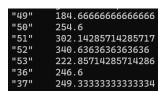


# **Assignment 1**

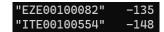
## Task 1:

- Steps: Define the sequence of steps in the MapReduce task. In this case, there is a step with a mapper function and a reducer function;
- Mapper Function (mapper\_get\_age): It produces a key-value pair, where the key is the age and the value is the number of friends (converted into an integer);
- Reducer Function (reducer\_average\_ages): It yields a key-value pair where the key is the age, and the value is the mean (average) of the values (number of friends) associated with that age.
- Result (extract):



## Task 2:

- Steps: Define the sequence of steps in the MapReduce task. In this case, there is one step with a mapper and a reducer function;
- Mapper Function (mapper\_get\_temp): It checks if the observation type is "TMIN" and,
  if so, yields a key-value pair where the key is the weather station and the value is the
  temperature (converted to an integer);
- Reducer Function (reducer\_min\_temp): It yields a key-value pair where the key is the
  weather station, and the value is the minimum temperature among the values
  associated with that weather station.
- Result:



## Task 3:

- "WORD\_RE = re.compile(r"[a-zA-Z0-9\_]+")": A regular expression used to find all the words in a character string. It matches alphanumeric characters and underscores;
- Steps: Define the sequence of steps in the MapReduce job. It consists of two steps: one for counting words and another for sorting and outputting the results. Each stage contains a mapper function and a reducer function;
- Counting words:
  - Mapper Function (mapper\_get\_words): It yields a key-value pair where the key is a lowercase word, and the value is 1;
  - Reducer Function (reducer\_count\_words): It yields a key-value pair where the key is a word, and the value is the sum of the counts (word frequency).
- Sorting and output of the results:
  - Mapper Function (mapper\_make\_counts\_key): It takes the key and value resulting from the last phase as input and produces a key-value pair, where the key is 'None' and the value is a tuple of the word count (converted to integer) and the word;
  - Reducer Function (reducer\_output\_words): It receives a key and an iterator
    of word-count pairs. Sorts the word-count pairs based on their count and
    produces key-value pairs in which the key is a word and the value is its
    count.

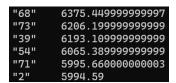


### Result (extract):



### Task 4:

- Steps: Define the sequence of steps in the MapReduce job. It consists of two steps: one for getting the total amount spent by each customer and another for sorting and outputting the results. Each stage contains a mapper function and a reducer function;
- Getting the total amount spent by each customer:
  - Mapper Function (mapper\_get\_amount): It yields a key-value pair where the key is the customer ID, and the value is the amount spent (converted to a float);
  - Reducer Function (reducer\_total\_amount): It yields a key-value pair where the key is the customer ID, and the value is the sum of the amount spent by that customer.
- Sorting and output of the results:
  - Mapper Function (mapper\_make\_amounts\_key): It takes the key and value resulting from the last phase as input and produces a key-value pair, where the key is 'None' and the value is a tuple of the total amount spent and the customer ID;
  - Reducer Function (reducer\_output\_amounts): It takes a key and an iterator of amount-customer pairs. It sorts the amount-customer pairs based on the amount and yields key-value pairs.
- Result (extract):



# **Assignment 2**

#### Task 1:

- Line 8: Division of RDDs into lists with the data we need (Weather Station, Observation Type, Temperature);
- Line 11: Filter only the lines that relate to the minimum temperature (as requested in the statement);
- Line 14: Convert the temperature to an integer, so that you can apply the min function:
- Line 17: For each capital (key), it returns only the lowest temperature value.
- Result:



# Task 2:

- "WORD\_RE = re.compile(r"[a-zA-Z0-9\_]+")": A regular expression used to find all the words in a character string. It matches alphanumeric characters and underscores.
- Line 8: Dividing RDDs into words (by spaces);

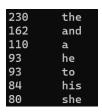


- Line 11: Map the value of each word to 1;
- Line 14: It groups the words (keys) and adds up their occurrences (as the value of each word defaults to 1, the sum of these values corresponds to the total number of occurrences).
- Result (extract):



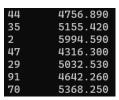
### Task 3:

- Until Line 14: Same as Task 2;
- Line 17: Reverse the key with the value, from (word, value) to (value, word), so that we can sort the occurrences, which are the key;
- Line 20: Sort the occurrences of the words, which are the key.
- Result (extract):



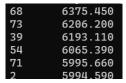
## Task 4:

- Line 8: Division of RDDs into lists with the data we need (Customer ID, Amount spent by customer);
- Line 11: Transform the amount spent by customers into a float, so that it can be added up:
- Line 13: It groups the customers (keys) and adds up the amount spent (values) for each.
- Result (extract):



## Task 5:

- Until Line 13: Same as Task 4;
- Line 15: Reverse the key with the value, from (Customer ID, Amount spent by customer) to (Amount spent by customer, Customer ID), so that we can sort the amount spent, these being the key;
- Line 17: Sort the amount customers spend, which is the key.
- Result (extract):





#### Task 6:

- Lines 10 and 12: Division of RDDs into lists with the data we need:
  - Marvel+Graph: Get the ID and the superhero name;
  - Marvel+Names: Each line represents a film and each ID represents a superhero who has appeared.
- Lines 14 and 15: Count the frequencies of each superhero ID;
- Line 17: Join the frequency RDD with the superhero names RDD using the superhero
   ID:
- Lines 19 and 21: Map the joined RDD to a new RDD with frequencies and superhero names, then sorts it in descending order by frequency.
- Lines 24 and 27: Retrieve the most popular frequency and filter superheroes with the same frequency.
- Result:

```
Most Popular Superhero: CAPTAIN AMERICA
Frequency: 1937
```

#### Task 7:

- Same as Task 6 but the sorting part is done in ascending order.
- Result (extract):

```
Least Popular Superhero: RED WOLF II
Frequency: 1

Least Popular Superhero: DEATHCHARGE
Frequency: 1

Least Popular Superhero: MARVEL BOY II/MARTIN
Frequency: 1
```

# **Assignment 3**

# Task 1:

- Line 8: Divide the read data into lists of what we need (Weather Station, Observation Type, Temperature);
- Line 11: Filter only the lines that relate to the minimum temperature;
- Line 14: Select only the weather station and temperatures from the three variables;
- Line 17: For each capital (grouped), return only the lowest temperature value.
- Result:



### Task 2:

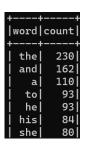
- Line 9: Replace non-alphanumeric characters with spaces;
- Line 12: Divide lines into words (by spaces);
- Line 15: Filtering the result from line 9 to remove any empty characters;
- Line 18: For each word (grouped), returns the number of occurrences (counting).
- Result (extract):





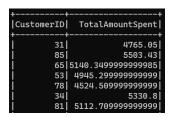
### Task 3:

- As in Task 2, but adding in line 21, sort the words in descending order according to their occurrences.
- Result (extract):



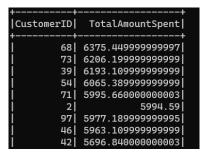
## Task 4:

- Line 9: Divide the data read into lists of what we need (Customer ID, amount spent by customer);
- Line 12: For each customer (grouped), return the amount spent (summing).
- Result (extract):



# Task 5:

- Up to Line 12: Same as Task 4;
- Line 15: Sort customers in descending order according to their expenses.
- Result (extract):



### Task 6:

- Line 10: Lines from "Marvel+Graph" are split into words using spaces as delimiters;
- Line 13: For "Marvel+Names" extract superhero IDs and names;
- Line 16: Count how many times each superhero appears in "Marvel+Graph";
- Line 19: Join frequency counts with superhero names using the IDs;
- Line 22: Sort the result by frequency in descending order;



- Lines 26 and 27: Retrieve the most popular frequency and filter superheroes with the same frequency.
- Result (extract):

```
Most Popular Superhero: "CAPTAIN AMERICA"
Frequency: 1937
```

### Task 7:

- Same as Task 6 but the sorting part is done in ascending order
- Result (extract):

```
Least Popular Superhero: "BERSERKER II"
Frequency: 1

Least Popular Superhero: "BLARE/"
Frequency: 1

Least Popular Superhero: "MARVEL BOY II/MARTIN"
Frequency: 1
```

# **Assignment 4**

## Task 1:

- Line 12: Dividing Lines (streams) into words (by the spaces)
- Line 15: Group the words and count their occurrences
- Line 17: Sort the occurrences of words
- Results (Input/Output):



