# 91.673 Advanced Database Systems Project

Chang Liu

chang_liu@student.uml.edu

01516178

November 25, 2015

## 1    Overview

For this project, I have configured the Hadoop environment as well as the Eclipse with Hadoop-eclipse-plugin, then setting up the namenode and other server on local machines, start the eclipse and create a Map-Reduce project in Eclipse with separate packages named 'ex', 'ex2', 'ex3' and 'ex4'.

For each package, I used different schemes to handle the source data file like 'city.txt', 'country.txt' and 'countrylanguage.txt'. These different tasks require some handling of mapping and reducing procedure in the Hadoop. By using different Map() and Reduce() function, I can parse the input text file and accumulate the lists that meet the specific requirement. I will give detailed description below about how it works.

I think I have finished all the projects that required by the description, but it seems that there is some mess unicode display problem, I think it's related with local machine's character set, and all of my handling of the string used the unicode version and should be no problem.

## 2    General Hadoop Project

For a general Hadoop project, There is actually four components we should do:

1) Map class that implement a map function. In this map function, the program receive each item from the raw data file like 'city.txt', and then we need to parse the input, add judgement or string handling, emit the necessary key-pair data as the source of reduce function.

2) Reduce class that implement a reduce function. Similarly as above, the reduce accepts the items that shared the same key, and then we add some manipulating of data like counting the number, concatenating the value string or other operations to generate the output.

3) A run function that configures the job, like the input or output of the Mapper/Reducer class, source data path.

4) A main function that starts the program and execute the run() function.

After these four components are done, I can start the running of the source program, with the hadoop server running in the backend. Then it will generate the results and save to the corresponding file we set in the source code.

# 3   Selection

For this first task, it's quite simple, I just need to parse the 'city.txt' file, get the each item in the lists, this is done by the java function 'String.split'. After we parse the data and get city name as well as the population, add a judgement for comparing the number with 300,000, only emits the items that are larger than the value, so the map function has finished.

In the reduce function, the program should accept all the item and just save the key value, since the key is the city name, and that's the only items that needs to be saved or output.

In the run function, set the configurations of input for Map and Reduce, here we both used the type 'Text' since it's the string.

# 4   Projection

For the second task, similar to the previous, the only different is that we emit the city name and the district as the output of map function, in the reduce function, we should connect all the string belonging to the same key, since we should output the city and its corresponding district.

# 5   Natural Joint

For natural joint, here there are two tables we need to parse, the 'country.txt' that stores the country code and its name and the 'countrylanguage.txt' that stores the country with its language.

First we judge the type of data, and know which table it is from, then emit only the English country with key-pair of the countrycode-NULL from language table. For the country table, emit all the countrycode-name pair for fetching the name when we know some English country.

In the reduce process, we have two types of data as the above key-pair shows sorted by the key, we only count the key that has more than 1 lists, because it means the it contains the country-NULL item and it's the English country, then by the value of the other list we know the name of the country, save it to the output. For only 1 list, it is from country table, and no English country is there, we just omit these records.

# 6  Aggregation

Similar to previous experiment, only use the 'city.txt' file and emit the district-city pair, then in the reduce function, by the key of district, count the number of the cities, output the district and its count, that forms the outputs of the total results.

# 7  Summary

I think these project are quite straight-forward. The main work lays on the setting up of the environment and getting a rough idea of how MapReduce program works, then by the knowledge from class, devise the correct map and reduce function to deal with data, we can get the correct outputs. I learned a lot from the books and it helps me a lot when starting the project and get to know how it works and how the program framework should be.

# 8  Reference

Install Guide: `https://hadoop.apache.org/docs/current/hadoop-project-dist/hadoop-common/SingleCluster.html`
Books: <Hadoop in Action>

# 9  Screenshot

All the running of program is similar so I just select the first output in eclipse, the final result is not output here, they're store under 'ex*/part-r-00000'. For the final result, please check the output folder in my submission, they have stored all the necessary items.