

BIG DATA & ANALYTICS

ASSIGNMENT 2: SPARK SQL

BACKGROUND.

It's been already a few weeks since you started your short-term internship in the Data Analytics Department of the start-up OptimiseYourJourney, which will enter the market next year with a clear goal in mind: *"leverage Big Data technologies for improving the user experience in transportation"*. Your contribution in Assignment 1 has proven the potential OptimiseYourJourney can obtain by applying MapReduce to analyse large-scale public transportation datasets as the one in the New York City Bike Sharing System: <https://www.citibikenyc.com/>

OptimiseYourJourney



In the department meeting that has just finished your boss was particularly happy, again.

- The very same dataset from Assignment 1 (let's call it my_dataset_1) provides an opportunity to leverage other large-scale data analysis libraries, such as Spark SQL.
- The streaming of a subset of the dataset allows you to explore the potential of Spark Structured Streaming, a library of Spark specialised in real-time data analysis.

FileStore/tables/DATASET 1: (folder my_dataset_1)

This dataset occupies ~80MB and contains 73 files. Each file contains all the trips registered the CitiBike system for a concrete day:

- 2019_05_01.csv => All trips registered on the 1st of May of 2019.
- 2019_05_02.csv => All trips registered on the 2nd of May of 2019.
- ...
- 2019_07_12.csv => All trips registered on the 12th of July of 2019.

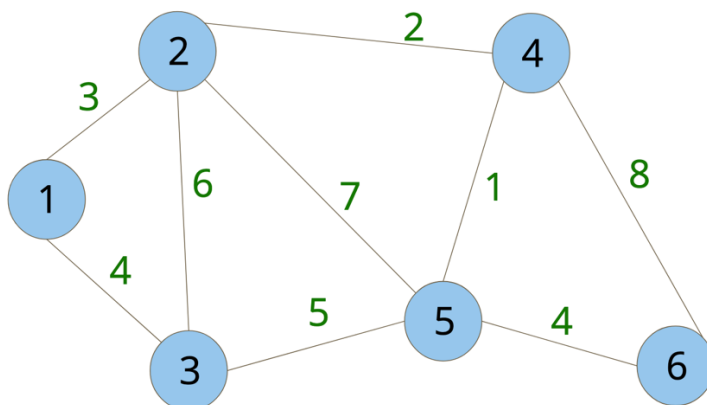
Altogether, the files contain 444,110 rows. Each row contains the following fields:
start_time , *stop_time* , *trip_duration* , *start_station_id* , *start_station_name* ,
start_station_latitude , *start_station_longitude* , *stop_station_id* , *stop_station_name* ,
stop_station_latitude , *stop_station_longitude* , *bike_id* , *user_type* , *birth_year* , *gender* ,
trip_id

- **(00) *start_time***
 - A String representing the time the trip started at.
<%Y/%m/%d %H:%M:%S>
 - Example: “2019/05/02 10:05:00”
- **(01) *stop_time***
 - A String representing the time the trip finished at.
<%Y/%m/%d %H:%M:%S>
 - Example: “2019/05/02 10:10:00”
- **(02) *trip_duration***
 - An Integer representing the duration of the trip.
 - Example: 300
- **(03) *start_station_id***
 - An Integer representing the ID of the CityBike station the trip started from.
 - Example: 150
- **(04) *start_station_name***
 - A String representing the name of the CitiBike station the trip started from.
 - Example: “E 2 St & Avenue C”.
- **(05) *start_station_latitude***
 - A Float representing the latitude of the CitiBike station the trip started from.
 - Example: 40.7208736
- **(06) *start_station_longitude***
 - A Float representing the longitude of the CitiBike station the trip started from.
 - Example: -73.98085795
- **(07) *stop_station_id***
 - An Integer representing the ID of the CityBike station the trip stopped at.
 - Example: 150
- **(08) *stop_station_name***
 - A String representing the name of the CitiBike station the trip stopped at.
 - Example: “E 2 St & Avenue C”.

- **(09) *stop_station_latitude***
 - A Float representing the latitude of the CitiBike station the trip stopped at.
 - Example: 40.7208736
- **(10) *stop_station_longitude***
 - A Float representing the longitude of the CitiBike station the trip stopped at.
 - Example: -73.98085795
- **(11) *bike_id***
 - An Integer representing the id of the bike used in the trip.
 - Example: 33882.
- **(12) *user_type***
 - A String representing the type of user using the bike (it can be either “Subscriber” or “Customer”).
 - Example: “Subscriber”.
- **(13) *birth_year***
 - An Integer representing the birth year of the user using the bike.
 - Example: 1990.
- **(14) *gender***
 - An Integer representing the gender of the user using the bike (it can be either 0 => Unknown; 1 => male; 2 => female).
 - Example: 2.
- **(15) *trip_id***
 - An Integer representing the id of the trip.
 - Example: 190.

FileStore/tables/DATASET 2: (folder my_dataset_2)

This dataset consists in the file `tiny_graph.txt`, which contains 18 edges (indeed, 9 edges, one on each direction) in a graph with 6 nodes.



MY CODE

This folder **my_code** contains 6 subfolders *A02_Part?*, one folder per exercise.

Each folder *A02_Part?* contains the file *A02_Part?.py*

This is the file specifying the exercise. You must complete it.

MY RESULTS

This folder **my_results** contains the following subfolders:

- **Assignment_Solutions:**

This is the folder with the expected solutions to the exercises.

- **Student_Solutions:**

This is the folder with the actual solutions to the exercises (based on what you have programmed).

- **test_checker.py:**

This is a program you can use to ensure your exercises in Parts 1, 2 and 3 are indeed producing the expected result. To do so, modify the line 74 of the program with the exercise you want to test (e.g., value 2 to test the exercise 2) and run the program. If it prints “True” means that you have passed the test; otherwise “False” means the result of your query differs from the expected one.

- **my_A02_Part4_check_results.py:**

This is a program you can use to ensure you exercise in Part 4 is indeed producing the expected result. To do so, run the program. If it prints “True” means that you have passed the test; otherwise “False” means the result of your query differs from the expected one.

Main Message

Use the programs **my_A02_check_results.py** and **my_A02_Part4_check_results.py** to ensure that all your Spark SQL and Spark Structured Streaming exercises produce the expected output (and in the right format!).

TASKS / EXERCISES.

The tasks / exercises to be completed as part of the assignment are described in the next pages:

- The following exercises are placed in the folder **my_code**:
 1. **A02_Part1/A02_Part1.py**
 2. **A02_Part2/A02_Part2.py**
 3. **A02_Part3/A02_Part3.py**
 4. **A02_Part4/A02_Part4.py**
 5. **A02_Part5/A02_Part5.py**
 6. **A02_Part6/A02_Part6.py**

Marks are as follows:

1. **A02_Part1/A02_Part1.py** => 15 marks
2. **A02_Part2/A02_Part2.py** => 15 marks
3. **A02_Part3/A02_Part3.py** => 15 marks
4. **A02_Part4/A02_Part4.py** => 20 marks
5. **A02_Part5/A02_Part5.py** => 15 marks
6. **A02_Part6/A02_Part6.py** => 20 marks

Tasks:

1. **A02_Part1/A02_Part1.py**
2. **A02_Part2/A02_Part2.py**
3. **A02_Part3/A02_Part3.py**
4. **A02_Part4/A02_Part4.py**
5. **A02_Part6/A02_Part6.py**

Complete the function **my_main** of the Python program.

Do not modify the name of the function nor the parameters it receives.

The entire work must be done within Spark SQL:

- The function **my_main** must start with the creation operation read above loading the dataset to Spark SQL.
- The function **my_main** must finish with an action operation collect, gathering and printing by the screen the result of the Spark SQL job.
- The function **my_main** must not contain any other action operation collect other than the one appearing at the very end of the function.
- The **resVAL** iterator returned by collect must be printed straight away, you cannot edit it to alter its format for printing.

6. **A02_Part5/A02_Part5.py**

Complete the function **my_spark_part compute graph** of the Python program.

Same rules as for the functions **my_main** of previous parts apply.

Complete the function **my_python part compute page rank** of the Python program.

This function is to be completed using pure Python (not Spark SQL) and, therefore, the rules described before do not apply anymore.

RUBRIC.

Exercises 1-6.

- 20% of the marks => Complete attempt of the exercise (even if it does not lead to the right solution or right format due to small differences).
- 40% of the marks => Right solution and format (following the aforementioned rules) for the provided dataset.
- 40% of the marks => Right solution and format (following the aforementioned rules) for any “Additional Dataset” test case we will generate. The marks will be allocated in a per test basis (i.e., if 4 extra test are tried, each of them will represent 10% of the marks).

TEST YOUR SOLUTIONS.

- The folder **my_results** contains the expected results for each exercise.
 - **A02_Part1**/result.txt
 - **A02_Part2**/result.txt
 - **A02_Part3**/result.txt
 - **A02_Part4**/result.txt
 - **A02_Part5**/result.txt
 - **A02_Part6**/result.txt
- Moreover, the subfolder **my_results/check_results** allows you to see if your code is producing the expected output or not.
 - The file **test_checker.py** needs two folders and compares if their files are equal or not.
 - When you have completed one part (e.g., **A02_Part1**), copy the folder **my_results/A02_Part1** into the folder **my_results/check_results/Student_Attempts/A02_Part1**.
 - Open the file **test_checker.py** and edit the line 104 with the value of the part you are attempting (e.g., `part = 1`).
 - Run the program **test_checker.py**. It will tell you whether your output is correct or not.

For example, as an example let's run the Python program **test_checker.py** to see if the solution attempt done by the student for **A02_Part1** and **A02_Part1** is correct or not.

➤ `python3.13 test_checker.py 1`

```
-----  
Checking :  
./Assignment_Solutions/A02_Part1/result.txt  
./Student_Attempts/A02_Part1/result.txt
```

```
Test passed!
```

```
-----  
Congratulations, the code passed all the tests!  
-----
```

As we can see, the code of the student is correct, and thus it gets the marks.

➤ python3.13 test_checker.py 4

Checking :

./Assignment_Solutions/A02_Part4/result.txt

./Student_Attempts/A02_Part4/result.txt

Test did not pass.

Sorry, the output of some files is incorrect!

As we can see, the code of the student is not correct, and thus it does not get the marks. The problem was that some output lines in some files were wrong.

Main Message

Use the program **test_checker.py** to ensure that all your exercises produce the expected output (and in the right format!).

SUBMISSION DETAILS / SUBMISSION CODE OF CONDUCT.

Submit to Canvas by the 2nd of May, 5pm.

- Submissions up to 1 week late will have 10 marks deducted.
- Submissions up to 2 weeks late will have 20 marks deducted.

On submitting the assignment you adhere to the following declaration of authorship. If you have any doubt regarding the plagiarism policy discussed at the beginning of the semester do not hesitate in contacting me.

Declaration of Authorship

I, ____ NAME____, declare that the work presented in this assignment, titled *Big Data & Analytics - Assignment 2: Spark SQL*. I confirm that:

- This work was done wholly by me as part of my BSc. in Software Development.
- Where I have consulted the published work and source code of others, this is always clearly attributed.
- Where I have quoted from the work of others, the source is always given. With the exception of such quotations, this assignment source code is entirely my own work.

EXERCISE – A02 Part1

(15 marks)

Technology:

Spark SQL.

Your task is to:

- Compute the amount of trips starting from and finishing at each station_name.

Complete the function my_main of the Python program.

- Do not modify the name of the function nor the parameters it receives.
- The entire work must be done within Spark SQL:
 - The function my_main must start with the creation operation 'read' above loading the dataset to Spark SQL.
 - The function my_main must finish with an action operation 'collect', gathering and printing by the screen the result of the Spark SQL job.
 - The function my_main must not contain any other action operation 'collect' other than the one appearing at the very end of the function.
 - The resVAL iterator returned by 'collect' must be printed straight away, you cannot edit it to alter its format for printing.

Results:

Output one Row per station_name. Rows must follow alphabetic order in the name of the station. Each Row must have the following fields:

```
Row(station, num_departure_trips, num_arrival_trips)
```


EXERCISE – A02 Part2

(15 marks)

Technology:

Spark SQL.

Your task is to:

- Compute the top_n_bikes with highest total duration time for their trips.

Complete the function my_main of the Python program.

- Do not modify the name of the function nor the parameters it receives.
- The entire work must be done within Spark SQL:
 - The function my_main must start with the creation operation 'read' above loading the dataset to Spark SQL.
 - The function my_main must finish with an action operation 'collect', gathering and printing by the screen the result of the Spark SQL job.
 - The function my_main must not contain any other action operation 'collect' other than the one appearing at the very end of the function.
 - The resVAL iterator returned by 'collect' must be printed straight away, you cannot edit it to alter its format for printing.

Results:

Output one Row per bike_id. Rows must follow decreasing order in highest total duration time for their trips. Each Row must have the following fields:

`Row(bike_id, totalTime, numTrips)`

EXERCISE – A02 Part3

(15 marks)

Technology:

Spark SQL.

Your task is to:

- Sometimes bikes are re-organised (moved) from station A to station B to balance the amount of bikes available in both stations. A truck operates this bike re-balancing service, and the trips done by-truck are not logged into the dataset. Compute all the times a given bike_id was moved by the truck re-balancing system.

Complete the function **my_main** of the Python program.

- Do not modify the name of the function nor the parameters it receives.
- The entire work must be done within Spark SQL:
 - The function **my_main** must start with the creation operation 'read' above loading the dataset to Spark SQL.
 - The function **my_main** must finish with an action operation 'collect', gathering and printing by the screen the result of the Spark SQL job.
 - The function **my_main** must not contain any other action operation 'collect' other than the one appearing at the very end of the function.
 - The resVAL iterator returned by 'collect' must be printed straight away, you cannot edit it to alter its format for printing.

Results:

Output one Row per moving trip. Rows must follow temporal order. Each Row must have the following fields:

`Row(start_time, start_station_name, stop_time, stop_station_name)`

For example, if the dataset **contains** the following 2 trips:

- **Trip1:** A user used bike_id to start a trip from Station1 on 2019/05/10 09:00:00 and finished the trip at Station2 on 2019/05/10 10:00:00
- **Trip2:** A user used bike_id to start a trip from Station3 on 2019/05/10 11:00:00 and finished the trip at Station4 on 2019/05/10 12:00:00

And the dataset **does not contain** any extra trip:

- **Trip3:** A user used bike_id to start a trip from Station2 and finish at Station3 anytime between 2019/05/10 10:00:00 and 2019/05/10 11:00:00

Then it is clear that the bike was moved from Station2 to Station3 by truck, and we output:

`Row(start_time=2019/05/10 10:00:00, start_station_name=Station2, stop_time=2019/05/10 11:00:00, stop_station_name=Station3)`

EXERCISE – A02 Part4.

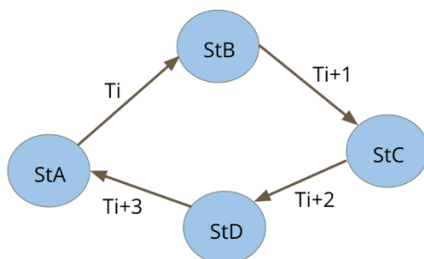
(20 marks)

Technology:

Spark SQL.

Your task is to:

- We say k consecutive trips carried out by a given `bike_id` represent a loop if the following conditions fullfil:
 - Given the n trips of a given `bike_id`, ordered in chronological order, $\langle T_1, T_2, \dots, T_n \rangle$ a subset of k consecutive trips $\langle T_i, T_{i+1}, \dots, T_{i+(k-1)} \rangle$ is selected.
 - The start station of T_i and the end station of $T_{i+(k-1)}$ are the same.
 - For each consecutive pair of trips (T_a, T_b) in $\langle T_i, T_{i+1}, \dots, T_{i+(k-1)} \rangle$ the end station of T_a and the start station of T_b are the same.
 - The total amount of different stations visited in $\langle T_i, T_{i+1}, \dots, T_{i+(k-1)} \rangle$ is k .
- Below there is an example of a loop of 4 consecutive trips.



- Given a parameter `loop_size` specifying the size of a loop, compute all loops of size `loop_size` performed by the bikes of the dataset.

Complete the function `my_main` of the Python program.

- Do not modify the name of the function nor the parameters it receives.
- The entire work must be done within Spark SQL:
 - The function `my_main` must start with the creation operation 'read' above loading the dataset to Spark SQL.
 - The function `my_main` must finish with an action operation 'collect', gathering and printing by the screen the result of the Spark SQL job.
 - The function `my_main` must not contain any other action operation 'collect' other than the one appearing at the very end of the function.
 - The `resVAL` iterator returned by 'collect' must be printed straight away, you cannot edit it to alter its format for printing.

Results:

Output one Row per loop (order the Rows by increasing `bike_id`), with the following fields:

```
Row(bike_id, start_time_ti, stop_time_ti, start_station_name_ti, stop_station_name_ti, ...,
start_time_ti+k-1, stop_time_ti+k-1, start_station_name_ti+k-1, stop_station_name_ti+k-1)
```

EXERCISE 5.

(15 marks)

Technology:

Spark SQL and Pure Python.

Your task is to:

- Compute your own sequential implementation of the PageRank algorithm for the graph formed by the stations (nodes) and trips (edges) of the NYC dataset.

Complete the function **my_spark_part_compute_graph** of the Python program.

- Do not modify the name of the function nor the parameters it receives.
- The entire work must be done within Spark SQL:
 - The function `my_main` must start with the creation operation 'read' above loading the dataset to Spark SQL.
 - The function `my_main` must finish with an action operation 'collect', gathering and printing by the screen the result of the Spark SQL job.
 - The function `my_main` must not contain any other action operation 'collect' other than the one appearing at the very end of the function.
 - The `resVAL` iterator returned by 'collect' must be printed straight away, you cannot edit it to alter its format for printing.

Results:

Output one Row per station, with the following fields:

```
Row(start_station_name, neighbours, num_neighbours)
```

Complete the function **my_python_part_compute_page_rank** of the Python program.

- Do not modify the name of the function nor the parameters it receives.
- The function must return a dictionary with (key, value) pairs, where:
 - Each key represents a node id.
 - Each value represents the pagerank value computed for this node id.

Results:

Given the requested dictionary, the program automatically outputs one (key, value) pair per line. Lines follow a decreasing order in the page rank value of the node. Each line has the following format:

```
id=key ; pagerank=value \n
```

EXERCISE 6.

(20 marks)

Technology:

Spark SQL.

Your task is to:

- Using Spark SQL, compute the shortest path distance from a source node to the remaining nodes of the graph.

Complete the function my_main of the Python program.

- Do not modify the name of the function nor the parameters it receives.
- The entire work must be done within Spark SQL:
 - The function my_main must start with the creation operation 'read' above loading the dataset to Spark SQL.
 - The function my_main must finish with an action operation 'collect', gathering and printing by the screen the result of the Spark SQL job.
 - The function my_main must not contain any other action operation 'collect' other than the one appearing at the very end of the function.
 - The resVAL iterator returned by 'collect' must be printed straight away, you cannot edit it to alter its format for printing.
- The difficulty of this exercise is on coming up with your own Spark SQL implementation of the Dijkstra shortest path algorithm. That is, you must implement a Spark SQL program following the steps of the Dijkstra algorithm explained in class.

Results:

Output one Row per bike_id. Rows must follow a decreasing order in the cost of the path. Each Row must have the following fields:

Row(id, cost, path)