# Big Data Systems: Assignment1

Deadline 11 07, 2024 at 23:59 PM

*LEE CHIH-PIN NDHU CSIE*

# Problem 1

In Hadoop cluster, we will implement the sorting algorithm. The input is a list of string

## Pseudo code

```
1: import sys
2: for line in sys.stdin do
3:     line ← line.strip()                              ▷ Remove leading and trailing whitespace
4:     words ← line.split()                                         ▷ Split line into words
5:     for word in words do
6:         print (word, 1)                              ▷ Output word with count 1, tab-separated
7:     end for
8: end for
```

Algorithm 1: Sort Mapper

```
1: import sys
2: word_count ← {}                              ▷ Initialize an empty dictionary for word counts
3: for line in sys.stdin do
4:     line ← line.strip()                              ▷ Remove leading and trailing whitespace
5:     parts ← line.split("\t", 1)                              ▷ Split line by tab into two parts
6:     if len(parts) < 2 then
7:         continue                                             ▷ Skip incomplete lines
8:     end if
9:     word, count ← parts                                      ▷ Extract word and count
10:    try:
11:        count ← int(count)                                   ▷ Convert count to integer
12:    except ValueError:
13:        continue                                             ▷ Skip lines with invalid counts
14:    if word in word_count then
15:        word_count[word] ← word_count[word] + count
16:    else
17:        word_count[word] ← count                             ▷ Add word with initial count
18:    end if
19: end for
20: for word in sorted(word_count) do
21:     print (word, word_count[word])               ▷ Output word and its total count, tab-separated
22: end for
```

Algorithm 2: Sort Reducer

## Explanation

In this section, we will explain the pseudo code of the sorting algorithm using MapReduce.

The algorithm consists of two main components: the Mapper and the Reducer. The Mapper reads the input data and emits key-value pairs, while the Reducer aggregates the key-value pairs and produces the final output.

---

The Mapper reads the input data line by line and splits each line into words. For each word, the Mapper outputs a key-value pair where the key is the word and the value is 1. This allows the Reducer to count the occurrences of each word.

The Reducer reads the key-value pairs emitted by the Mapper and aggregates the counts for each word. It uses a dictionary to store the word counts and updates the count for each word as it processes the input. Finally, the Reducer sorts the words in lexicographical order and outputs the word and its total count.

The sorting algorithm using MapReduce is an example of the MapReduce programming model, which is commonly used for processing large datasets in parallel. By dividing the input data into smaller chunks and processing them in parallel, MapReduce allows for efficient processing of large datasets on distributed systems.

## Result

The result of sorting in terminal and Hadoop website is shown in Figure 1 and Figure 2 respectively.



Figure 1: The result of sorting in terminal



Figure 2: The result of sorting is successful status in Website

# Problem 2

Now, we will implement the Searching algorithms. The input is a list of string and the output is the word that we want to search.

## Pseudo code

1: **import** sys, os
2: search_word ← os.getenv('SEARCH_WORD', 'default_word')      ▷ Get search word from environment variable
3: **for** line **in** sys.stdin **do**
4:      filepath ← os.environ.get('map_input_file', 'unknown_file')
5:      filename ← os.path.basename(filepath)      ▷ Get filename from file path
6:      line ← line.strip()      ▷ Remove leading and trailing whitespace
7:      line_num ← 0      ▷ Initialize line number (for demonstration purposes)
8:      positions ← [ ]      ▷ Initialize empty list to store positions of search word
9:      start ← 0
10:      **while** True **do**
11:          start ← line.find(search_word, start)      ▷ Find position of search word
12:          **if** start = -1 **then**
13:              **break**      ▷ Exit loop if word not found
14:          **end if**
15:          positions.append(start)      ▷ Add position to list
16:          start ← start + len(search_word)      ▷ Move start to next position after found word
17:      **end while**
18:      **if** positions is not empty **then**
19:          **print** (filename, "_line_", line_num, line, positions) ▷ Output filename, line number, line content, and positions of word
20:      **end if**
21: **end for**

Algorithm 3: Search Mapper

## Explanation

In this section, we will explain the pseudo code of the searching algorithm using MapReduce.

The algorithm consists of two main components: the Mapper and the Reducer. The Mapper reads the input data and emits key-value pairs, while the Reducer aggregates the key-value pairs and produces the final output.

The Mapper reads the input data line by line and searches for the search word in each line. If the search word is found, the Mapper outputs a key-value pair where the key is the filename and the value is a tuple containing the line content and the positions of the search word in the line.

The Reducer reads the key-value pairs emitted by the Mapper and aggregates the matches by filename. It uses a dictionary to store the matches for each file and updates the matches as it processes the input. Finally, the Reducer outputs the filename and all matches for that file.

     4

```
1: import sys
2: from collections import defaultdict
3: import ast
4: file_matches ← defaultdict(list)                    ▷ Initialize a dictionary to store matches by filename
5: for line in sys.stdin do
6:     filename, content, positions ← line.strip().split("\t", 2) ▷ Split input line into filename, content, and
   positions
7:     positions ← ast.literal_eval(positions)                      ▷ Convert positions from string to list
8:     file_matches[filename].append((content, positions))    ▷ Add content and positions to the file's entry
9: end for
10: for current_file, matches in file_matches.items() do
11:     print (current_file, matches)                              ▷ Output filename and all matches
12: end for
```

Algorithm 4: Search Reducer

The searching algorithm using MapReduce is an example of the MapReduce programming model, which is commonly used for processing large datasets in parallel. By dividing the input data into smaller chunks and processing them in parallel, MapReduce allows for efficient processing of large datasets on distributed systems.

## Result

The result of searching in terminal and Hadoop website is shown in Figure 3 and Figure 4 respectively.



Figure 3: The result of search in terminal



Figure 4: The result of search is successful in website

# Problem 3

Futhermore, we will implement the TF-IDF computation algorithm. The input is a list of string and the output is the TF-IDF value of each word.

## Definition

TF-IDF (Term Frequency-Inverse Document Frequency) is a statistical measure used to evaluate the importance of a word in a document relative to a collection of documents. It is commonly used in information retrieval and text mining to identify the most relevant words in a document.

The TF-IDF value of a word is calculated based on two components: Term Frequency (TF) and Inverse Document Frequency (IDF). The TF component measures the frequency of a word in a document, while the IDF component measures the rarity of the word in a collection of documents.

The TF-IDF algorithm computes the TF-IDF value of each word in a document by combining the TF and IDF components. The TF component is calculated as the frequency of the word in the document divided by the total number of words in the document. The IDF component is calculated as the logarithm of the total number of documents divided by the number of documents containing the word.

1. **Input:** Corpus of documents $D = \{d_1, d_2, \ldots, d_n\}$ and terms $T = \{t_1, t_2, \ldots, t_m\}$

2. **Output:** TF-IDF score for each term in each document.

3. **For each document** $d \in D$:

   (a) **For each term** $t \in T$:

       i. Calculate Term Frequency (TF) of $t$ in $d$:

   $$TF(t, d) = \frac{\text{Number of times } t \text{ appears in } d}{\text{Total number of terms in } d}$$

       ii. Calculate Document Frequency (DF) of $t$ in $D$:

   $$DF(t) = \text{Number of documents containing } t$$

       iii. Calculate Inverse Document Frequency (IDF) of $t$:

   $$IDF(t, D) = \log\left(\frac{|D|}{DF(t)}\right)$$

       iv. Calculate TF-IDF score for $t$ in $d$:

   $$TF\text{-}IDF(t, d, D) = TF(t, d) \times IDF(t, D)$$

4. **End For**

## Pseudo code

## Explanation

In this section, we will explain the pseudo code of the TF-IDF computation algorithm using MapReduce.

```
1:  import sys, os, re
2:  from collections import defaultdict
3:  file_name ← os.getenv('mapreduce_map_input_file', 'unknown_file').split('/')[-1]      ▷ Get filename from
    environment variable
4:  word_counts ← defaultdict(int)                                          ▷ Initialize a dictionary for word counts
5:  total_words ← 0
6:  for line in sys.stdin do
7:      line ← line.strip().lower()              ▷ Remove leading/trailing whitespace and convert to lowercase
8:      words ← re.findall(r'\b\w+\b', line)                          ▷ Extract words using regular expressions
9:      for word in words do
10:         word_counts[word] ← word_counts[word] + 1
11:         total_words ← total_words + 1                                      ▷ Increment total word count
12:     end for
13: end for
14: for word, count in word_counts.items() do
15:     tf ← count / total_words                                      ▷ Calculate term frequency
16:     print (word, file_name, tf)                          ▷ Output word, filename, and term frequency
17: end for
```

<div align="center">Algorithm 5: TF-IDF Mapper</div>

```
1:  import sys, math
2:  from collections import defaultdict
3:  tf_values ← defaultdict(list)                          ▷ Dictionary to store term frequencies by word
4:  df_counts ← defaultdict(int)                   ▷ Dictionary to store document frequency counts for each word
5:  documents ← set()                                      ▷ Set to track unique document names
6:  for line in sys.stdin do
7:      word, file_name, tf ← line.strip().split("\t")
8:      tf ← float(tf)                                      ▷ Convert tf to a floating-point number
9:      tf_values[word].append((file_name, tf))                   ▷ Store tf value for each word and document
10:     df_counts[word] ← df_counts[word] + 1            ▷ Increment document frequency count for word
11:     documents.add(file_name)                                      ▷ Add document to the set
12: end for
13: total_docs ← len(documents)                          ▷ Calculate the total number of documents
14: print("DEBUG: Total documents =", total_docs)
15: for word, file_tfs in tf_values.items() do
16:     print("DEBUG: DF count for word '", word, "':", df_counts[word])
17:     idf ← math.log(total_docs / df_counts[word]) if df_counts[word] > 0 else 0      ▷ Calculate inverse
    document frequency
18:     for file_name, tf in file_tfs do
19:         tf_idf ← tf * idf                                      ▷ Calculate TF-IDF
20:         print(word, file_name, ":.6f".format(tf_idf))          ▷ Output word, filename, and TF-IDF score
21:     end for
22: end for
```

<div align="center">Algorithm 6: TF-IDF Reducer</div>

7

The algorithm consists of two main components: the Mapper and the Reducer. The Mapper reads the input data and emits key-value pairs, while the Reducer aggregates the key-value pairs and produces the final output.

The Mapper reads the input data line by line and extracts words using regular expressions. For each word, the Mapper calculates the term frequency (TF) by counting the number of occurrences of the word in the document and dividing by the total number of words in the document. The Mapper outputs a key-value pair where the key is the word, the value is the filename, and the term frequency.

The Reducer reads the key-value pairs emitted by the Mapper and aggregates the TF values by word. It also calculates the document frequency (DF) for each word by counting the number of documents containing the word. The Reducer then calculates the inverse document frequency (IDF) for each word and combines it with the TF values to compute the TF-IDF score. Finally, the Reducer outputs the word, filename, and TF-IDF score for each word in each document.

The TF-IDF computation algorithm using MapReduce is an example of the MapReduce programming model, which is commonly used for processing large datasets in parallel. By dividing the input data into smaller chunks and processing them in parallel, MapReduce allows for efficient processing of large datasets on distributed systems.

## Result

The result of TF-IDF in terminal and Hadoop website is shown in Figure 5 and Figure 6 respectively.
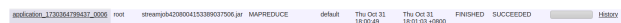


Figure 5: The result of TF IDF in terminal



Figure 6: The result of TF IDF is successful in website

# Problem 4

Finally, we will implement the Activity Mining algorithm. The input is a list of string and the output is the activity of each user.

## Definition

1. **Input:** Set of activity sequences $S = \{s_1, s_2, \ldots, s_n\}$, minimum support threshold $\theta$.

2. **Output:** Frequent activity patterns.

3. **Step 1: Preprocessing**

   (a) **For each sequence $s \in S$:**

      i. Tokenize each sequence $s$ into individual activities $A = \{a_1, a_2, \ldots, a_k\}$.

      ii. Store each unique activity in a list of candidate patterns $C$.

4. **Step 2: Pattern Generation and Counting**

   (a) Initialize an empty set $F$ for storing frequent patterns.

   (b) **For each candidate pattern $c \in C$:**

      i. **For each sequence $s \in S$:**

         A. **If $c$ is found as a subsequence in $s$ then**:
            - Increment the count of $c$.

      ii. **If** count of $c \geq \theta$:
            - Add $c$ to frequent patterns set $F$.

5. **Step 3: Pattern Expansion**

   (a) **For each pattern $f \in F$:**

      i. Generate expanded patterns by adding one additional activity to $f$.

      ii. Check support of each expanded pattern and retain patterns meeting $\theta$.

6. **Step 4: Output Results**

   (a) Return the set of frequent patterns $F$ as identified activity patterns.

## Pseudo code

## Explanation

In this section, we will explain the pseudo code of the Activity Mining algorithm using MapReduce.

The algorithm consists of two main components: the Mapper and the Reducer. The Mapper reads the input data and emits key-value pairs, while the Reducer aggregates the key-value pairs and produces the final output.

The Mapper reads the input data line by line and generates all possible subsets of activities of varying lengths. For each subset, the Mapper outputs a key-value pair where the key is the subset of activities and the value is 1. This allows the Reducer to count the occurrences of each subset.

```
1:  import sys
2:  from itertools import combinations
3:  MIN_LENGTH ← 2                                              ▷ Minimum length for subsets
4:  for line in sys.stdin do
5:      line ← line.strip().lower()          ▷ Remove leading/trailing whitespace and convert to lowercase
6:      words ← line.split()                                         ▷ Split line into words
7:      for length in range(MIN_LENGTH, len(words) + 1) do
8:          for subset in combinations(words, length) do
9:              print("-¿".join(subset), 1)              ▷ Output the subset as a joined string with count 1
10:         end for
11:     end for
12: end for
```

Algorithm 7: Activity Mining Mapper

```
1:  import sys
2:  SUPPORT_THRESHOLD ← 2                                          ▷ Define the support threshold
3:  current_combination ← None
4:  current_count ← 0
5:  for line in sys.stdin do
6:      line ← line.strip()
7:      if not line then
8:          continue                                               ▷ Skip empty lines
9:      end if
10:     combination, count ← line.split("\t")
11:     count ← int(count)                                         ▷ Convert count to an integer
12:     if combination = current_combination then
13:         current_count ← current_count + count          ▷ Accumulate count for the same combination
14:     else
15:         if current_combination and current_count ≥ SUPPORT_THRESHOLD then
16:             print(current_combination, current_count)      ▷ Output if count meets support threshold
17:         end if
18:         current_combination ← combination
19:         current_count ← count
20:     end if
21: end for
22: if current_combination and current_count ≥ SUPPORT_THRESHOLD then
23:     print(current_combination, current_count)             ▷ Final output if count meets threshold
24: end if
```

Algorithm 8: Activity Mining Reducer

10

The Reducer reads the key-value pairs emitted by the Mapper and aggregates the counts for each subset. It accumulates the counts for the same subset and outputs the subset and its count if the count meets the support threshold. The support threshold is a predefined value that determines the minimum number of occurrences required for a subset to be considered frequent.

## Result

The result of Activity Mining in terminal and Hadoop website is shown in Figure 7 and Figure 8 respectively.
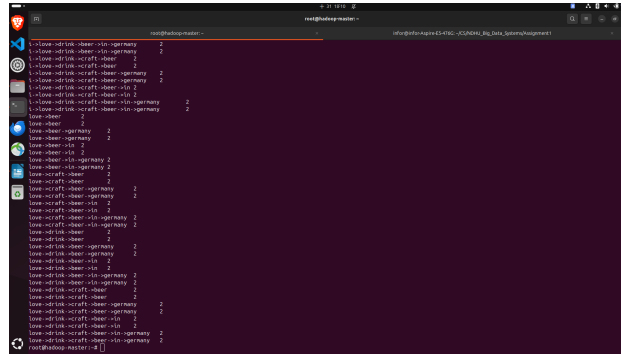


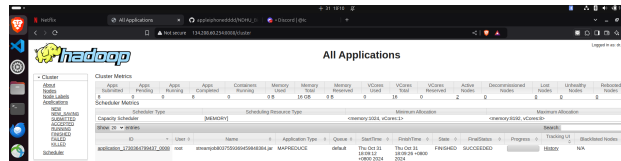Figure 7: The result of Activity Mining in terminal



Figure 8: The result of Activity Mining is successful in website