
 - The map() function in TermFrequencyMapper will read each line, split it into 3 parts based on the space, then return the key-value pairs with:
 - **Key:** docID (document index).
 - **Value:** Combine termID and frequency with tab (termID + "\t" + frequency).



```

1  public static class TermFrequencyMapper extends Mapper<LongWritable, Text, Text, Text> {
2      private int row = 0;
3
4      @Override
5      public void map(LongWritable key, Text value, Context context)
6          throws IOException, InterruptedException {
7          // skip the first two lines
8          if (row < 2) {
9              row++;
10         } else {
11             String[] parts = value.toString().split("\\s+");
12             // key: docID, value: termID+frequency
13             context.write(new Text(parts[1]), new Text(parts[0] + "\t" + parts[2]));
14         }
15     }
16 }
17

```

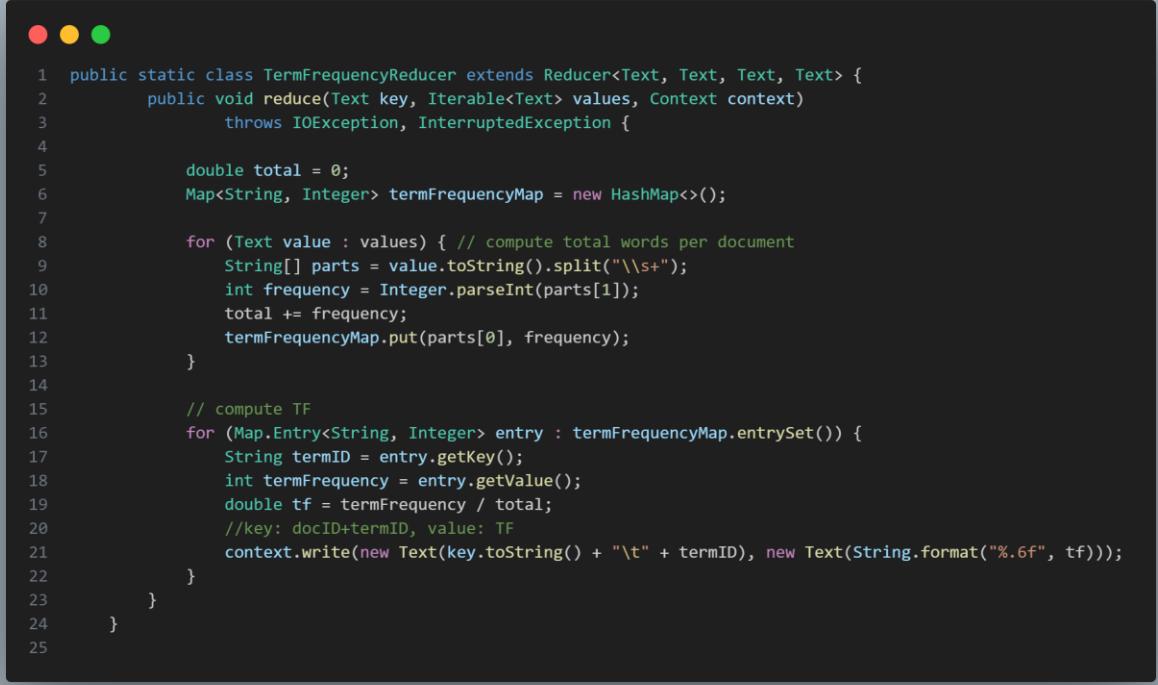
Figure 25: Job 1 Map function

- **Reduce Phase:**

- The purpose of reduce phase is to calculate the term frequency value for each term in each document, according to the following formula:

$$tf(t, d) = \frac{f(t, d)}{\sum\{f(w, d) : (w \in d)\}} = \frac{\text{frequency of term}}{\text{total number of terms in the document}}$$

- Input data: The output of Mapper with key as docID and value as termID combined with frequency.
- First, the reducer extracts the termID and frequency from each value and calculates the total number of terms in that document by summing up all the frequency values.
- After determining the total number of terms in the document, the reducer calculates the TF for each termID in that document by dividing the frequency of the term by the total number of terms in the document.
- The output will be key-value pairs where
 - **Key:** docID + “\t” + termID
 - **Value:** the computed TF value.
 - The output form: [docID] [termID] [TF]



```

1  public static class TermFrequencyReducer extends Reducer<Text, Text, Text, Text> {
2      public void reduce(Text key, Iterable<Text> values, Context context)
3          throws IOException, InterruptedException {
4
5          double total = 0;
6          Map<String, Integer> termFrequencyMap = new HashMap<>();
7
8          for (Text value : values) { // compute total words per document
9              String[] parts = value.toString().split("\\s+");
10             int frequency = Integer.parseInt(parts[1]);
11             total += frequency;
12             termFrequencyMap.put(parts[0], frequency);
13         }
14
15         // compute TF
16         for (Map.Entry<String, Integer> entry : termFrequencyMap.entrySet()) {
17             String termID = entry.getKey();
18             int termFrequency = entry.getValue();
19             double tf = termFrequency / total;
20             //key: docID+termID, value: TF
21             context.write(new Text(key.toString() + "\t" + termID), new Text(String.format("%.6f", tf)));
22         }
23     }
24 }
25

```

Figure 26: Job 1 Reduce function

2. Job 2

- Map Phase:

- The input for the map phase is the output from the reduce phase of Job 1, where each line consists of a docID, termID, and TF value.
- The mapper's job is to process this data and split each line into key-value pairs where:
 - **Key:** termID
 - **Value:** docID+ “\t” + TF value



```

1  public static class TFIDFMapper extends Mapper<LongWritable, Text, Text, Text> {
2      @Override
3      public void map(LongWritable key, Text value, Context context)
4          throws IOException, InterruptedException {
5          String[] parts = value.toString().split("\\s+");
6
7          // key: termID, value: docID+TF
8          context.write(new Text(parts[0]), new Text(parts[1] + "\\t" + parts[2]));
9      }
10 }

```

Figure 27: Job 2 Map function

- **Reduce Phase:**

- The purpose of the Reduce phase in Job 2 is to calculate the TF-IDF score for each term in each document.
- Input data: the output from the Map phase above with the key: termID and the value: docID and TF.
- Before going into the Reduce phase, we will need to read the value of the total number of documents that I introduced at the beginning, this value is saved in the configuration and is used to calculate the IDF score below.
- First, we will calculate Inverse Document Frequency (IDF) for each term in the document. It is calculated as the logarithm of the ratio of the total number of documents to the number of documents containing the term.

$$\begin{aligned}
 idf(t, D) &= \log\left(\frac{|D|}{|\{d \in D : t \in d\}|}\right) \\
 &= \log\left(\frac{\text{total number of documents}}{\text{number of document contains the term}}\right)
 \end{aligned}$$

- Then, for each pair of docID and TF in the list, we will calculate the TF-IDF score value by multiplying TF by IDF:

$$tfidf(t, d, D) = tf(t, d) \times idf(t, D)$$

- The final result will be returned as key value pairs in the following format:

- **Key:** termID + “\\t” + docID

- **Value:** TF-IDF score
- The output form: [termID] [docID] [TFIDF]

```

1  public static class TFIDFReducer extends Reducer<Text, Text, Text, Text> {
2      private int totalDocs;
3
4      @Override
5      protected void setup(Context context) throws IOException, InterruptedException {
6          totalDocs = context.getConfiguration().getInt("totalDocs", 1); // get totalDocs
7      }
8
9      @Override
10     public void reduce(Text key, Iterable<Text> values, Context context)
11         throws IOException, InterruptedException {
12         double documentCount = 0;
13         List<String> docTFList = new ArrayList<>();
14
15         for (Text value : values) { // get the document count: number of document contains term
16             documentCount += 1;
17             docTFList.add(value.toString());
18         }
19
20         double idf = Math.log(totalDocs / documentCount); // compute IDF
21
22         for (String docTF : docTFList) {
23             String[] parts = docTF.split("\\s+");
24
25             String docId = parts[0];
26             double tf = Double.parseDouble(parts[1]);
27
28             double tfidf = tf * idf; // compute TFIDF
29
30             // key: term ID+docID, value: TFIDF score
31             context.write(new Text(key.toString() + "\t" + docId), new Text(String.format("%.6f", tfidf)));
32         }
33     }
34 }

```

Figure 28: Job 2 Reduce Function

- The description of the map-reduce will be show in the figures below (illustrative data):

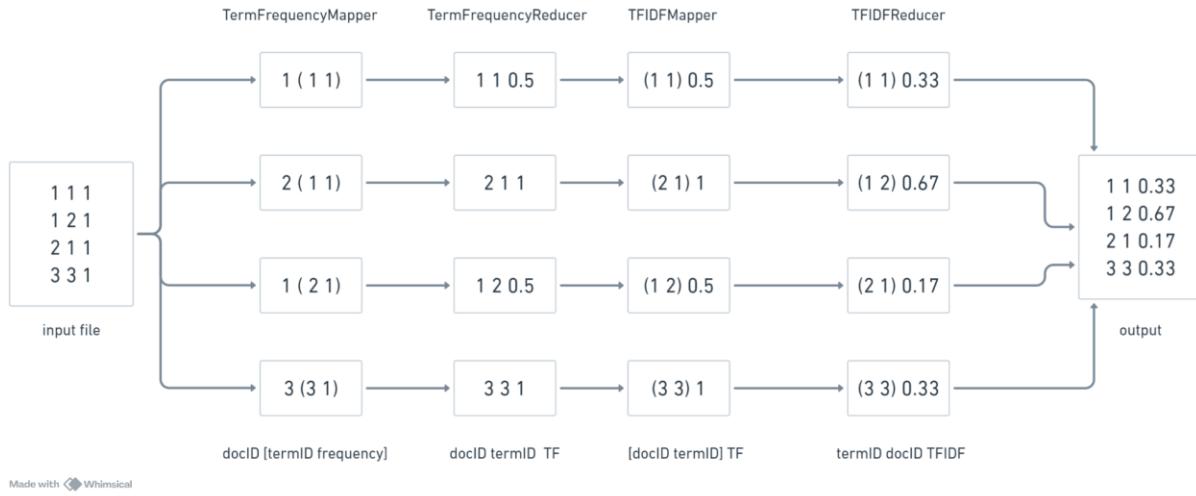
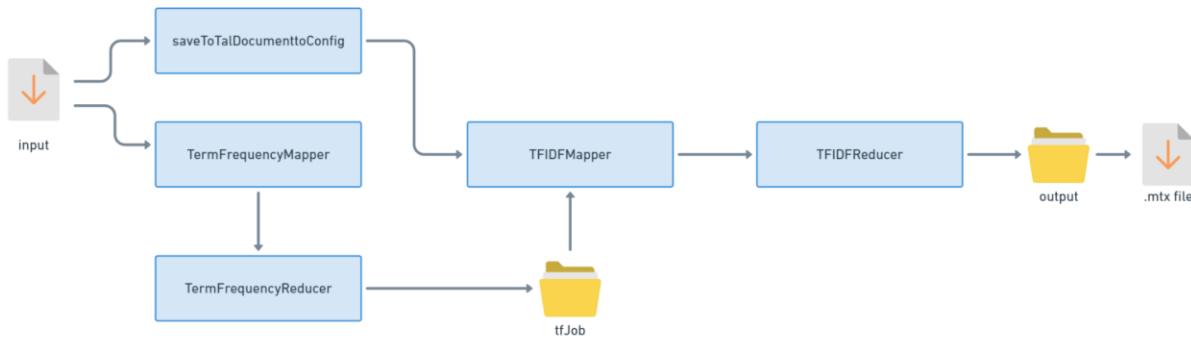


Figure 29: Description of the map-reduce flow

c. Task Flow:

- Step 1: Read data file to save total document quantity value into config via ***saveTotalDocumenttoConfig*** function.
- Step 2: Perform Job 1 to calculate the **TF** for each term
 - TermFrequencyMapper:
 - Read input data file and return key-value pairs corresponding to [docID] and [termID frequency]
 - TermFrequencyReducer:
 - Calculate TF value for each term in each document, return key-value pairs corresponding to [docID termID] and [TF]
- Step 3: Save the output from Job 1 to be the input for Job 2
- Step 4: Perform Job 2 to calculate the **TF-IDF** score for each term
 - TFIDFMapper:
 - Read values from input file which is output of job 1, return key value pairs corresponding to [docID termID] and [TF]
 - TFIDFReducer:
 - Using the previously set documentCount and TotalDocs, calculates the IDF value.

- Compute the TF-IDF score and generate the key value pairs corresponding to [termID docID] [TF]
- Step 5: Save the final output and convert it to an MTX file.



Made with Whimsical

Figure 30: Overall task flow of task 1.4

d. How to run:

- **Configure this section:** We have to define the output file path for Job 1 - TF Job, this will also be the input path for the Mapper in Job 2.

```

 1 FileOutputFormat.setOutputPath(tfJob, new Path("/tfJob"));
  
```

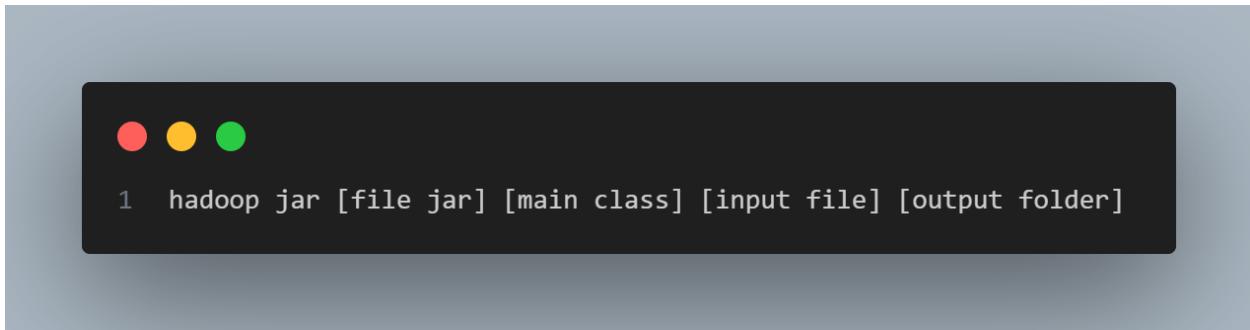
Figure 31: Config output path for Job 1 (Term Frequency Job)

```

 1 FileInputFormat.addInputPath(tfidfJob, new Path("/tfJob"));
  
```

Figure 32: Config input path for Job 2 (TFIDF Job)

- **Command to run the program:**

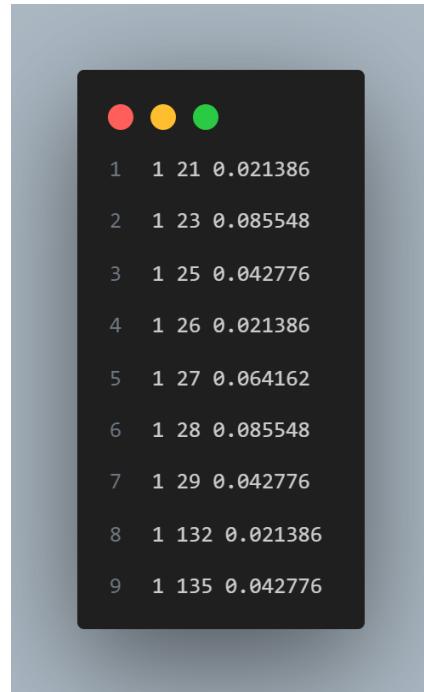


```
● ● ●  
1 hadoop jar [file jar] [main class] [input file] [output folder]
```

Figure 33: Command

e. Result:

- Sample of file task_1_4.mtx



```
● ● ●  
1 1 21 0.021386  
2 1 23 0.085548  
3 1 25 0.042776  
4 1 26 0.021386  
5 1 27 0.064162  
6 1 28 0.085548  
7 1 29 0.042776  
8 1 132 0.021386  
9 1 135 0.042776
```

Figure 34: Sample output of task 1.4

3.5. Task 1.5 Highest average TF-IDF

- **Formular:**

$$\text{avg}_{tfidf_{t,C_i,D}} = \frac{1}{|C_i|} \sum_{d \in C_i} \text{tdidf}_{t,d,D}$$

where: $|C_i|$ is total of files in class C_i

- To complete Task 1.5, we divided the MapReduce process into two distinct jobs:
 - **Job 1** calculates the average TF-IDF value for each term in each document class.
 - **Job 2** then identifies and outputs the top 5 terms with the highest average TF-IDF values for each class

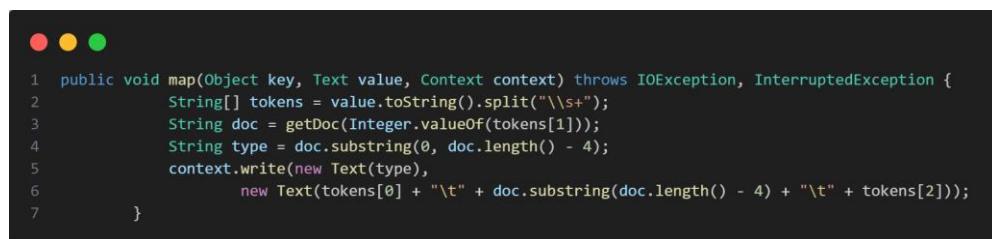
a. Data Description:

- **Job 1:** Calculate Average TF-IDF for Each Term in Each Class
 - The input for Job 1 is an MTX file from Task 1.4. This file follows the Matrix Market format and contains:
 - *termID*: Identifier for the term (word).
 - *docID*: Identifier for the document containing the term.
 - *TF-IDF*: TF-IDF value of the term in the document.
- **Job 2:** Identify Top 5 Terms with Highest Average TF-IDF for Each Class
 - The input for Job 2 is the output from Job 1, which contains termID, class, and average TF-IDF values.

b. MapReduce Job Implementation:

1. Job 1: Calculate Average TF-IDF for Each Term in Each Class

- **Map Phase:**



```

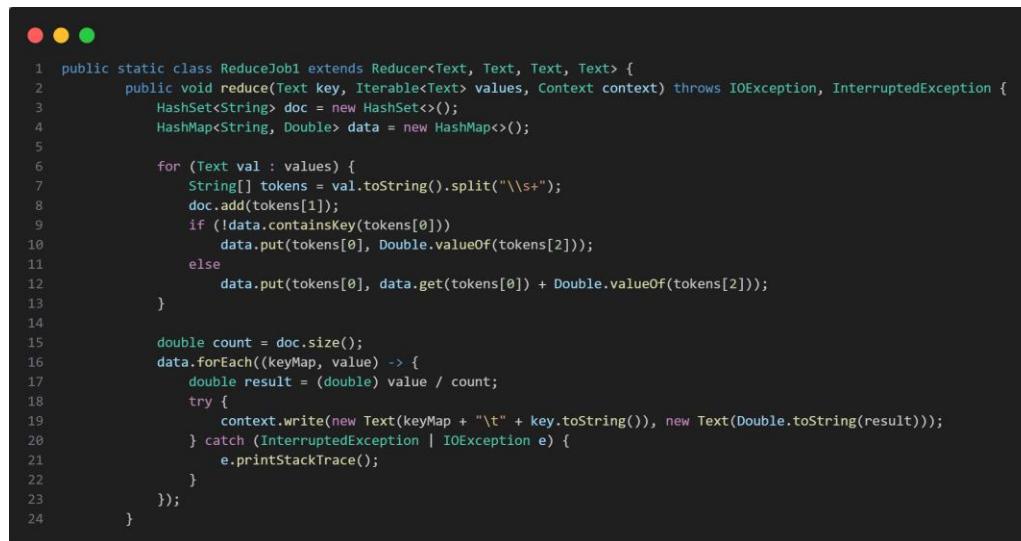
1 public void map(Object key, Text value, Context context) throws IOException, InterruptedException {
2     String[] tokens = value.toString().split("\\s+");
3     String doc = getDoc(Integer.valueOf(tokens[1]));
4     String type = doc.substring(0, doc.length() - 4);
5     context.write(new Text(type),
6                 new Text(tokens[0] + "\t" + doc.substring(doc.length() - 4) + "\t" + tokens[2]));
7 }

```

Figure 35: Map code (job 1)

- Input: Each line of the MTX file represents a term-document pair with a TF-IDF value.

- Processing:
 - Step 1: Read the **bbc.docs** file: Read the list of documents from the **bbc.docs** file and store them in a list for reference.
 - Step 2: Extract classID from docID: For each docID, split it to obtain the first part as classID.
 - Step 3: Emit key-value pairs:
 - **Key:** docID (class name corresponding to the docID).
 - **Value:** termID + "\t" + docID + "\t" + TF-IDF.
 - Output: A list of key-value pairs where the key is classID and the value is a combination of termID, docID, and the TF-IDF value.
- Reduce Phase:



```

1  public static class ReduceJob1 extends Reducer<Text, Text, Text, Text> {
2      public void reduce(Text key, Iterable<Text> values, Context context) throws IOException, InterruptedException {
3          HashSet<String> doc = new HashSet<>();
4          HashMap<String, Double> data = new HashMap<>();
5
6          for (Text val : values) {
7              String[] tokens = val.toString().split("\\s+");
8              doc.add(tokens[1]);
9              if (!data.containsKey(tokens[0]))
10                  data.put(tokens[0], Double.valueOf(tokens[2]));
11              else
12                  data.put(tokens[0], data.get(tokens[0]) + Double.valueOf(tokens[2]));
13          }
14
15          double count = doc.size();
16          data.forEach((keyMap, value) -> {
17              double result = (double) value / count;
18              try {
19                  context.write(new Text(keyMap + "\t" + key.toString()), new Text(Double.toString(result)));
20              } catch (InterruptedException | IOException e) {
21                  e.printStackTrace();
22              }
23          });
24      }
}

```

Figure 36: Reduce code (job 1)

- Input: The Reducer receives key-value pairs with classID as the key and termID + docID + TF-IDF as the value.
- Processing:
 - Step 1: Create a HashSet to store the documents: A HashSet helps count the number of documents in each class (classID).
 - Step 2: Use a HashMap to store cumulative TF-IDF values: The HashMap stores the total TF-IDF for each termID within a class.

- Step 3: Calculate the average TF-IDF: For each *termID* in a class, calculate the average TF-IDF by dividing the total TF-IDF by the number of documents in the class.
- Step 4: Emit the average TF-IDF with termID and classID as the key.

Key: termID + "\t" + classID.

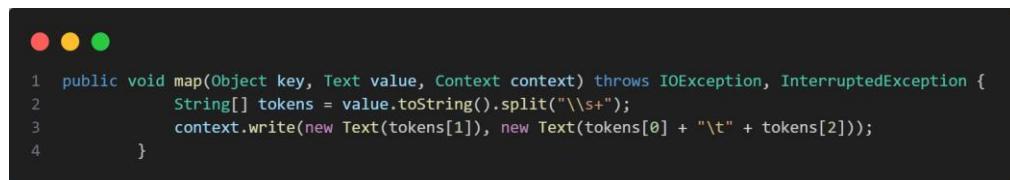
Value: The average TF-IDF value.

- Output: Key-value pairs with:

- Key: *TermID* + *classID*
- Value: The average TF-IDF of the termID in the classID.

2. Job 2: Identify Top 5 Terms with Highest Average TF-IDF for Each Class

- Map Phase:



```

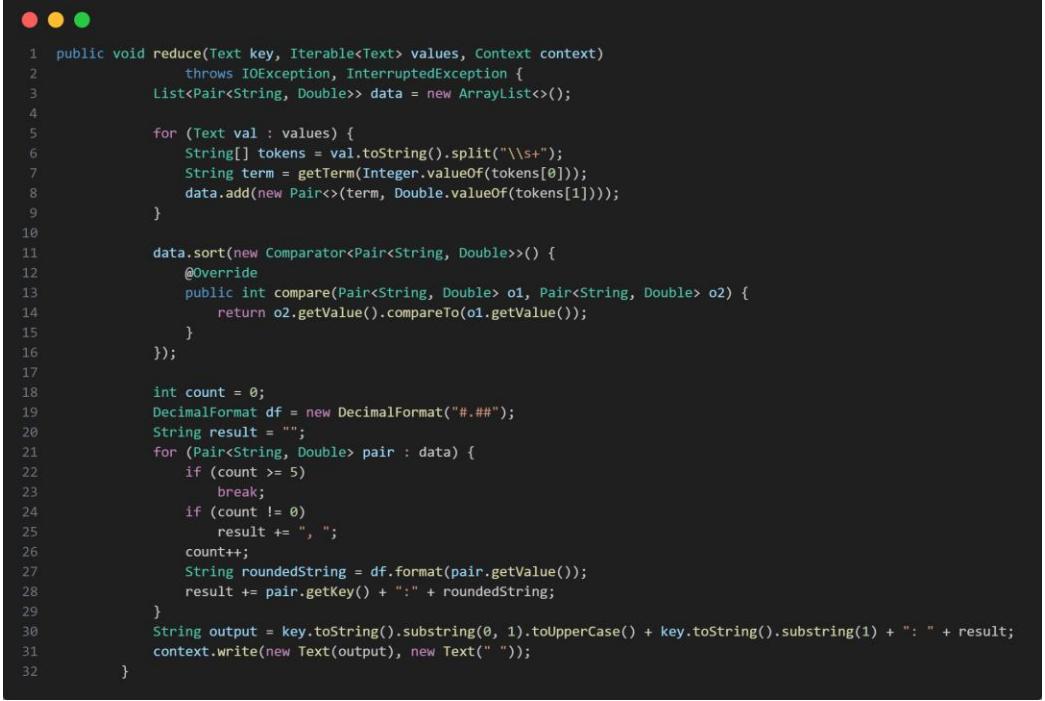
1  public void map(Object key, Text value, Context context) throws IOException, InterruptedException {
2      String[] tokens = value.toString().split("\\s+");
3      context.write(new Text(tokens[1]), new Text(tokens[0] + "\t" + tokens[2]));
4  }

```

Figure 37: Map code (job 2)

- Input: Each key-value pair from Job 1 with *termID* and *average TF-IDF value*.
- Processing: Create key-value pairs:
 - Key: classID (class name).
 - Value: termID + "\t" + average TF-IDF value.
- Output: Key-value pairs where the key is classID, and the value is termID with the average TF-IDF value.

- Reduce Phase:



```

1  public void reduce(Text key, Iterable<Text> values, Context context)
2      throws IOException, InterruptedException {
3      List<Pair<String, Double>> data = new ArrayList<>();
4
5      for (Text val : values) {
6          String[] tokens = val.toString().split("\\s+");
7          String term = getTerm(Integer.valueOf(tokens[0]));
8          data.add(new Pair<>(term, Double.valueOf(tokens[1])));
9      }
10
11     data.sort(new Comparator<Pair<String, Double>>() {
12         @Override
13         public int compare(Pair<String, Double> o1, Pair<String, Double> o2) {
14             return o2.getValue().compareTo(o1.getValue());
15         }
16     });
17
18     int count = 0;
19     DecimalFormat df = new DecimalFormat("#.##");
20     String result = "";
21     for (Pair<String, Double> pair : data) {
22         if (count >= 5)
23             break;
24         if (count != 0)
25             result += ", ";
26         count++;
27         String roundedString = df.format(pair.getValue());
28         result += pair.getKey() + ":" + roundedString;
29     }
30     String output = key.toString().substring(0, 1).toUpperCase() + key.toString().substring(1) + ":" + result;
31     context.write(new Text(output), new Text(" "));
32 }

```

Figure 38: Reduce code (job 2)

- Input: The Reducer receives key-value pairs where the key is classID, and the values are termID + average TF-IDF.
- Processing:
 - Step 1: Create a Pair class to map termID to actual term names: Use the bbc.terms file to map termID to the actual term name.
 - Step 2: Sort terms by their average TF-IDF values: Sort the terms by their average TF-IDF values in descending order.
 - Step 3: Select the top 5 terms with the highest average TF-IDF: Choose the top 5 terms with the highest average TF-IDF values for each class.
 - Step 4: Emit the top 5 terms with their average TF-IDF for each class.
 - **Key:** The name of the class (e.g., “Business,” “Tech”).
 - **Value:** A string containing the top 5 terms formatted as: term:average TF-IDF, with terms sorted by their average TF-IDF values in descending order.

- Output:
 - Key: The name of the class.
 - Value: A string containing the top 5 terms and their average TF-IDF for each class.
- The description of the map reduce will be show in the figures below (illustrative data):

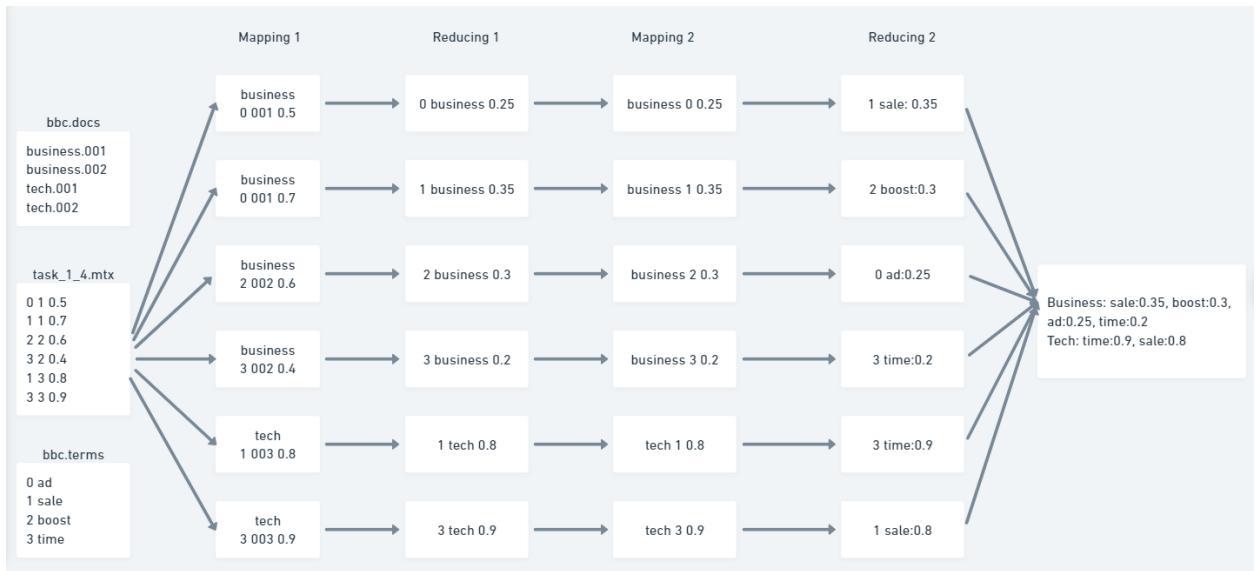


Figure 39: Description of the map reduce flow

c. Task Flow:



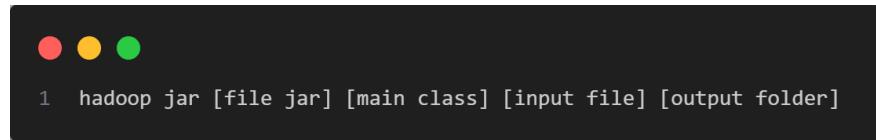
Figure 40: Description of task 1.5 flow

- **Job 1:** Calculate Average TF-IDF for Each Term in Each Class
 - **Data Preparation:** Use the MTX file from Task 1.4.

- **Map Phase:** Process term-document pairs to emit termID and TF-IDF values with docID as the key.
- **Reduce Phase:** Compute average TF-IDF values for each termID within each class and emit results.
- **Job 2:** Identify Top 5 Terms with Highest Average TF-IDF for Each Class
 - **Data Preparation:** Use the output from Job 1.
 - **Map Phase:** Emit class and termID with average TF-IDF values.
 - **Reduce Phase:** Sort terms by average TF-IDF, select the top 5 terms for each class, and emit the final results.

d. Result:

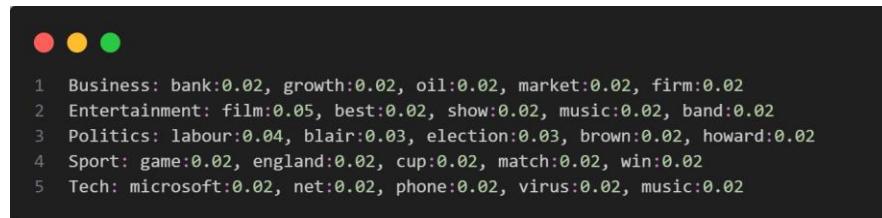
- **Command to Run the Program:**



```
1 hadoop jar [file jar] [main class] [input file] [output folder]
```

Figure 41: command

- The final output is written to *task_1_5.txt*. This file contains the top 5 terms with the highest average TF-IDF values for each class.
- Each line in the file is formatted as:
 - **Class: term: average IF - IDF value**
- After completing the program, the results will return to the *task_1_3.txt* file. The resule of the *task_1_3.txt* file are as follows:



```
1 Business: bank:0.02, growth:0.02, oil:0.02, market:0.02, firm:0.02
2 Entertainment: film:0.05, best:0.02, show:0.02, music:0.02, band:0.02
3 Politics: labour:0.04, blair:0.03, election:0.03, brown:0.02, howard:0.02
4 Sport: game:0.02, england:0.02, cup:0.02, match:0.02, win:0.02
5 Tech: microsoft:0.02, net:0.02, phone:0.02, virus:0.02, music:0.02
```

Figure 42: Output (*task_1_5.txt*)

4. Part 2: K-Means Algorithm

4.1. Task 2.1 K-Means on 2D Data

- a. **Data Description:** A sample file containing 2D data points (one point per line, separated by spaces) with cluster, datapoint x, datapoint y.
- b. **MapReduce Job Implementation:**

First, we will randomly initialize k centroids based on selecting k random points from the given dataset. This will be used on the *initializeCentroids* function.

- This function will read the data file, then randomly select k data points from the file to be centroids.



```

1  public static List<Point> initializeCentroids(int k, FileSystem fs, String inputPath) throws IOException {
2      List<Point> allPoints = new ArrayList<>();
3      List<Point> centroids = new ArrayList<>();
4
5      FileStatus[] files = fs.listStatus(new Path(inputPath));
6      for (FileStatus file : files) {
7          try (BufferedReader br = new BufferedReader(new InputStreamReader(fs.open(file.getPath())))) {
8              String line;
9              while ((line = br.readLine()) != null) {
10                  try {
11                      String[] parts = line.trim().split(",");
12                      double x = Double.parseDouble(parts[1]);
13                      double y = Double.parseDouble(parts[2]);
14                      allPoints.add(new Point(x, y));
15                  } catch (NumberFormatException e) {
16                      continue;
17                  }
18              }
19          }
20      }
21
22      // shuffle the list to get random centroids
23      Collections.shuffle(allPoints);
24
25      for (int i = 0; i < k; i++) {
26          centroids.add(allPoints.get(i));
27      }
28
29      return centroids;
30  }

```

Figure 43: *initializeCentroids* function

We have added a parameter called **THRESHOLD** with the value of **0.00001** to set the **early stopping** mechanism for the algorithm. If the distance between the newly generated centroids and the old centroids is not greater than this threshold, the algorithm will stop.

In this task, we have **1 Map phase** and **1 Reduce phase**. Specifically, the execution steps are as follows:

- **In Map Phase:**

- Before the mapping phase begins, the Mapper loads the centroids from the previous iteration (or from the initial centroids if it's the first iteration) using the *setup* method. These centroids are stored in the Hadoop configuration via *saveCentroidsToConfig* and *readCentroidsFromConfig* functions (check it via our source code).
- The input data for the Mapper consists of data points in the format “cluster, datapoint_x, datapoint_y”. We will parse this string to extract the x and y coordinates of each data point, which are essential to determine the location of each point on the 2D plane. We intentionally skip over any lines that start with the letter "c" (these lines typically store centroids in the format "centroids:cluster, centroid_x, centroid_y", we will explain in detail in the Reduce Phase).
- Then, we will find the nearest cluster for the datapoint by using the *getNearestCluster* function.
- Finally, it will generates pairs of key-value corresponding to the following format:
 - **Key:** the nearest cluster of the datapoint
 - **Value:** datapoint (x,y)
 - The output form: “[cluster],[datapoint_x],[datapoint_y]”.

```

1 public static class ClusterMapper extends Mapper<Object, Text, IntWritable, Text> {
2     private List<Point> centroids;
3     private IntWritable clusterKey = new IntWritable();
4
5     @Override
6     protected void setup(Context context) {
7         centroids = readCentroidsFromConfig(context.getConfiguration());
8     }
9
10    @Override
11    protected void map(Object key, Text value, Context context) throws IOException, InterruptedException {
12        String[] parts;
13        Point point;
14        try {
15            parts = value.toString().trim().split(",");
16            double x, y;
17            x = Double.parseDouble(parts[1]);
18            y = Double.parseDouble(parts[2]);
19            if (parts[0].charAt(0) == 'c') { // skip if it is centroids
20                return;
21            }
22            point = new Point(x, y);
23        } catch (Exception e) {
24            return;
25        }
26
27        int nearestCluster = getNearestCluster(point);
28        clusterKey.set(nearestCluster);
29
30        // key: cluster, value: datapoint
31        context.write(clusterKey, new Text(point.toString()));
32    }
33
34    private int getNearestCluster(Point point) {
35        int nearestCluster = 0;
36        double minDistance = Double.MAX_VALUE;
37
38        for (int i = 0; i < centroids.size(); i++) {
39            double distance = point.distanceTo(centroids.get(i));
40            if (distance < minDistance) {
41                minDistance = distance;
42                nearestCluster = i;
43            }
44        }
45
46        return nearestCluster;
47    }
48 }

```

Figure 44: Map function

- **In Reduce Phase:**

- After the "shuffle and sort" stage, Reducer receives as input the cluster index (key) and the list of data points belonging to that cluster (values). Then it will convert the input text datapoint values into lists of Points.

- Next, it will calculate to create new centroids values for the cluster based on averaging the points in the cluster.
- Because we need to save the datapoint values for future iterations, we will save the data points and the cluster they belong to as key-value pairs of the form:
 - **Key:** cluster
 - **Value:** datapoint (x,y)
 - The output form: “[cluster],[datapoint_x],[datapoint_y]”.
- Finally we will save the newly calculated centroids values, however to mark and distinguish them from the data points saved above we will save them as: "centroids:"[cluster],[centroid_x],[centroid_y]. That's why there are data lines starting with the char “c” that we need to ignore in the Map phase above. Therefore, the Reducer output key-value pairs will be of the form:
 - **Key:** “centroids:”cluster
 - **Value:** centroid (x,y)
 - The output form: "centroids:[cluster],[centroid_x],[centroid_y]"

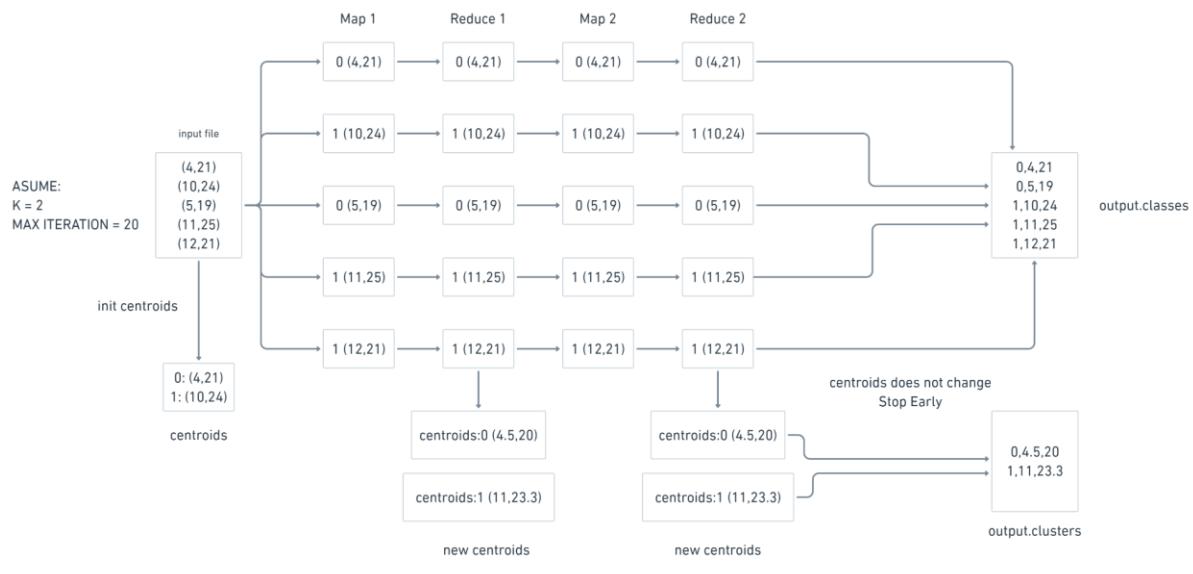
```

1  public static class ClusterReducer extends Reducer<IntWritable, Text, Text, Text> {
2      @Override
3      protected void reduce(IntWritable key, Iterable<Text> values, Context context)
4          throws IOException, InterruptedException {
5          List<Point> points = new ArrayList<>();
6          for (Text value : values) {
7              points.add(Point.fromString(value.toString()));
8          }
9
10         Point newCentroid = Point.average(points);
11
12         for (Point point : points) {
13             // save the cluster & datapoints: key: cluster, value: datapoint
14             context.write(new Text(key.toString() + ","), new Text(point.toString()));
15         }
16     }
17
18     // save the centroids, key: "centroids:"cluster, value: centroid
19     context.write(new Text("centroids:" + key.toString() + ","), new Text(newCentroid.toString()));
20   }
21 }

```

Figure 45: Reduce Function

- The description of the map reduce will be show in the figures below (illustrative data):



Made with Whimsical

Figure 46: The description of the map reduce

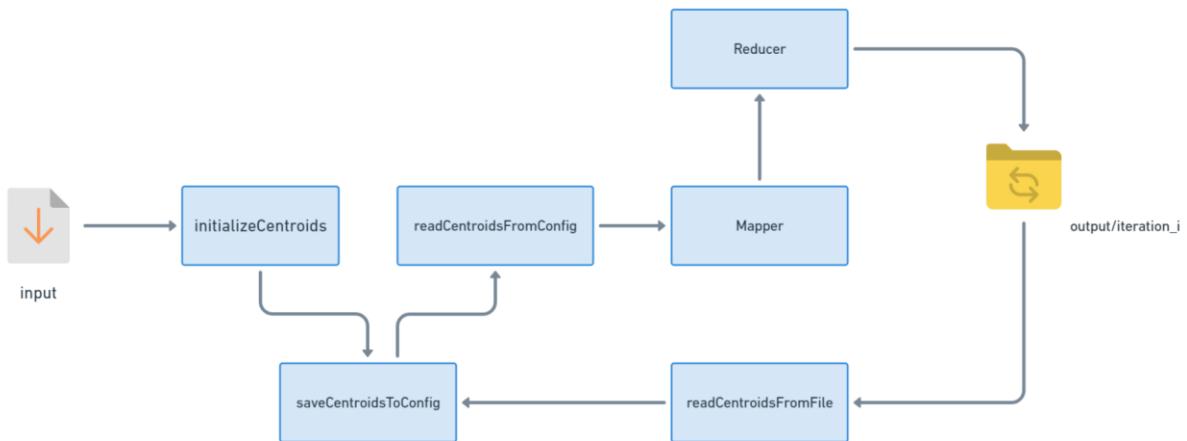
c. Task Flow:

- Step 1: Init randomly k centroids by *initializeCentroids*. Save it to the configuration by *saveCentroidsToConfig* function.
- Step 2: With K = 3 and Iterations = 20, for each iteration:
 - Start Map Phase:
 - Load the centroids from configuration via *readCentroidsFromConfig*.
 - Find the nearest cluster for the datapoint.
 - Return the key-value pairs of cluster and datapoint.
 - Start Reduce Phase
 - Load the output of Map Phase.
 - Compute new centroids by averaging all data points in the cluster.
 - Return the key-value pairs of cluster and centroids and key-value pairs of cluster and data points.
 - Save the Reduce Phase output.
 - Check if the distance between new centroids and old centroids does not larger than THRESHOLD -> stop early
 - Read the output Reduce Phase from HDFS to get the new centroid via *readCentroidsFromFile* function and then save it to configuration using *saveCentroidsToConfig*.
- Step 3: Save the last iteration output into **task_2_1.classes** and **task_2_1.cluster**.



```
1 public static void runKmeans(Configuration conf, String inputPath, String outputPath) throws Exception {
2     FileSystem fs = FileSystem.get(conf);
3
4     // init randomly K centroids
5     List<Point> centroids = initializeCentroids(K, fs, inputPath);
6     saveCentroidsToConfig(conf, centroids);
7     Path centroidOutputPath = new Path(outputPath);
8
9     for (int iteration = 0; iteration < MAX_ITERATIONS; iteration++) {
10         Job job = Job.getInstance(conf, "KMeans Iteration " + iteration);
11         job.setJarByClass(task2_1.class);
12
13         job.setMapperClass(ClusterMapper.class);
14         job.setReducerClass(ClusterReducer.class);
15
16         job.setMapOutputKeyClass(IntWritable.class);
17         job.setMapOutputValueClass(Text.class);
18         job.setOutputKeyClass(Text.class);
19         job.setOutputValueClass(Text.class);
20
21         FileInputFormat.addInputPath(job, new Path(inputPath));
22         FileOutputFormat.setOutputPath(job, new Path(outputPath + "/iter_" + iteration));
23
24         job.waitForCompletion(true);
25
26         // update centroids
27         centroidOutputPath = new Path(outputPath + "/iter_" + iteration + "/part-r-00000");
28         List<Point> newCentroids = readCentroidsFromFile(fs, centroidOutputPath);
29
30         // check if can stop early
31         double maxDistance = 0;
32         for (int j = 0; j < centroids.size(); j++) {
33             double distance = centroids.get(j).distanceTo(newCentroids.get(j));
34             if (distance >= maxDistance) {
35                 maxDistance = distance;
36             }
37         }
38
39         if (maxDistance <= THRESHOLD) {
40             break;
41         }
42
43         centroids = newCentroids;
44         saveCentroidsToConfig(conf, centroids);
45     }
46
47     // save the output into task_2_1.cluster & task_2_1.classes
48     writeOutput(fs, K, centroidOutputPath, outputPath + "/task_2_1.clusters", outputPath + "/task_2_1.classes");
49 }
```

Figure 47: Function to set up flow for run KMeans program



The overall job architecture for task 2.1

Made with Whimsical

Figure 48: Overall task flow of task 2.1

d. How to run:

- Command to run the program:

```

● ● ●
1 hadoop jar [file jar] [main class] [input file] [output folder]
  
```

Figure 49: Command

e. Result:

- File task_2_1.cluster



Figure 50: File task_2_1.clusters

- Part of file task_1_2.classes



Figure 51: Sample of task_2_1.classes

4.2. Task 2.2 K-Means on Preprocessed Data

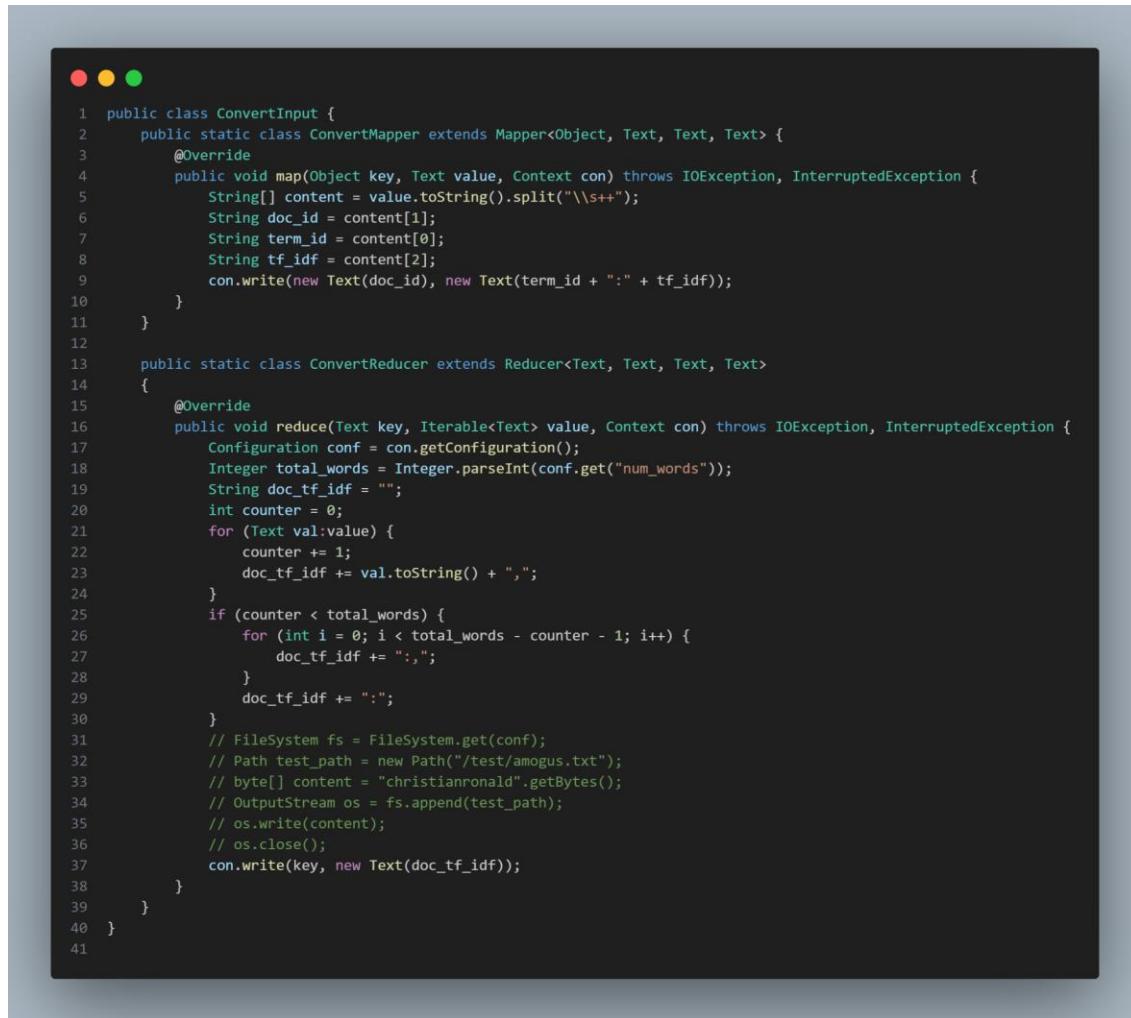
a. Data Description:

- The input data from Task 1.4 consists of a .mtx file containing tf-idf scores for words, which have been filtered by frequency and organized line by line. Due to the complexity of applying MapReduce logic directly to this format, a reformatting stage is necessary to convert the file into a format more suitable for efficient MapReduce processing.
- This reformatting stage involves two sequential MapReduce jobs:

- **Counting Unique Words:** The first MapReduce job processes each line of the input file from Task 1.4 to count the number of unique words. This step is essential as it determines the dimensionality of the vector representing each document in the dataset. The map phase operates similarly to the standard word count example, using the term_id as the key. During the reduce phase, identical term_id values are combined and a global counter variable is incremented. This counter represents the vector dimension, which will be used in subsequent calculations.
- **Conversion to the Desired Format:** The second MapReduce job focuses on converting the input file into a specified format as outlined in the project guidelines. During the map phase, the document_id is used as the key, and the term_id pair as the value. The reduce phase then aggregates all term_id pairs into a single string, with each pair separated by a comma. The structure of the output file will adhere to the following format:

doc_id term_id_1:tf_idf_1,term_id_2:tf_idf_2,...

- The outputs of these two MapReduce jobs will be stored in a /tmp subdirectory within the output directory. The structure and details of this directory will be elaborated upon in the following section of the report.



```

1 public class ConvertInput {
2     public static class ConvertMapper extends Mapper<Object, Text, Text, Text> {
3         @Override
4         public void map(Object key, Text value, Context con) throws IOException, InterruptedException {
5             String[] content = value.toString().split("\\s+");
6             String doc_id = content[1];
7             String term_id = content[0];
8             String tf_idf = content[2];
9             con.write(new Text(doc_id), new Text(term_id + ":" + tf_idf));
10        }
11    }
12
13    public static class ConvertReducer extends Reducer<Text, Text, Text, Text>
14    {
15        @Override
16        public void reduce(Text key, Iterable<Text> value, Context con) throws IOException, InterruptedException {
17            Configuration conf = con.getConfiguration();
18            Integer total_words = Integer.parseInt(conf.get("num_words"));
19            String doc_tf_idf = "";
20            int counter = 0;
21            for (Text val:value) {
22                counter += 1;
23                doc_tf_idf += val.toString() + ",";
24            }
25            if (counter < total_words) {
26                for (int i = 0; i < total_words - counter - 1; i++) {
27                    doc_tf_idf += ":";;
28                }
29                doc_tf_idf += ":";;
30            }
31            // FileSystem fs = FileSystem.get(conf);
32            // Path test_path = new Path("/test/amogus.txt");
33            // byte[] content = "christianronald".getBytes();
34            // OutputStream os = fs.append(test_path);
35            // os.write(content);
36            // os.close();
37            con.write(key, new Text(doc_tf_idf));
38        }
39    }
40 }
41

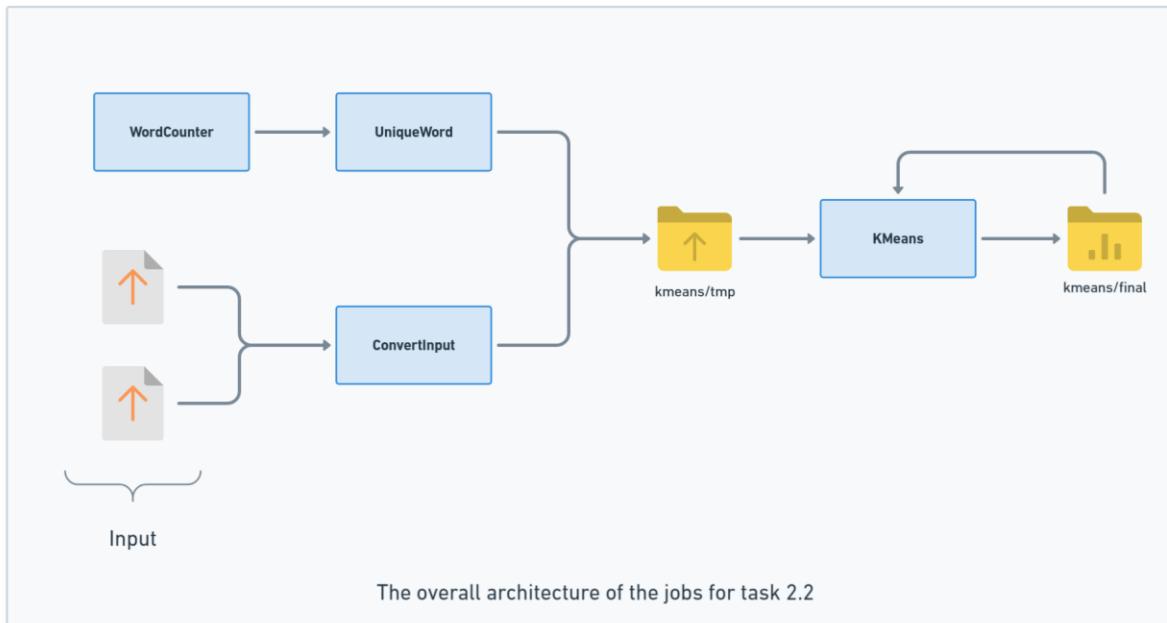
```

Figure 52: ConvertInput function

b. MapReduce Job Implementation:

- **In Map Phase:**
- Takes in an object key, Text value, and Context con.
- Splits the input text (value) by whitespace into an array.
- Extracts doc_id, term_id, and tf_idf (term frequency-inverse document frequency) from the array.
- Writes the key-value pair as <Text(doc_id), Text(term_id + ":" + tf_idf)> to the context.
- **In Reduce Phase:**

- Retrieves the total number of words (total_words) from the configuration.
 - Iterates over the values associated with a given key (Text type), constructing a string doc_tf_idf by concatenating the values with ";" as a separator.
 - If there are fewer values than total_words, it appends additional ";" to ensure the string has the correct format.
 - Writes the resulting doc_tf_idf string back to the context associated with the key.
- **KMeans Implementation:**
- The implementation of the K-Means algorithm will follow a structured sequence of steps:
 - **Initialization of Centroids:** The process begins by randomly initializing k centroids.
 - **Assignment of Data Points:** Each data point (document) in the dataset is then assigned to the cluster with the highest cosine similarity between the data point and the centroid.
 - **Centroid Update:** For each cluster, the centroids are recalculated based on the mean position of all data points assigned to that cluster.
 - **Iteration:** Steps 2 and 3 are repeated until the algorithm converges, meaning the centroids no longer change significantly.
 - While these steps outline the basic workflow of the K-Means algorithm, understanding its implementation using MapReduce requires further exploration, which is illustrated in the following graphs



Made with  Whimsical

Figure 53: The overall architecture

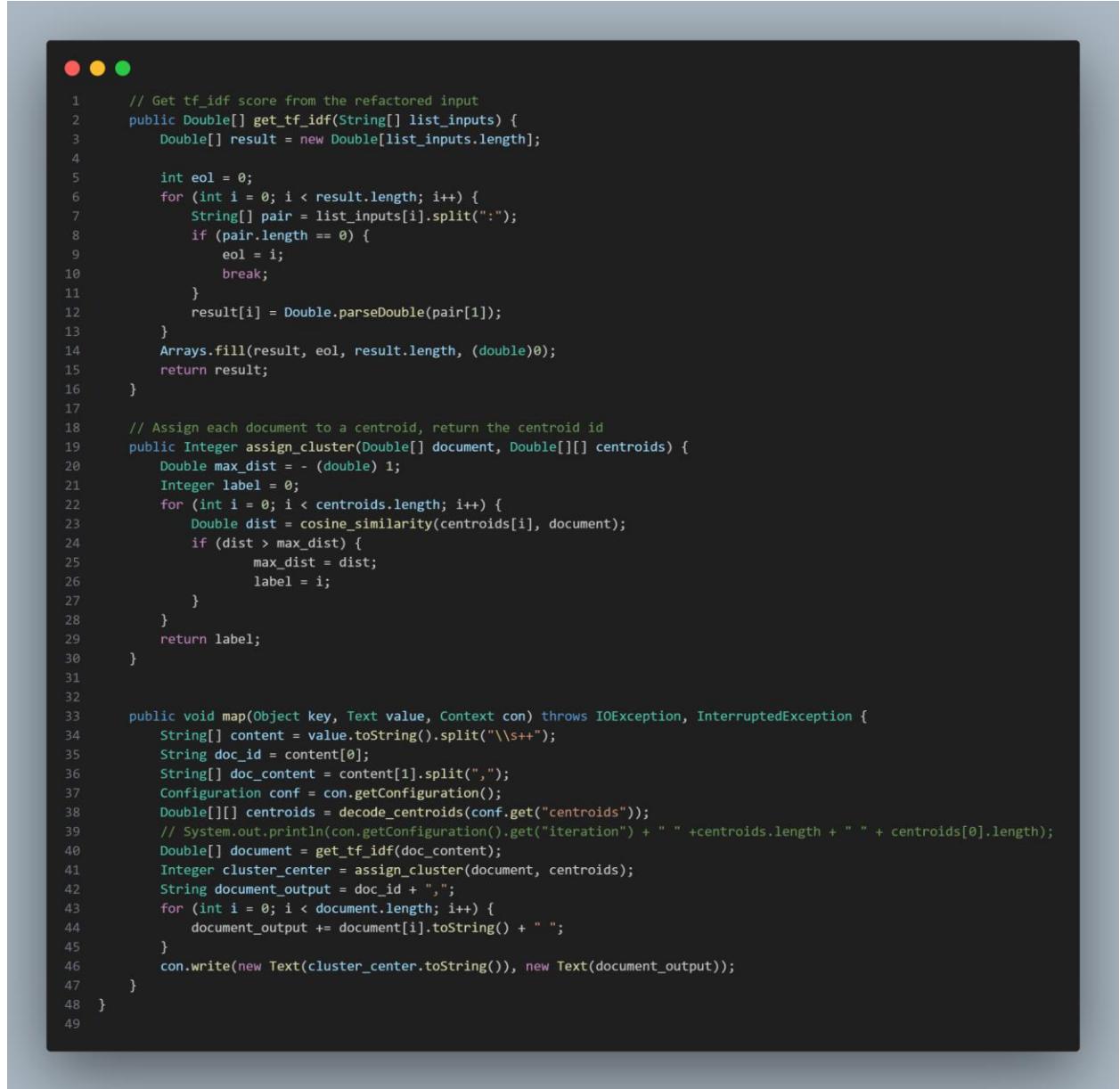
- The diagram illustrates the architecture of the job sequence for Task 2.2, which involves the implementation of the K-Means algorithm using MapReduce. The process begins with three distinct preprocessing jobs: **UniqueWord**, **WordCounter**, and **ConvertInput**. Each of these jobs performs a specific function to prepare the input data.
 - **UniqueWord** is responsible for identifying and isolating unique words within the dataset.
 - **WordCounter** counts the occurrences of words or documents, providing essential frequency data.
 - **ConvertInput** converts the data into a format suitable for subsequent processing by the K-Means algorithm.
- The outputs from these preprocessing jobs are directed into a temporary storage location, represented by the cylinder labeled */kmeans/tmp*. This intermediary

storage acts as a staging area where the preprocessed data is consolidated and made ready for the core K-Means processing.

- The KMeans job then processes the data from */kmeans/tmp*. This job involves the iterative steps of the K-Means algorithm, where data points are clustered, and centroids are recalculated until the algorithm converges.
- Finally, the results of the K-Means processing are stored in the final output location, labeled */kmeans/final*. The diagram also highlights the iterative nature of the K-Means algorithm, as indicated by the feedback loop from the KMeans job back to */kmeans/tmp*, allowing for continuous refinement of the clustering results until convergence is achieved.

```
● ● ●  
1  public class KMeansMapper extends Mapper<Object, Text, Text, Text>{  
2      // Returns the 2D array representation of the centroids  
3      public Double[][] decode_centroids(String centroid_string) {  
4          String[] initial_centroids = centroid_string.split(",");  
5          Integer k = initial_centroids.length;  
6          String[] sample_centroid = initial_centroids[0].split(" ");  
7          Double[][] centroids = new Double[k][sample_centroid.length];  
8          for (int i = 0; i < k; i++) {  
9              String[] ith_centroid = initial_centroids[i].split(" ");  
10             for (int j = 0; j < sample_centroid.length; j++) {  
11                 centroids[i][j] = Double.parseDouble(ith_centroid[j]);  
12             }  
13         }  
14         return centroids;  
15     }  
16  
17     // Utility function for cosine similarity  
18     public Double vector_length(Double[] vector) {  
19         Double squared = (Double) 0.0;  
20         for (int i = 0; i < vector.length; i++) {  
21             squared += vector[i] * vector[i];  
22         }  
23  
24         return (Double) Math.sqrt((double)squared);  
25     }  
26  
27     // Calculate cosine similarity between centroid and sample  
28     public Double cosine_similarity(Double[] vector_a, Double[] vector_b) {  
29         if (vector_a.length != vector_b.length) {  
30             return - (Double) 9999.0;  
31         }  
32         if (vector_a.equals(vector_b)) {  
33             return (Double) 1.0;  
34         }  
35  
36         Double dot_product = (Double) 0.0;  
37         Double sum_vector_length = (Double) 0.0;  
38         for (int i = 0; i < vector_a.length; i++) {  
39             dot_product += vector_a[i] * vector_b[i];  
40         }  
41         sum_vector_length = vector_length(vector_a) + vector_length(vector_b);  
42         if (sum_vector_length == 0.0)  
43             return - (Double) 9999.0;  
44         return (Double) dot_product / sum_vector_length;  
45     }
```

Figure 54: Map Phase



```

1 // Get tf_idf score from the refactored input
2 public Double[] get_tf_idf(String[] list_inputs) {
3     Double[] result = new Double[list_inputs.length];
4
5     int eol = 0;
6     for (int i = 0; i < result.length; i++) {
7         String[] pair = list_inputs[i].split(":");
8         if (pair.length == 0) {
9             eol = i;
10            break;
11        }
12        result[i] = Double.parseDouble(pair[1]);
13    }
14    Arrays.fill(result, eol, result.length, (double)0);
15    return result;
16 }
17
18 // Assign each document to a centroid, return the centroid id
19 public Integer assign_cluster(Double[] document, Double[][] centroids) {
20     Double max_dist = - (double) 1;
21     Integer label = 0;
22     for (int i = 0; i < centroids.length; i++) {
23         Double dist = cosine_similarity(centroids[i], document);
24         if (dist > max_dist) {
25             max_dist = dist;
26             label = i;
27         }
28     }
29     return label;
30 }
31
32
33 public void map(Object key, Text value, Context con) throws IOException, InterruptedException {
34     String[] content = value.toString().split("\\s+");
35     String doc_id = content[0];
36     String[] doc_content = content[1].split(",");
37     Configuration conf = con.getConfiguration();
38     Double[][] centroids = decode_centroids(conf.get("centroids"));
39     // System.out.println(conf.getConfiguration().get("iteration") + " " + centroids.length + " " + centroids[0].length);
40     Double[] document = get_tf_idf(doc_content);
41     Integer cluster_center = assign_cluster(document, centroids);
42     String document_output = doc_id + ",";
43     for (int i = 0; i < document.length; i++) {
44         document_output += document[i].toString() + " ";
45     }
46     con.write(new Text(cluster_center.toString()), new Text(document_output));
47 }
48 }
49

```

Figure 55: Map Phase

Define decode_centroids function:

Input: A string representing the centroids, separated by commas.

Output: A 2D array of Double representing the centroids.

Steps:

- Split the centroid string by commas to get individual centroids.
- Determine the number of centroids (k) and the dimensionality of each centroid (from the first centroid).

- Initialize a 2D array centroids of size k x dimensionality.
- For each centroid:
 - Split the centroid string by spaces to get its components.
 - Parse each component to Double and store it in the 2D array.
- Return the 2D array centroids.

Define vector_length function:

Input: A vector of Double.

Output: The length (magnitude) of the vector.

Steps:

- Initialize squared as 0.
- For each component in the vector:
 - Add the square of the component to squared.
- Return the square root of squared.

Define cosine_similarity function:

Input: Two vectors of Double, vector_a and vector_b.

Output: The cosine similarity between the two vectors.

Steps:

- If the lengths of the two vectors are not equal, return an error value (-9999.0).
- If the two vectors are identical, return a similarity of 1.0.
- Initialize dot_product as 0.
- For each pair of components in the two vectors:
 - Multiply the components and add the result to dot_product.
- Calculate the sum of the lengths of the two vectors.
- If the sum of the lengths is 0, return an error value (-9999.0).
- Return the ratio of dot_product to the sum of the vector lengths as the cosine similarity.

Define get_tf_idf function:

Input: An array of strings, each representing a term and its tf-idf value separated by a colon.

Output: An array of Double representing the tf-idf values.

Steps:

- Initialize an array result of Double of the same length as the input.
- For each string in the input:
- Split the string by the colon to separate the term and the tf-idf value.
- If splitting fails, break the loop and fill the rest of the result array with 0.
- Parse the tf-idf value to Double and store it in the result array.
- Return the result array.

Define assign_cluster function:

Input: A Double array representing a document, and a 2D array of Double representing the centroids.

Output: An integer representing the index of the closest centroid.

Steps:

- Initialize max_dist as -1 and label as 0.
- For each centroid:
- Calculate the cosine similarity between the centroid and the document.
- If the similarity is greater than max_dist, update max_dist and set label to the current centroid index.
- Return the label.

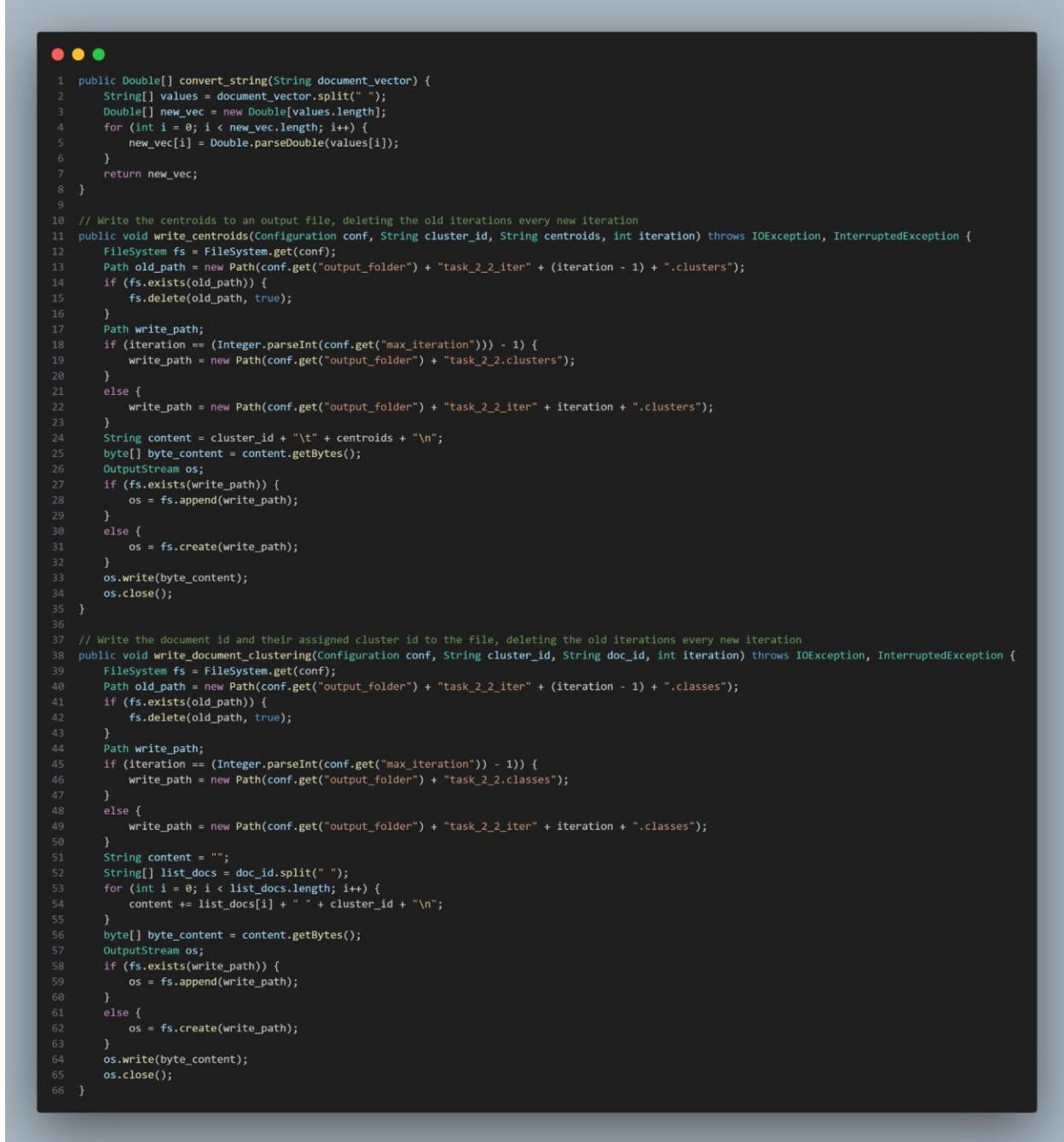
Define map function:

Input: A key, a text value (representing a document), and a context object.

Steps:

- Split the value into the document ID and document content.
- Get the centroids from the configuration using decode_centroids.
- Convert the document content to an array of tf-idf values using get_tf_idf.
- Assign the document to the closest centroid using assign_cluster.

- Prepare the document output as a string with the document ID followed by its tf-idf values.
- Write the centroid index and document output to the context.



```

1 public Double[] convert_string(String document_vector) {
2     String[] values = document_vector.split(" ");
3     Double[] new_vec = new Double[values.length];
4     for (int i = 0; i < new_vec.length; i++) {
5         new_vec[i] = Double.parseDouble(values[i]);
6     }
7     return new_vec;
8 }
9
10 // Write the centroids to an output file, deleting the old iterations every new iteration
11 public void write_centroids(Configuration conf, String cluster_id, String centroids, int iteration) throws IOException, InterruptedException {
12     FileSystem fs = FileSystem.get(conf);
13     Path old_path = new Path(conf.get("output_folder") + "task_2_2_iter" + (iteration - 1) + ".clusters");
14     if (fs.exists(old_path)) {
15         fs.delete(old_path, true);
16     }
17     Path write_path;
18     if (iteration == (Integer.parseInt(conf.get("max_iteration")) - 1)) {
19         write_path = new Path(conf.get("output_folder") + "task_2_2.clusters");
20     }
21     else {
22         write_path = new Path(conf.get("output_folder") + "task_2_2_iter" + iteration + ".clusters");
23     }
24     String content = cluster_id + "\t" + centroids + "\n";
25     byte[] byte_content = content.getBytes();
26     OutputStream os;
27     if (fs.exists(write_path)) {
28         os = fs.append(write_path);
29     }
30     else {
31         os = fs.create(write_path);
32     }
33     os.write(byte_content);
34     os.close();
35 }
36
37 // Write the document id and their assigned cluster id to the file, deleting the old iterations every new iteration
38 public void write_document_clustering(Configuration conf, String cluster_id, String doc_id, int iteration) throws IOException, InterruptedException {
39     FileSystem fs = FileSystem.get(conf);
40     Path old_path = new Path(conf.get("output_folder") + "task_2_2_iter" + (iteration - 1) + ".classes");
41     if (fs.exists(old_path)) {
42         fs.delete(old_path, true);
43     }
44     Path write_path;
45     if (iteration == (Integer.parseInt(conf.get("max_iteration")) - 1)) {
46         write_path = new Path(conf.get("output_folder") + "task_2_2.classes");
47     }
48     else {
49         write_path = new Path(conf.get("output_folder") + "task_2_2_iter" + iteration + ".classes");
50     }
51     String content = "";
52     String[] list_docs = doc_id.split(" ");
53     for (int i = 0; i < list_docs.length; i++) {
54         content += list_docs[i] + " " + cluster_id + "\n";
55     }
56     byte[] byte_content = content.getBytes();
57     OutputStream os;
58     if (fs.exists(write_path)) {
59         os = fs.append(write_path);
60     }
61     else {
62         os = fs.create(write_path);
63     }
64     os.write(byte_content);
65     os.close();
66 }

```

Figure 56: Reduce Phase

```
1 // Write the top 10 tf_idf values from the centroid
2 public void write_top_10(Configuration conf, String cluster_id, String top_10, int iteration) throws IOException, InterruptedException{
3     FileSystem fs = FileSystem.get(conf);
4     String content = "";
5     if (conf.get("has_written_mean").toString() == "0") {
6         content += "\nIteration " + (iteration + 1) + ": \n";
7         conf.set("has_written_mean", "1");
8     }
9     Path write_path = new Path(conf.get("output_folder") + "task_2_2.txt");
10    content += top_10 + "\n";
11    byte[] byte_content = content.getBytes();
12    OutputStream os;
13    if (fs.exists(write_path)) {
14        os = fs.append(write_path);
15    }
16    else {
17        os = fs.create(write_path);
18    }
19    os.write(byte_content);
20    os.close();
21 }
22
23 // Write WCSS loss to the file
24 public void write_losses(Configuration conf, String cluster_id, Double loss, int iteration) throws IOException, InterruptedException {
25     FileSystem fs = FileSystem.get(conf);
26     String content = "";
27     if (conf.get("has_written_loss").toString() == "0") {
28         content += "\nIteration " + (iteration + 1) + ": \n";
29         conf.set("has_written_loss", "1");
30     }
31     Path write_path = new Path(conf.get("output_folder") + "task_2_2.losses");
32     content += loss.toString();
33     content += "\n";
34     byte[] byte_content = content.getBytes();
35     OutputStream os;
36     if (fs.exists(write_path)) {
37         os = fs.append(write_path);
38     }
39     else {
40         os = fs.create(write_path);
41     }
42     os.write(byte_content);
43     os.close();
44 }
```

Figure 57: Reduce Phase

```
1  public Double calculate_loss(Double[] data, Double[] centroid) {
2      Double wcss = (double) 0;
3      for (int i = 0; i < data.length; i++) {
4          wcss += Math.pow(centroid[i] - data[i], 2);
5      }
6      return wcss;
7  }
8
9  public String get_top_10(Double[] centroid) {
10     TreeMap<Double, Integer> top_10 = new TreeMap<Double, Integer>();
11     for (int i = 0; i < centroid.length; i++) {
12         top_10.put(centroid[i], i);
13     }
14     String top_10_mean = "";
15     NavigableMap<Double, Integer> descending_top_10 = top_10.descendingMap();
16     Iterator tmp = descending_top_10.entrySet().iterator();
17     for (int i = 0; i < 10; i++) {
18         top_10_mean += tmp.next().toString() + ", ";
19     }
20     return top_10_mean;
21 }
22
23
24 public void reduce(Text key, Iterable<Text> value, Context con) throws IOException, InterruptedException {
25     Configuration conf = con.getConfiguration();
26     Integer iteration = Integer.parseInt(conf.get("iteration"));
27     String old_centroids = conf.get("centroids");
28     String[] num_centroids = old_centroids.split(",");
29     int elements = num_centroids[0].split(" ").length;
30     Double[] new_centroid = new Double[elements];
31     Arrays.fill(new_centroid, Double.valueOf(0));
32     int counter = 0;
33     String list_docs = "";
34     String cluster_id = key.toString();
35     Double[] data = new Double[elements];
36     Double wcss = (double) 0;
37     for (Text v: value) {
38         String[] true_val = v.toString().split(",");
39         String doc_id = true_val[0];
40         list_docs += doc_id + " ";
41         String doc_content = true_val[1];
42         Double[] tmp = convert_string(doc_content);
43         for (int i = 0; i < tmp.length; i++) {
44             data[i] = tmp[i];
45         }
46         counter += 1;
47         for (int i = 0; i < new_centroid.length; i++) {
48             new_centroid[i] += tmp[i];
49         }
50         wcss += calculate_loss(data, new_centroid);
51     }
52     String new_centroid_str = "";
53     for (int i = 0; i < new_centroid.length; i++) {
54         new_centroid[i] /= (counter + 1);
55         new_centroid_str += new_centroid[i].toString() + " ";
56     }
57     String top_10 = get_top_10(new_centroid);
58     write_centroids(conf, cluster_id, new_centroid_str, iteration);
59     write_document_clustering(conf, cluster_id, list_docs, iteration);
60     write_losses(conf, cluster_id, wcss, iteration);
61     write_top_10(conf, cluster_id, top_10, iteration);
62     con.write(key, new Text(new_centroid_str));
63 }
```

Figure 58: Reduce Phase

Define convert_string function:

Input: A string representing a document vector with values separated by spaces.

Output: An array of Double representing the document vector.

Steps:

- Split the input string by spaces to get individual values.
- Initialize a Double array of the same length as the values.
- Convert each value to Double and store it in the array.
- Return the Double array.

Define write_centroids function:

Input: Configuration object, cluster ID, centroid string, and iteration number.

Output: Writes centroids to an output file.

Steps:

- Initialize the file system using the configuration.
- Delete the old centroid file if it exists (from the previous iteration).
- Determine the write path based on the current iteration:
- If it's the last iteration, use the final cluster file path.
- Otherwise, use the iteration-specific file path.
- Convert the centroid data to bytes.
- Write the centroid data to the determined path, either appending to the file if it exists or creating a new file.
- Close the file output stream.

Define write_document_clustering function:

Input: Configuration object, cluster ID, document ID, and iteration number.

Output: Writes the document ID and their assigned cluster ID to a file.

Steps:

- Initialize the file system using the configuration.
- Delete the old document clustering file if it exists (from the previous iteration).
- Determine the write path based on the current iteration:

- If it's the last iteration, use the final document cluster file path.
- Otherwise, use the iteration-specific file path.
- For each document ID in the list, prepare content by appending the document ID and cluster ID.
- Convert the content to bytes.
- Write the document clustering data to the determined path, either appending to the file if it exists or creating a new file.
- Close the file output stream.

Define write_top_10 function:

Input: Configuration object, cluster ID, top 10 TF-IDF values string, and iteration number.

Output: Writes the top 10 TF-IDF values for centroids to a file.

Steps:

- Initialize the file system using the configuration.
- If not already noted, write the iteration header to the file (only once per iteration).
- Determine the file path for writing top 10 values.
- Convert the top 10 data to bytes.
- Write the top 10 data to the determined path, either appending to the file if it exists or creating a new file.
- Close the file output stream.

Define write_losses function:

Input: Configuration object, cluster ID, loss value, and iteration number.

Output: Writes the Within-Cluster Sum of Squares (WCSS) loss to a file.

Steps:

- Initialize the file system using the configuration.
- If not already noted, write the iteration header to the file (only once per iteration).
- Determine the file path for writing loss values.
- Convert the loss data to bytes.

- Write the loss data to the determined path, either appending to the file if it exists or creating a new file.
- Close the file output stream.

Define calculate_loss function:

Input: Document vector and centroid vector as arrays of Double.

Output: The WCSS loss (a Double value).

Steps:

- Initialize wcss as 0.
- For each component in the document vector and corresponding component in the centroid vector:
 - Calculate the squared difference between the centroid and document components and add it to wcss.
- Return the wcss value.

Define get_top_10 function:

Input: Centroid vector as an array of Double.

Output: A string representing the top 10 TF-IDF values from the centroid.

Steps:

- Initialize a TreeMap to store centroid values and their indices.
- For each value in the centroid vector, add it to the TreeMap.
- Get a descending order of the TreeMap and iterate through it.
- Extract the top 10 values from the map and convert them to a string format.
- Return the formatted string of the top 10 values.

Define reduce function:

Input: Key (Text), Iterable of values (Text), and context object.

Output: Writes the updated centroid information and other metrics.

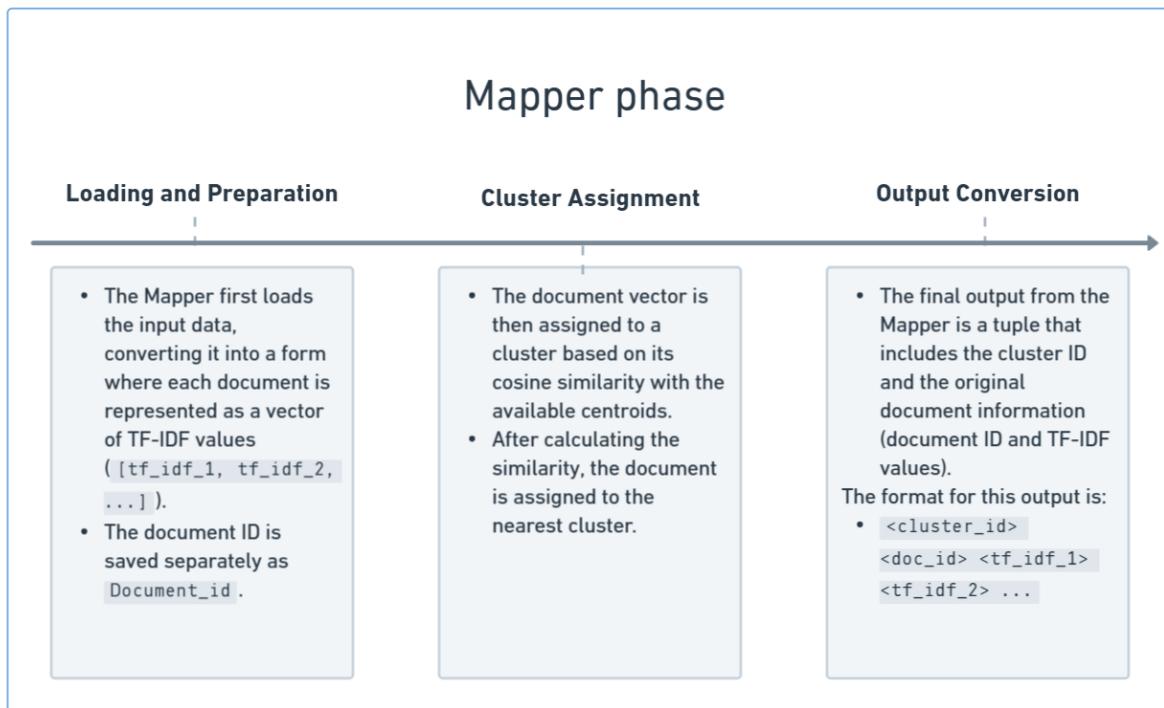
Steps:

- Initialize the configuration and retrieve the current iteration.
- Get the old centroids from the configuration.
- Determine the number of elements in each centroid.

- Initialize a new centroid array and set all elements to 0.
- Initialize a counter and an empty string for storing document IDs.
- Initialize a WCSS value to 0.
- For each value in the input:
 - Parse the document ID and its content.
 - Convert the content string to a document vector.
 - Update the new centroid by adding the document vector components to it.
 - Increment the counter.
 - Calculate the WCSS loss and add it to the total WCSS.
- Calculate the final centroid by dividing the sum by the number of documents.
- Convert the centroid to a string format.
- Get the top 10 TF-IDF values from the centroid.
- Write the centroid, document clustering, WCSS loss, and top 10 TF-IDF values to the appropriate files.
- Write the new centroid information to the context.

c. Task Flow:

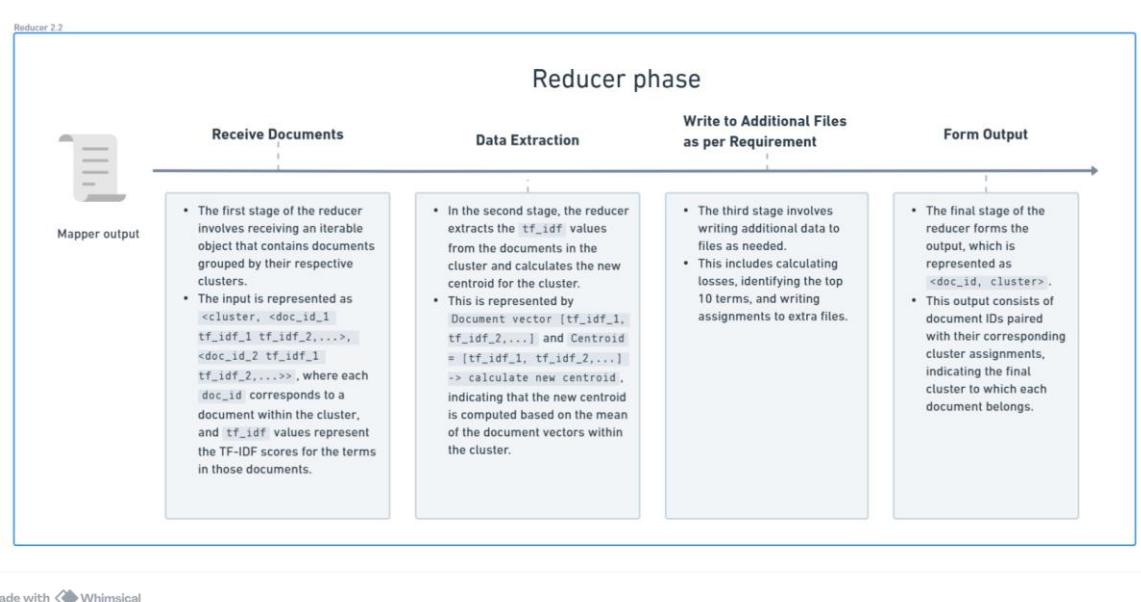
- The **map phase** operates as follows:



Made with  Whimsical

Figure 59: Maper phase for task 2.2

- In the initial iteration, the process begins by using randomly generated centroids. These centroids are serialized into a string and stored in the job configuration as a variable, enabling the job to keep track of centroid coordinates. After reading the input produced by the conversion job, the map function deserializes the string into an array of Doubles. It then calculates the similarity between each document and the centroids, assigning each document to a cluster. The cluster is used as the key for the reduce phase, and the value is formatted as: `cluster_id doc_id tf_idf_1 tf_idf_2`



Made with Whimsical

Figure 60: Reducer phase for task 2.2

- The reduce phase, which processes the output from the map phase, functions as follows: Upon receiving all documents belonging to the same cluster, the reduce phase iterates through them, adds them to a new centroid array, and keeps track of the document count. The new centroid array is then divided by this count, completing the centroid updating step described in the general plan. The losses are calculated as the sum of all squared Euclidean distances between each document in the cluster and the newly formed centroid. The top 10 terms are also identified, sorted, and written to a separate file. During this phase, all four required output files are generated without using the `context.write` function provided by MapReduce, in order to minimize follow-up tasks. However, a minor issue arises with the top 10 terms task, as no terms are transferred from the map phase to the reduce phase. This results in the output file `task_2_2.txt` only containing the index of words in the document vector, requiring an additional step to trace back to the actual term index.
- From the second iteration onward, the centroid configuration variable is updated using the `task_2_2.clusters` file. Each iteration has its own output directory,

allowing users to easily monitor the execution process. However, the final four output files can be found in the /kmeans/final/ directory as configured by the program.

d. How to run

A screenshot of a terminal window with a dark background. At the top left are three small colored circles: red, yellow, and green. Below them is a command line prompt followed by a single digit '1'. To the right of the prompt is a blue command line entry.

```
1 $(sudo) hadoop jar /path/to/jar package.path.to.KMeans /path/to/input /path/to/output.
```

Figure 61: Command

e. Result



Figure 62: task_2_2.cluster

```

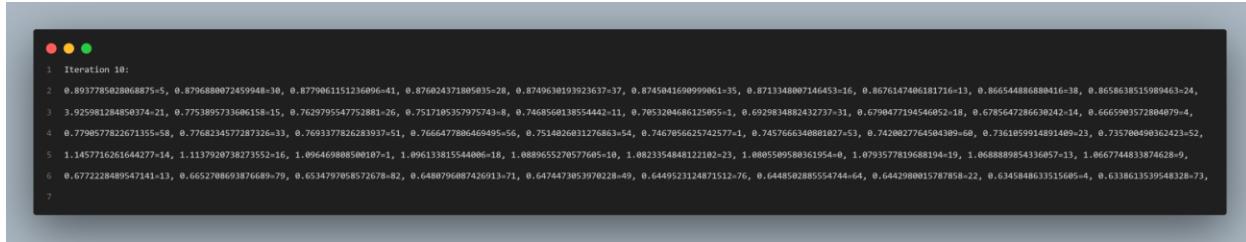
● ● ●
1 Iteration 1:
2 2.310286454603482E9
3 3.5724283593192964E9
4 9.87175495449976E8
5 4.95287530239782E8
6 7478950.737230793
7
8 Iteration 2:
9 1.652864188791038E10
10 4.260079409940085E8
11 5.204232008023394E8
12 432707.216830549
13 4.0428898376282483E8
14
15 Iteration 3:
16 2.6385916895248405E10
17 1345649.9017333952
18 8.426304159984188E8
19 57297.31143758421
20 7.186884609013995E8
    
```

Figure 63: task_2_2.loss

```

● ● ●
0 0.8553311247464271 0.8228767734227889 0.81332280773292 0.7964381208378652 0.8412551354722899 0.8937795828868875 0.8315268124782974 0.81614095005301569 0.841256926532928 0.7909328270281297 0.8393378656889882 0.8572497746668953 0.818251064336181 0.86761474685181716 0.78889531442270368
    
```

Figure 64: task_2_2.classes



```

1 Iteration 10:
2 0.8937785928068875=5, 0.8796880072459948=30, 0.8779061151236096=41, 0.876024371805035=28, 0.8749630193923637=37, 0.8745041690999061=35, 0.8713348007146453=16, 0.8676147406181716=13, 0.866544886880416=38, 0.8658638515989463=24,
3 3.925981284850374=21, 0.7753895733606158=15, 0.762979554752881=26, 0.7517105357975743=8, 0.7468568138554442=11, 0.7053204686125955=1, 0.6929834882432737=21, 0.6790477194546052=18, 0.6785647286630242=14, 0.6665983572884079=4,
4 0.7798577822671355=58, 0.776823457727326=33, 0.7693377826283937=51, 0.7666477886469495=56, 0.751426831276863=54, 0.74678562742577=5, 0.7457666340881027=53, 0.7428027764504309=60, 0.7361859914891409=23, 0.735708498362423=52,
5 1.1457716261644277=14, 1.1137920738273552=16, 1.096469808500107=1, 1.096133815544006=18, 1.0889655270577605=10, 1.0823354848122102=23, 1.0805509580361954=0, 1.0793577819688194=19, 1.0688889854336957=13, 1.0667744833874628=9,
6 0.6772228489547141=13, 0.6652708693876689=29, 0.6534797058572678=82, 0.648079687426913=71, 0.6474473053970228=49, 0.6449523124871512=76, 0.6448502885554744=64, 0.64429800515787858=22, 0.6345848633515605=4, 0.6338613539548328=73,
7

```

Figure 65: task_2_2.txt

4.3. Task 2.3 Scalable K-Means++ Initialization

a. Data description:

- In this task 2.3, we used the same method to process data as task 2.2

b. MapReduce Job Implementation:

- The K-Means++ initialization algorithm can be summarized as follows:

- **Random Sampling:** A random sample point is selected from the entire dataset, referred to as X . This sample is added to a set C and subsequently removed from the dataset.

Note: Each time a point is added to set C , it is removed from the dataset.

- **Distance Calculation:** The distance between every point in X and each point in C is calculated. These distances are then summed and stored in a variable, such as *sum_distances*.
- **Iterative Sampling:** A for-loop is initiated, with the number of iterations set according to the value of *log(sum_distances)*:

- a. Each point x_i in X is independently sampled with a probability

$$\frac{l \times \minDist(x_i, C)}{\sum \text{distances}} \quad .$$

- b. Then, the sampled points are added to set C .

- **Weight Calculation:** An empty array w , matching the size of set C , is created. For each sample point in X , the closest centroid in C is identified, and the corresponding index (e.g., j) in w is incremented by 1.

- **Centroid Selection:** The weighted array w is returned, and to select k centroids, k elements are chosen from C , sorted in descending order.

c. Challenges and Solutions:

- A significant challenge arises when handling large datasets, as the algorithm requires calculating the distance between every point in X and the potential centroids in C . A simple solution might involve generating a random key and processing all data from the map phase on a single node during the reduce phase. However, this approach is unsuitable for larger datasets due to the risk of memory overflow.
- To mitigate this issue, an ensemble technique was implemented. The job configuration allows for specifying the number of nodes (e.g., 5 nodes) to work on this task. During the map phase, each line of data is assigned a random key, ranging from 0 to the number of configured nodes. This approach effectively splits dataset X into smaller subsets, enabling each node to process its mini-dataset independently. After completing their respective reduce tasks, which involve applying the K-Means \parallel algorithm to the smaller datasets, the output from each node is compiled into a single file. The results are then averaged to obtain the desired k centroids.
- A visual representation of the full MapReduce process for this task is provided in the accompanying graph.

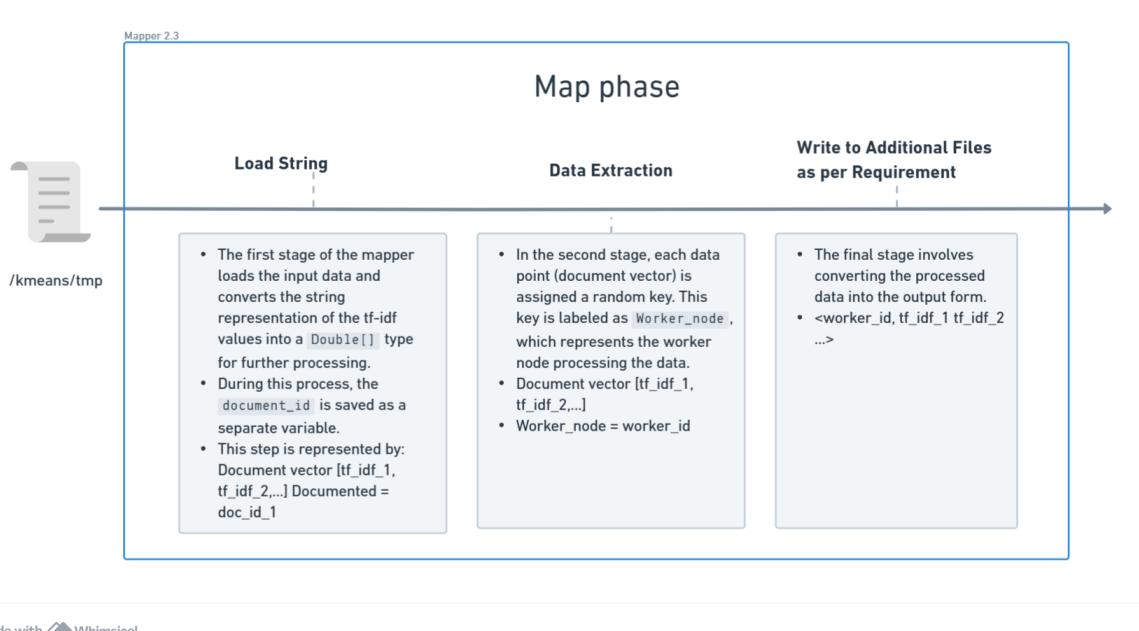


Figure 66: Map Phase for task 2.3

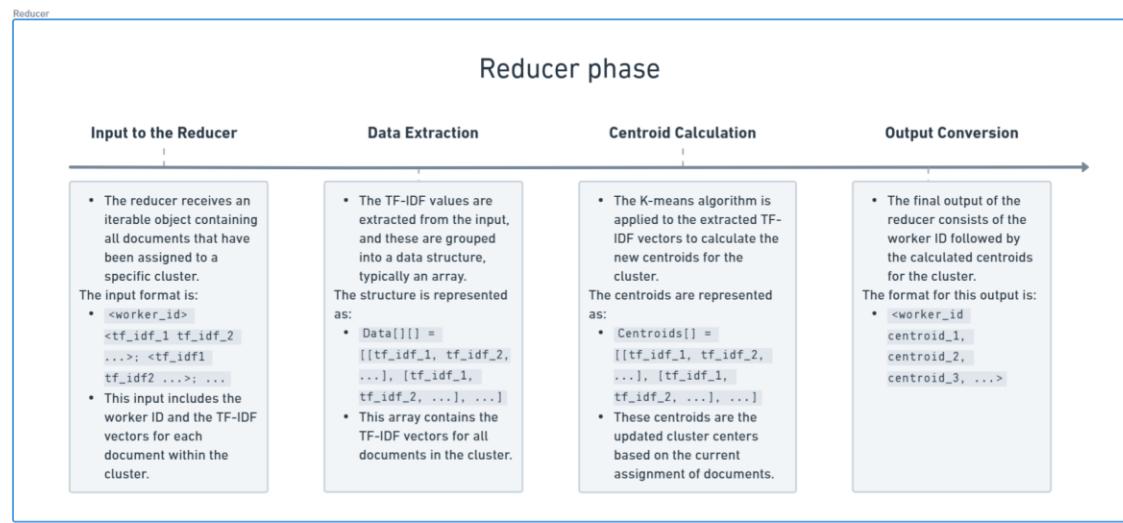


Figure 67: Reducer Phase for task 2.3

d. How to run:



```
1 $(sudo) hadoop jar /path/to/jar package.path.to.KMeans /path/to/input /path/to/output.
```

Figure 68: Command

e. Result:



```
1 Iteration 8:  
2 5.901838850016779E9  
3 2.230024966961907E10  
4 1.2941278959564645E7  
5 149116.60700060264  
6 135761.18430515437  
7  
8 Iteration 9:  
9 5.616679424753652E9  
10 2.2341630846910957E10  
11 1.6060629475042786E7  
12 143602.26711959386  
13 177725.7443308732  
14  
15 Iteration 10:  
16 5.529911594810488E9  
17 2.0789031648855812E10  
18 3.706785827790922E7  
19 124220.28765376497  
20 199968.42493400522
```

Figure 9: task_2_3.loss

```

1 Iteration 10:
2 0.6888444692245171=61, 0.6799104889462771=54, 0.6777465672286088=56, 0.6754884737735082=13, 0.6750986865108283=33, 0.6716387425596393=49, 0.6681148669337812=58, 0.6650919715862481=50, 0.6646342275489378=26, 0.6638396125468137=53,
3 0.9235828795683888=21, 0.9199573499254021=21, 0.9872906383832491=16, 0.0938711451380522=49, 0.9824111370614811=13, 0.8984543527225972=5, 0.89438958262425624=8, 0.8938339837982497=30, 0.8902231361351493=22, 0.8901912800847696=1,
4 2.4869431344378554=12, 1.1567842897655687=18, 1.0857204878186699=15, 1.0356853798859105=11, 1.0275163931912168=5, 1.01338802557180422=1, 1.08488841148496=17, 1.0022377361887955=21, 0.9994668845771387=14, 0.9941652642878671=7,
5 3.2630391229361893=63, 0.8048839832245164=25, 0.7391886436987047=9, 0.7089176964606829=48, 0.6784038537291337=13, 0.6771237190387108=51, 0.6628235028297742=24, 0.6512318259371458=32, 0.6453424884600705=52, 0.6433134687582589=42,
6 3.507438283125435=60, 0.7904634285976168=47, 0.7452391138156684=33, 0.6981959202987297=58, 0.6837828476767289=24, 0.6710583749562588=5, 0.6453820819577459=43, 0.630372513968637=48, 0.630125947802431=48, 0.6227272907185638=49,
7

```

Figure 70: task_2_3.txt

5. Conclusion

Through this lab exercise, our group significantly deepened our understanding of the **MapReduce mechanism** and its **practical applications** in solving complex problems.

Initially, we encountered **challenges**, particularly with the programming aspects of MapReduce and the scarcity of related documentation. However, by engaging in **collaborative discussions** and **meticulously studying the MapReduce framework**, we were able to overcome these difficulties. As we became more familiar with the concepts, the implementation process became progressively smoother.

In summary, this lab was invaluable for our group. It not only enhanced our **comprehension of data processing techniques** but also provided us with hands-on experience in **implementing the K-means clustering algorithm using Hadoop MapReduce**. The skills and insights gained from this exercise will undoubtedly benefit us in future projects involving large-scale data processing.

6. References

- [1] [Apache Hadoop MapReduce Tutorial](#)
- [2] [MapReduce Algorithms for k-means Clustering](#)
- [3] [How to calculate TF-IDF with MapReduce - Slogix](#)