



Middle East Technical University
Department of Computer Engineering

CENG 495
Cloud Computing
Spring 2018-2019 Homework 3

Due Date: 19.05.2019, 23:55

This homework aims to get you familiar with MapReduce paradigm. You are going to develop and deploy a MapReduce application by using Apache Hadoop Packages and Java language.

Keywords: Cloud Computing, Hadoop, Apache, MapReduce, Java

1. Apache Hadoop

- Download and install the latest stable release of Apache Hadoop.
- You must have the required JDK version to use Hadoop.
- See the useful links section.

2. Specifications

- You will implement a Java code with Hadoop environment to analyze the input files consisting of football game results.
- You will be given a folder containing input text files.
- The inputs will have the following form:

`<HOME_TEAM > <AWAY_TEAM> <RESULT>`

- RESULT can have three different values. “0” means draw, “1” means home win, and “2” means away team win.
- If a game ends as a draw, both teams get 1 point each. Otherwise, winner gets 3 points, and loser gets no points from that game.

- Your program will execute the following tasks:
 - a. List the number of games played for each team. (**game**)
 - b. List the total points gained by each team. (**point**)
 - c. List the average points per game for each team. (**avg**)
 - d. List the total number of “DRAW”, “HOME_WIN” and “AWAY_WIN” on different files. i.e. number of draws on “part-r-00000”, number of home team wins on “part-r-00001”, and number of away team wins on “part-r-00002”. (**stat**)
- Teams names will be consisted of lowercase ASCII characters from ‘a’ to ‘z’. Thus, you do not need to handle uppercase letters or any other characters.
- The outputs of MapReduce are sorted according to the keys by default, thus you do not need to change anything for the order of the outputs.
- There can be more than one input file. Your program should read all the files in the input folder.
- You can see the input and output formats on the sample input and output files. Since black-box testing will be used for grading, be sure to stick to the format.
- Your code must be in Java language using the Apache Hadoop library.
- Your codes will be evaluated automatically in Local (Standalone) Mode of Hadoop. Assuming that all of the Java files of your solution exist in the current directory, the command sequence below will be executed in order to build the solution:

hadoop com.sun.tools.javac.Main *.java

jar cf Hw3.jar *.class

- The output jar file will be tested with commands given below with different inputs.

hadoop jar Hw3.jar Hw3 game input output_g

hadoop jar Hw3.jar Hw3 point input output_p

hadoop jar Hw3.jar Hw3 avg input output_a

hadoop jar Hw3.jar Hw3 stat input output_s

3. Useful Links

- Apache Hadoop: <http://hadoop.apache.org/>
- To download: <http://kozyatagi.mirror.guzel.net.tr/apache/hadoop/common/stable/>
- Install guide: https://hadoop.apache.org/docs/r3.2.0/hadoop-project-dist/hadoop-common/SingleCluster.html#Installing_Software (Note that the most common problem is to forget to set the environment variables on file “hadoop-env.sh”)
- You can look at the following tutorial and use the corresponding code as a base for your work: <https://hadoop.apache.org/docs/r3.2.0/hadoop-mapreduce-client/hadoop-mapreduce-client-core/MapReduceTutorial.html>

4. Submission

- In this assignment, you are expected to submit your Java code(s) to ODTÜClass. For submission on ODTÜClass, a tar.gz archive file (named hw3.tar.gz) that contains all your source code files.
- The work you submit should be implemented by only you and genuine.
- We have zero tolerance policy for cheating. There is no teaming up! People involved in cheating will be punished according to the university regulations and will get 0. You can discuss design choices or language preferences, but sharing code between each other or submitting third party code as a whole is strictly forbidden. In case a match is found, this will be considered as cheating.