**MapReduce for word count problem on Hadoop:**

In MapReduce word count example, we find out the frequency of each word. Here, the role of Mapper is to map the keys to the existing values and the role of Reducer is to aggregate the keys of common values. So, everything is represented in the form of key-value pair.

**Example:**

Let's solve a word count problem using MapReduce on Hadoop.

**Step 1:** Open Cloudera Quickstart VM.



**Step 2:** Create a .txt data file inside **/home/cloudera** directory that will be passed as an input to MapReduce program. For simplicity purpose, we name it as **word_count_data.txt.**

**Step 3:** Create `mapper.py` and `reducer.py` files inside `/home/cloudera` directory.

## mapper.py

```
#!/usr/bin/python
# import sys because we need to read and write data to STDIN and
STDOUT

import sys
# reading entire line from STDIN (standard input)

for line in sys.stdin:
    # to remove leading and trailing whitespace
    line = line.strip()

    # split the line into words
    words = line.split()

    # we are looping over the words array and printing the word
    # with the count of 1 to the STDOUT

    for word in words:

        # write the results to STDOUT (standard output);
        # what we output here will be the input for the
        # Reduce step, i.e., the input for reducer.py
        print("%s\t%s" % (word, 1))
```

## reducer.py

```
#!/usr/bin/python

from operator import itemgetter
import sys

current_word = None
current_count = 0
word = None

# read the entire line from STDIN
for line in sys.stdin:
    # remove leading and trailing whitespace
    line = line.strip()

    # spliting the data on the basis of tab we have provided in mapper.py
    word, count = line.split('\t', 1)

    # convert count (currently a string) to int
    try:
        count = int(count)
    except ValueError:
        # count was not a number, so silently
        # ignore/discard this line
        continue

    # this if-switch only works because Hadoop sorts map output
    # by key (here: word) before it is passed to the reducer
```

```
        if current_word == word:
            current_count += count

        else:
            if current_word:

                # write result to STDOUT
                print("%s\t%s" % (current_word, current_count))

            current_count = count
            current_word = word


# do not forget to output the last word if needed

if current_word == word:
    print("%s\t%s" % (current_word, current_count))
```
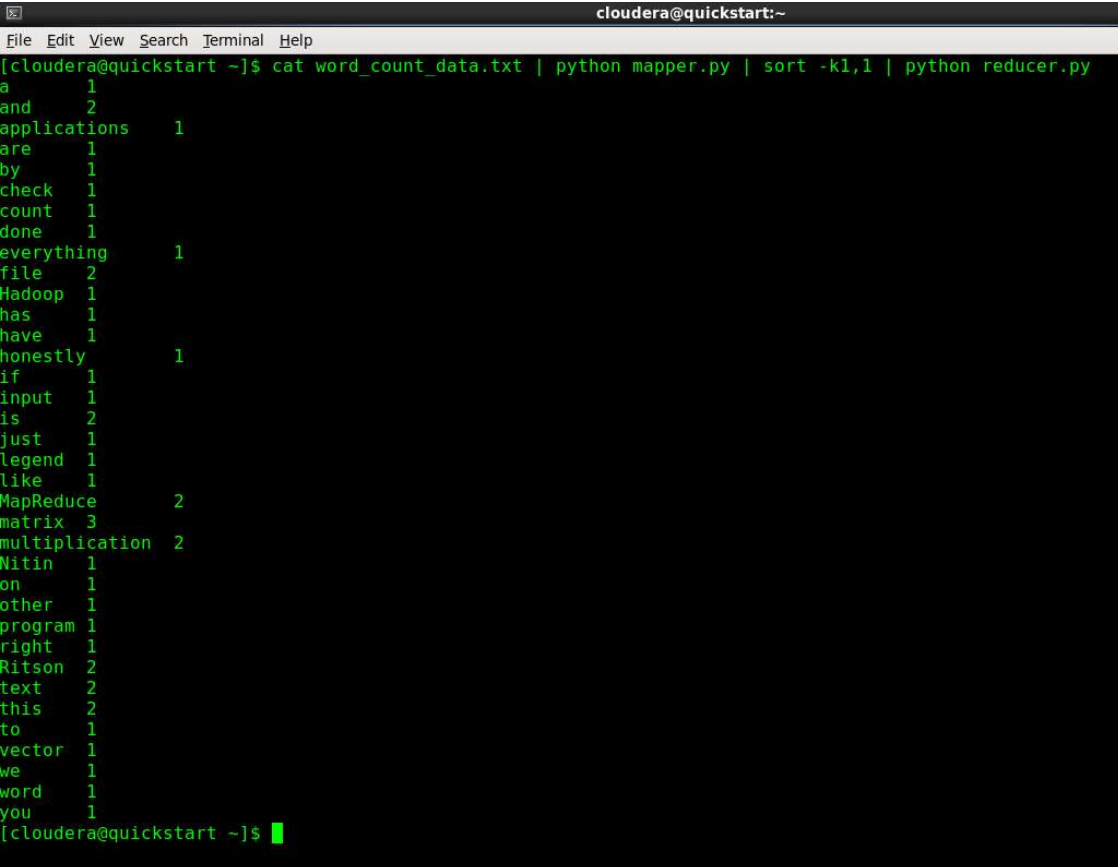
**Step 4:** Test the MapReduce program(s) locally to check if everything works properly before running on Hadoop.

**cat word_count_data.txt | python mapper.py | sort -k1,1 | python reducer.py**



For the above example, the output obtained is exactly the same as expected.

If you see all the words correctly mapped, sorted and reduced to their respective counts, then your program is good to be tested on Hadoop.

**Step 5:** Configure Hadoop services and settings.

Now, we need to configure certain settings on Hadoop before we run the MapReduce program for word count.

**5a: Login to Cloudera Quickstart.**

Open browser on Cloudera Quickstart VM and open quickstart.cloudera:7180/cmf/login. Login by entering the credentials as cloudera for both, username and password.



**Note:** If you see the error "Unable to connect" while logging in to quickstart.cloudera:7180/cmf/login, try restarting the CDH services.



Restart CDH services by typing the following command:

sudo /home/cloudera/cloudera-manager --express --force

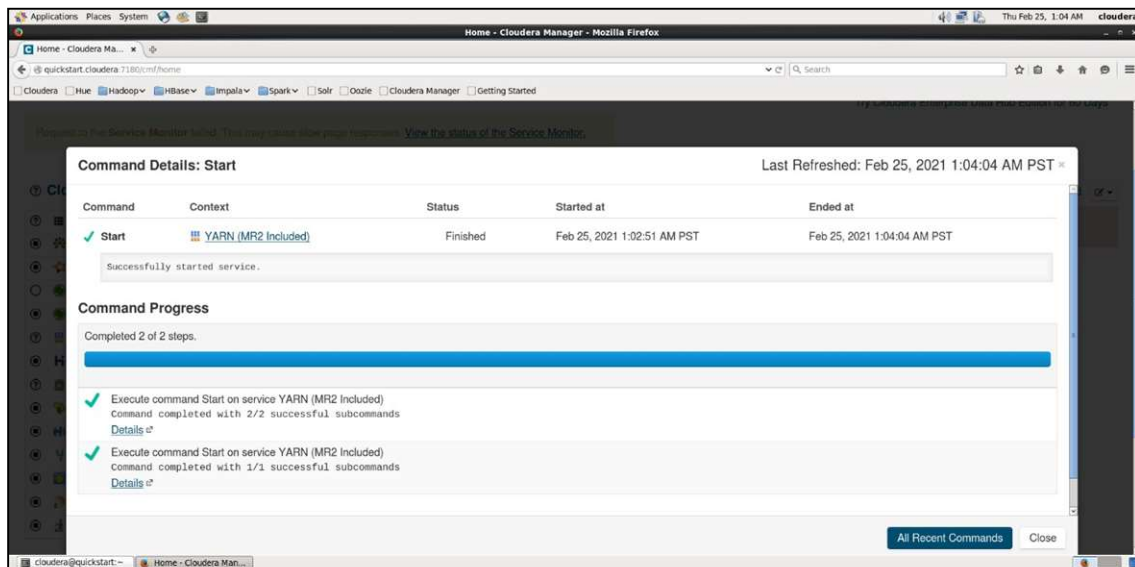## 5b: Start HDFS and YARN services.

Click the dropdown arrow and choose Start option for HDFS and YARN services.



You'll see the following if both; HDFS and YARN services are started successfully.



HDFS service started successfully.

YARN service started successfully.

## Step 6: Create a directory on HDFS

Now, we create a directory named `word_count_map_reduce` on HDFS where our input data and its resulting output would be stored. Use the following command for it.
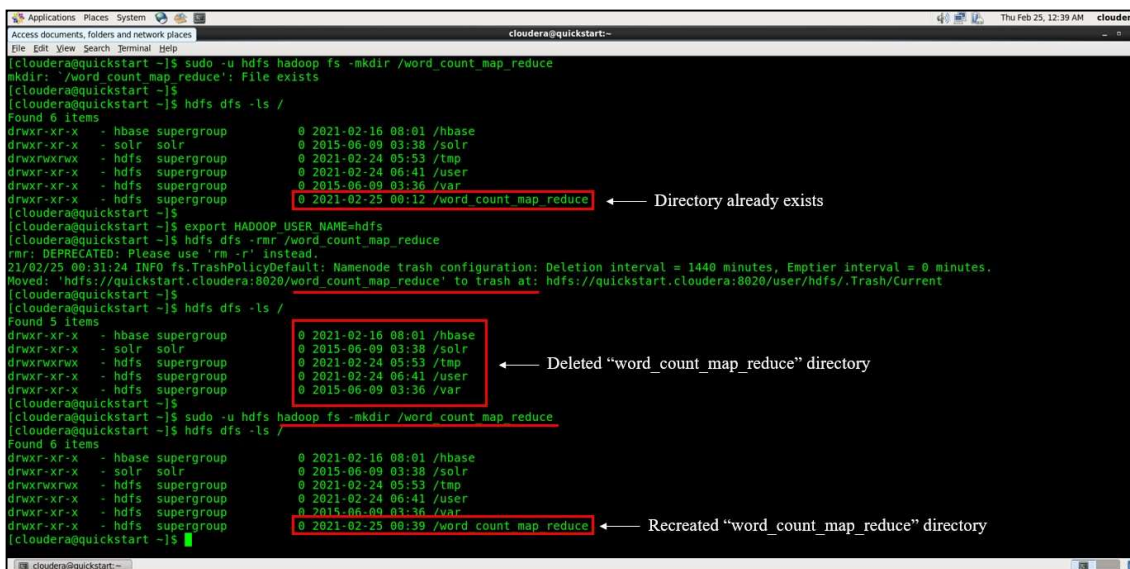
sudo -u hdfs hadoop fs -mkdir /word_count_map_reduce

> **Note:** If the directory already exists, then either create a directory with new name or delete the existing directory using the following command.
>
> export HADOOP_USER_NAME=hdfs
> hdfs dfs -rmr /word_count_map_reduce

List HDFS directory items using the following command.
hdfs dfs -ls /

**Step 7: Move input data file to HDFS.**

Copy the `word_count_data.txt` file to `word_count_map_reduce` directory on HDFS using the following command.

sudo -u hdfs hadoop fs -put /home/cloudera/word_count_data.txt /word_count_map_reduce

Check if file was copied successfully to the desired location.
hdfs dfs -ls /word_count_map_reduce



**Step 8: Download hadoop-streaming JAR 2.7.3.**

Open browser and go to https://jar-download.com/artifacts/org.apache.hadoop/hadoop-streaming/2.7.3/source-code and download hadoop-streaming JAR 2.7.3 file.



Once the file is downloaded, unzip it inside `/home/cloudera` directory.

Double-check if the JAR file was unzipped successfully and is present inside `/home/cloudera` directory.

ls



**Step 9: Configure permissions to run MapReduce on Hadoop.**

We're almost ready to run our MapReduce job on Hadoop but before that, we need to give permission to read, write and execute the Mapper and Reducer programs on Hadoop.

We also need to provide permission for the default user (cloudera) to write the output file inside HDFS.

Run the following commands to do so:

chmod 777 mapper.py reducer.py
sudo -u hdfs hadoop fs -chown cloudera /word_count_map_reduce



**Step 10: Run MapReduce on Hadoop.**

We're at the ultimate step of this program. Run the MapReduce job on Hadoop using the following command.

hadoop jar /home/cloudera/hadoop-streaming-2.7.3.jar \
> -input /word_count_map_reduce/word_count_data.txt \
> -output /word_count_map_reduce/output \
> -mapper /home/cloudera/mapper.py \
> - reducer /home/cloudera/reducer.py

If you see the output on terminal as shown in above two images, then the MapReduce job was executed successfully.

**Step 11: Read the MapReduce output.**

Now, finally run the following command to read the output of MapReduce for word count of the input data file you had created.

hdfs dfs -cat /word_count_map_reduce/output/part-00000

```
[cloudera@quickstart ~]$ hdfs dfs -cat /word_count_map_reduce/output/part-00000
Hadoop  1
MapReduce       2
Nitin   1
Ritson  2
a       1
and     2
applications    1
are     1
by      1
check   1
count   1
done    1
everything      1
file    2
has     1
have    1
honestly        1
if      1
input   1
is      2
just    1
legend  1
like    1
matrix  3
multiplication  2
on      1
other   1
program 1
right   1
text    2
this    2
to      1
vector  1
we      1
word    1
you     1
[cloudera@quickstart ~]$
```

Congratulations, the output for MapReduce on Hadoop is obtained exactly as expected. All the words in the input data file have been mapped, sorted and reduced to their respective counts.