Univ.-Prof. Dr.-Ing. Matthias Boehm

Graz University of Technology Computer Science and Biomedical Engineering Institute of Interactive Systems and Data Science BMVIT endowed chair for Data Management

4 Data Management WS19: Exercise 04 – Large-Scale Data Analysis

Published: Dec 31, 2019 (last update: Dec 31) Deadline: Jan 21, 2020, 11.59pm

This exercise on large-scale data analysis aims to provide practical experience with distributed data management and large-scale data analysis on top of Apache Spark. The expected result is a zip archive named DB_Exercise04_<student_ID>.zip, submitted in TeachCenter.

4.1 Apache Spark Setup (4/25 points)

As a preparation step, setup Apache Spark and necessary Hadoop client APIs inside an IDE (integrated development environment) of your language choice. This exercise can be done with the Spark language bindings Java, Scala, or Python. For example in Java, you could simply include the maven dependencies spark-core and spark-sql into your project. On Windows, please download winutils.exe from https://github.com/steveloughran/winutils/tree/master/ hadoop-2.7.1/bin, put it into a directory <some-path>/hadoop/bin, and create a new environment variable HADOOP_HOME=<some-path>/hadoop. The input data for this exercise is available at https://mboehm7.github.io/teaching/ws1920_dbs/data.zip (from Exercise 3, based on the schema from Exercise 2).

Partial Results: N/A.

4.2 SQL Query Processing (5/25 points)

In order to further practice basic SQL query processing, please create the following two SQL queries. You get 2.5 points per query as this is primarily a repetition but note that these queries are the input for tasks 4.3 and 4.4.

- **Q09**: What are the top-5 cities—considering all their airports—by total number of route departures? (return city name and number of departures in descending order of number of departures).
- **Q10:** Which plane types are used on more than 2048 routes? (return plane type name and number of routes it is used on).

Partial Results: SQL script Queries.sql.

4.3 Query Processing via Spark RDDs (10/25 points)

Spark's fundamental abstraction for distributed collections are so-called Resilient Distributed Datasets (RDDs). In this task, you should implement queries Q09 and Q10 from task 4.2 via RDD operations, collect the results in the driver and print the result list to stdout. Please implement these queries as two self-contained functions/methods executeQ09RDD() and executeQ10RDD() that internally create a SparkContext sc, read the files via sc.textFile(), and use only RDD operations to compute the query results.

Partial Results: Source file QueriesRDD.*.

4.4 Query Processing via Spark SQL (6/25 points)

Spark also provides the high-level APIs Dataframe and Dataset for SQL processing. In this task, you should implement queries Q09 and Q10 from task 4.2 via Dataset operations, and write the outputs to JSON files out09.json and out10.json. Please implement these queries as two self-contained functions/methods executeQ09Dataset() and executeQ10Dataset() that internally create a SparkSession sc, read the inputs files via sc.read().format("csv"), and use only SQL or Dataset operations to compute and write the query results. You might either (1) register the individual input Datasets as temporary views and compute the results directly via SQL, or (2) alternatively use the functional API of Datasets. Both specifications share a common query optimization and processing pipeline.

Partial Results: Source file QueriesDataset.*.

4.5 Extra Credit (5 points)

Please provide the following SQL Query as a SQL script ExtraQueries.sql:

• Q11: What is the longest route in km, as computed via the Haversine distance from the longitude and latitude of departure and arrival airports? (return departure city name, arrival city name, distance in km).