# Indian Institute of Technology Madras
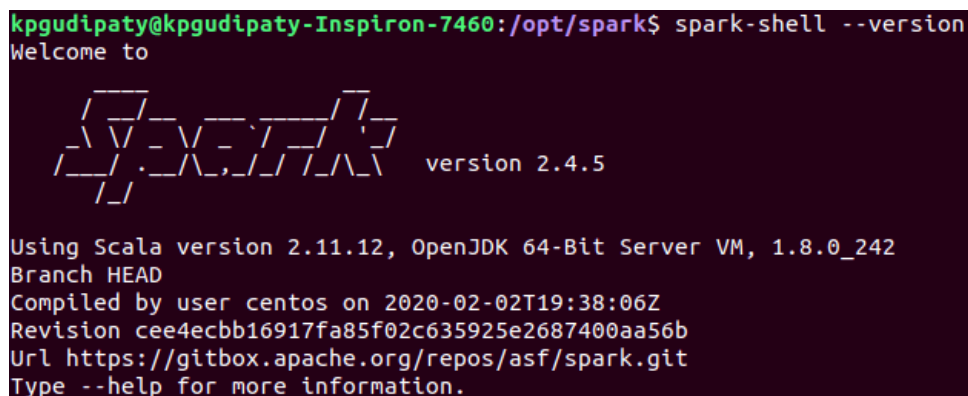## CS6847: Cloud Computing

Krishna Praneet - MM16B029

April 19, 2020

## Assignment-3 Report

### Setup

Apache spark has been setup on Ubuntu system. The *tar* file was downloaded from the Apache website and extracted to suitable folder. *Scala* was already installed, along with *Java*. Successful installation can be verified by running `spark-shell --version` which displays the version of the installation as well. Figure 2 shows the configuration file for spark which shows the setup of one master node i.e. the localhost.



Figure 1: Spark successful installation



Figure 2: Spark configuration - spark-defaults.conf.template file

The downloaded data sets are put onto the hdfs file system under **/data_mm16b029/** directory. This completes the initials setup for running the programs.



Figure 3: Datasets on hdfs system

# ALS

The ALS algorithm provided by spark was run on the given dataset. The fine tuning was done for the range of iterations from 1 to 20. The regularization parameter was run for values $0.001, 0.002, 0.005, 0.01, 0.02, 0.05, 0.1, 0.2, 0.5$. The corresponding *RMSE* are recorded in `ALS_out.txt` file. The minimised *RMSE* for current experimental ranges was approx. 0.90144 achieved at *9* iterations and a *0.1* value for the regularization parameter. The default train-test split of 0.8-0.2 was used.

These values were used to minimize *RMSE* wrt to the train-test split of data. It was done from range 0.1 to 0.9 split ratio of training data with 0.1 increments. The minimum was achieved at 0.9 with an *RMSE* of approx. 0.89861. These can be found in `ALS_out_train_split.txt` file.

# FP Growth

The FPGrowth example given by spark was used on both the datasets given. `FP_Part-1.csv` required just a bit of pre-pocessing. This preprocessing involved removing multiple instances of items in a transaction and to hold two-worded items together by joining them with an *underscore* between the words. For e.g. `french fries` became `french_fries`. The `FP_Part-2.csv` required additional data cleaning as well. The pre-processing and data cleaning used can be found in the `DataCleaner.py` file.

Five frequently occuring pairs are reported in `FP_out1.txt` and `FP_out2.txt`. The minimum support is 0.04 for `FP_Part-1.csv` and 0.02 for `FP_Part-2.csv`

*************************************************************