LINKS

---

TASKS:

1. Counting the frequency of words in the given input with MapReduce algorithm

   a. Create Java WordCount class

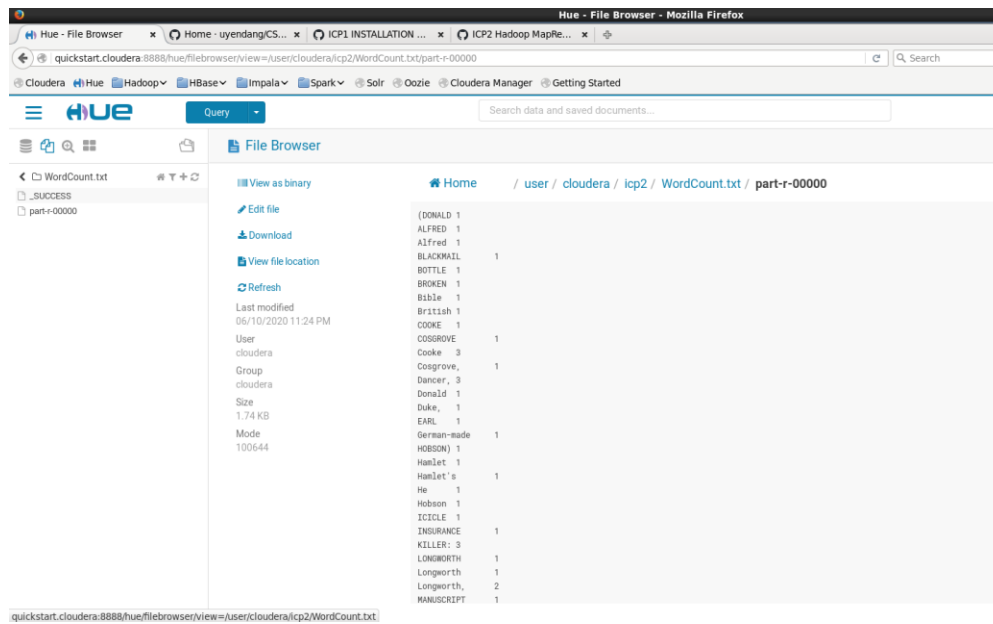   b. Add external libraries JARS

   c. Export the Jar

   d. Input data file in hdfs

   e. Run MapReduce Job

   hadoop jar /home/cloudera/WordCount.jar WordCount
   /user/cloudera/icp2/sample.txt /user/cloudera/icp2/WordCount.txt



Map-Reduce Job

Result

2. Counting the frequency of words in given text file that starts with letter 'a'

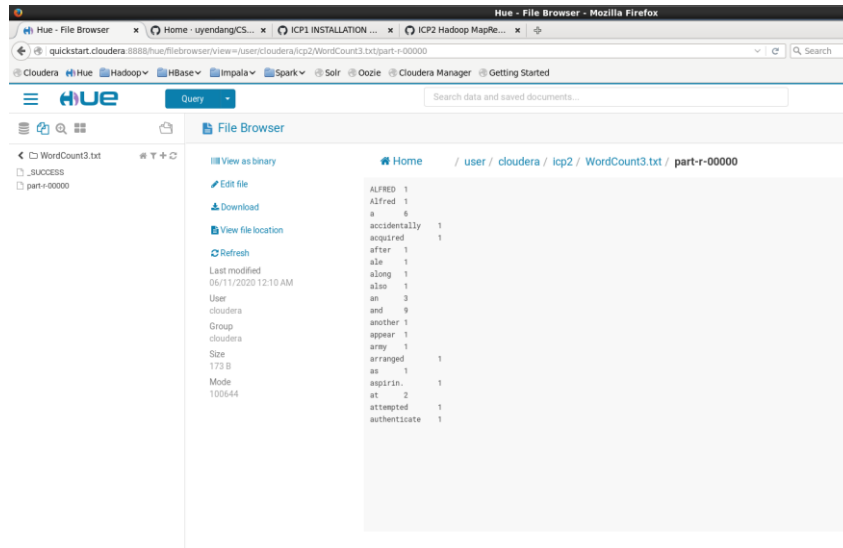      a. Add AWordCount function

      b. Export the Jar

      c. Input data file in hdfs

      d. Run MapReduce Job

      hadoop jar /home/cloudera/AWordCount.jar WordCount /user/cloudera/icp2/sample.txt /user/cloudera/icp2/WordCount3.txt

Result

---

REFERENCES: https://umkc.app.box.com/s/xk4jj0do3p7fa3swx6utcuazx3wny8j8