# BDA Lab 6 Tasks

**Notes:**
- Keep in mind that output directory name will be your roll number else no marks would be awarded
- You are required to run your WordCount on the text file provided with this lab named "TheseWords.txt" which contains all the words in the tutorial.
- Must change "i19-1207" with your roll number at all places.
- Must change "SaadNaeem" at all places with your name.

# Task-1

Creating your first WordCount MapReduce program in python. Use the pdf provided with this lab or follow the tutorial at:

- https://www.michael-noll.com/tutorials/writing-an-hadoop-mapreduce-program-in-python/
- You need to make changes in the code according to your problem statement and fix any bugs that are there in the code (Because there are few points that need to be fixed in the code).

**Extra Help related to Task-1:**

**Information** about hadoop streaming can be found at:
https://hadoop.apache.org/docs/r1.2.1/streaming.html

**Hint-1:** you need to get the streaming api jar file for your hadoop version
**Hint-2:** I have put the streaming api jar file in my hadoop installation at:
/usr/local/Cellar/hadoop/3.2.1/libexec/share/hadoop/mapreduce/hadoop-streaming-3.2.1.jar
so when i run the MapReduce program in python i give the above path to "hadoop jar /usr/local/Cellar/hadoop/3.2.1/libexec/share/hadoop/mapreduce/hadoop-streaming-3.2.1.jar ................."

**Debugging Tips:**
No Spaces in file path
echo "foo foo quux labs foo bar quux" | python mapper.py
echo "foo foo quux labs foo bar quux" | python mapper.py | sort -k1,1 | python reducer.py

# Task-2 (20% Marks)



Requirement 1:i19-1207-image-1.png

You are required to submit i19-1207-Image-1.png where i19-1207 is your roll number that shows that your installation is up and running and shows your unique block id and cluster id.

# Task-3 (20% Marks)



Requirement 2:i19-1207-image2.png

Submit "i19-1207-Image-2.png" that shows where were you when you executed the command to run map reduce in python and what was the command itself.
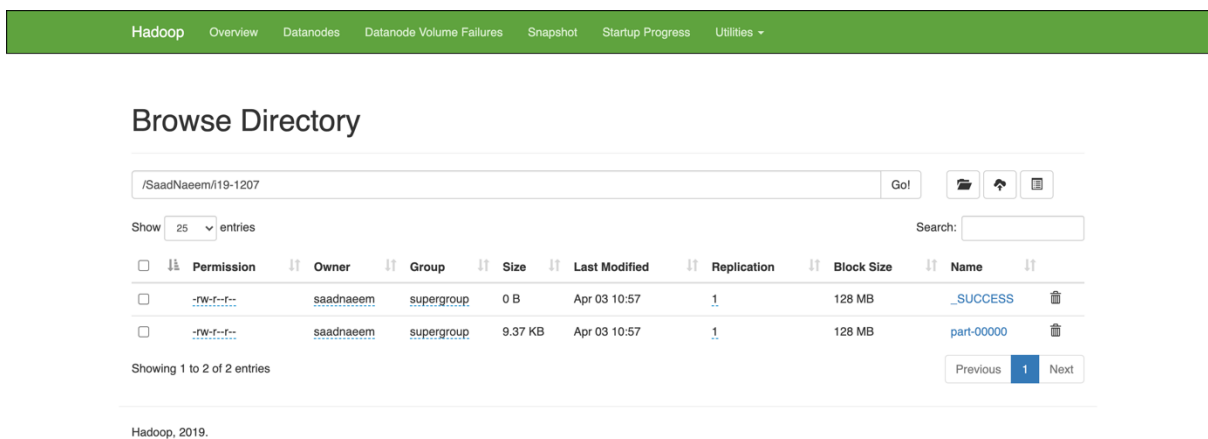
# Task-4 (20% Marks)

Submit "i19-1207-Image-3.png" that shows the user name of the system, the success of map reduce on terminal, and if you notice carefully the streaming API shows the output directory where the results are stored.
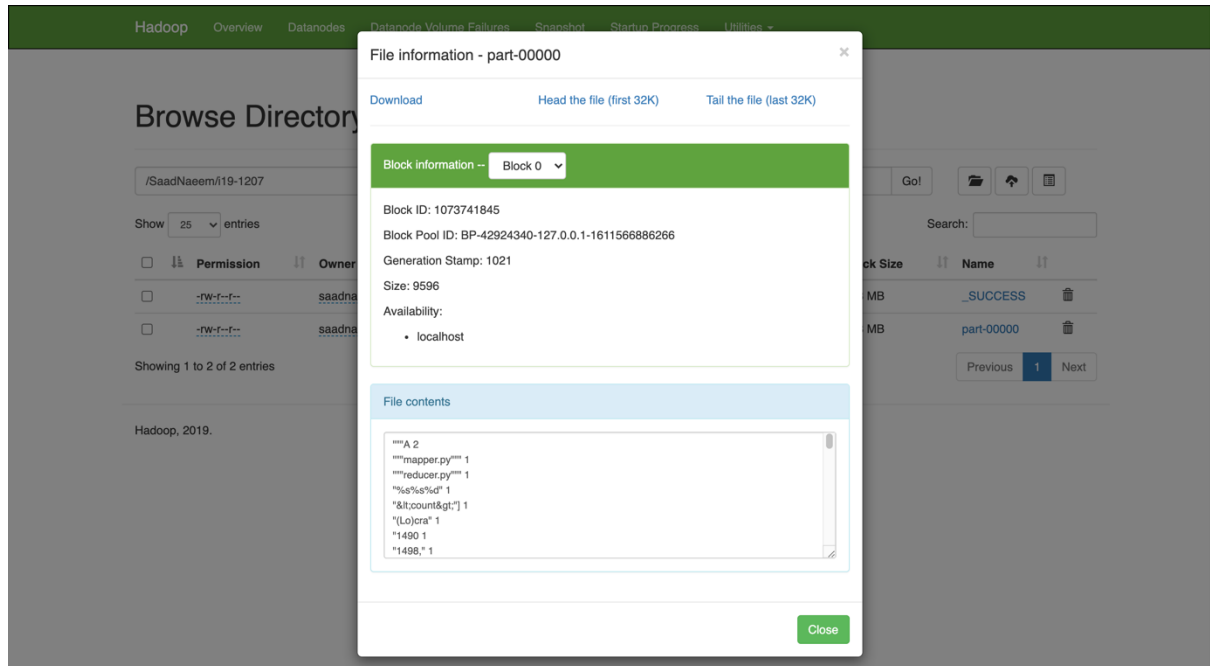
# Task-5 (20% Marks)



*Requirement 4:i19-1207-image-4.png*

i19-1207-image-4.png shows that your Datanode is working properly and the MapReduce success in terminal has indeed reflected here.

# Task-6 (20% Marks)



*Requirement 5: i19-1207-image-5*

Moment of truth " i19-1207-image-5.png" shows the actual results and verifies your MapReduce logic.

# Bonus Task

**Task:** There is a redundancy in Mapper/Reducer logic you need to find it and fix it. This requires and understanding of what "Mapper" does and what "Reducer" does. And what "streaming.jar" is doing.

**Requirement: You can only do this task if you have done all the mandatory tasks &** Only submit your "mapper.py" and "reducer.py" if you completed this task.