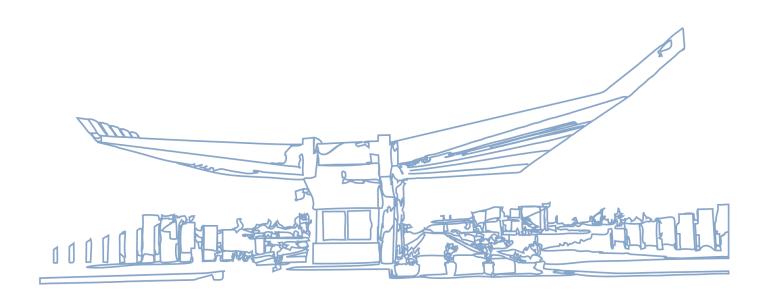


CEN 571 – Data Mining

Assignment 01



PREPARED:

Baftjar TABAKU

12.04.2020

Epoka University Tirana, ALBANIA **ACCEPTED:**

Prof.Dr. Arben Asllani

Assignment tasks and notes

- 1. Install VirtualBox and Setup Cloudera QuickStart VM (50)
- 2. Add Ubuntu in VirtualBox and Setup a One Node Hadoop Cluster (50)
- 3. (optional/extra credit): Setup a Four Node Hadoop Cluster using AWS (20)

For each of the above, submit a word document that includes a list of commands accompanied by screenshots. Each document must look like a tutorial.

Hadoop and Cloudera

The core of Apache Hadoop consists of a storage part, known as Hadoop Distributed File System (HDFS), and a processing part which is a MapReduce programming model. Hadoop splits files into large blocks and distributes them across nodes in a cluster. It then transfers packaged code into nodes to process the data in parallel. This approach takes advantage of data locality, where nodes manipulate the data they have access to. This allows the dataset to be processed faster and more efficiently than it would be in a more conventional supercomputer architecture that relies on a parallel file system where computation and data are distributed via high-speed networking. Hadoop's ability to process and store different types of data makes it a particularly good fit for big data environments. They typically involve not only large amounts of data, but also a mix of structured transaction data and semi structured and unstructured information, such as internet clickstream records, web server and mobile application logs, social media posts, customer emails and sensor data from the internet of things.

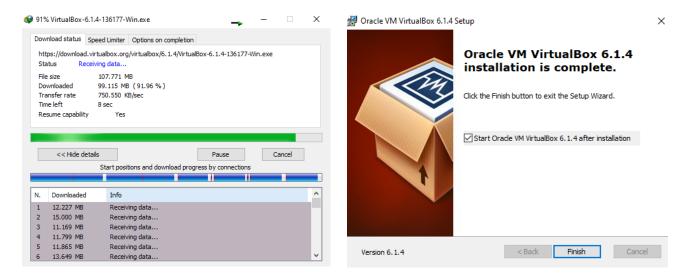
Ref: https://searchdatamanagement.techtarget.com

The Cloudera delivers the modern platform for machine learning and advanced analytics built on the latest open source technologies. The world's leading organizations trust Cloudera to help solve their most challenging business problems by efficiently capturing, storing, processing and analyzing vast amounts of data.

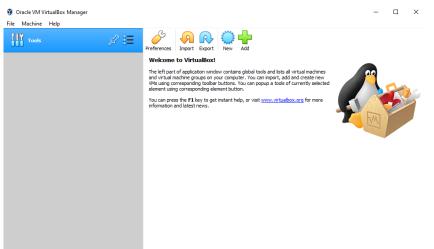
Ref: https://cloudian.com

Installing VirtualBox and Setting up Cloudera QuickStart VM (50)

I downloaded VirtualBox and set up, as the screenshots shows,

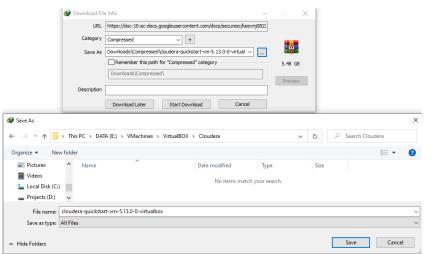


After installing



Next step was to download the Cloudera QuickStart VM, provided in the

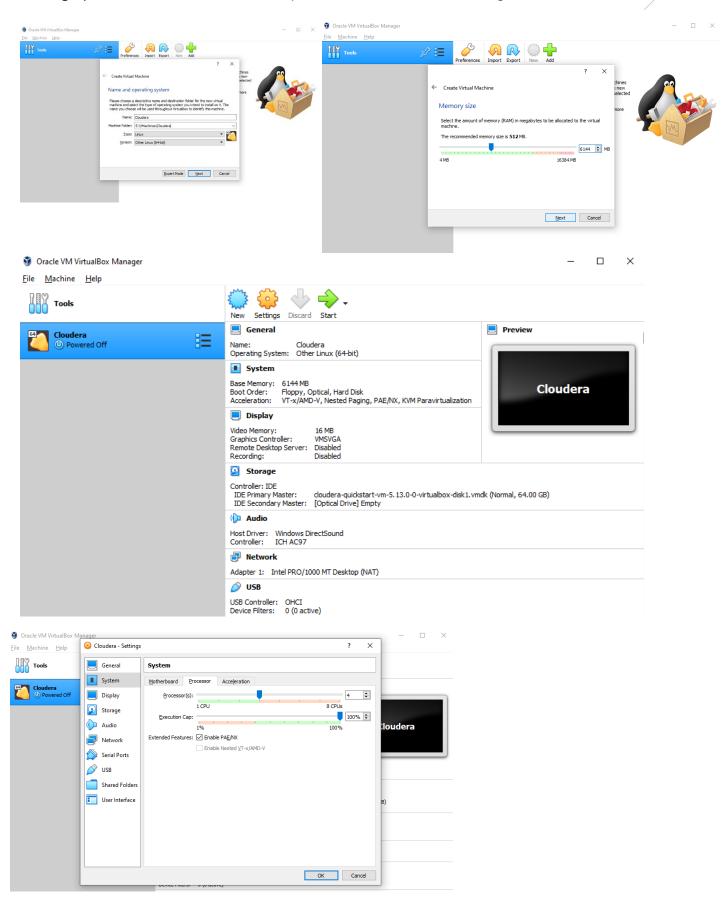
link: https://drive.google.com/file/d/1IK-ZfiKfKaY8LrU0oDk--HH4ijgVpFrB/view



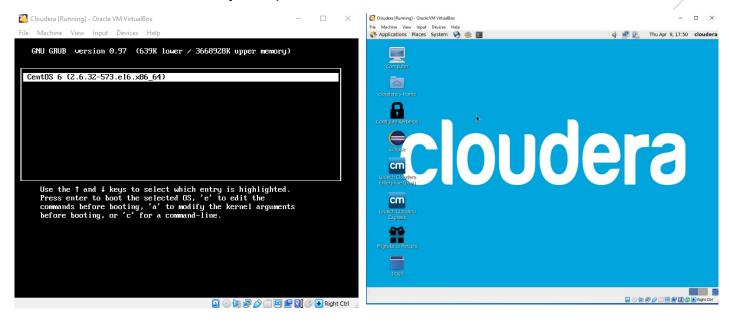
That took approximately 1h 15min.

After downloading it, In set it up on VirtualBox.

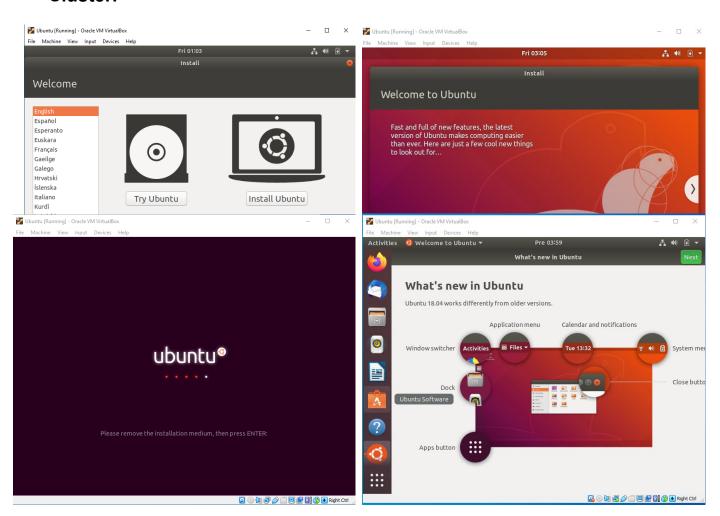
Setting up the Cloudera QuickStart VM process is shown in the following screenshots,



The Cloudera successfully set up on VirtualBox.

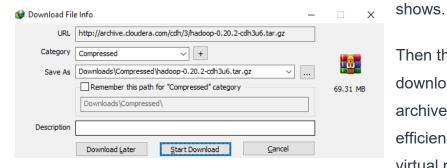


Adding Ubuntu in VirtualBox and setting up a One Node Hadoop Cluster.



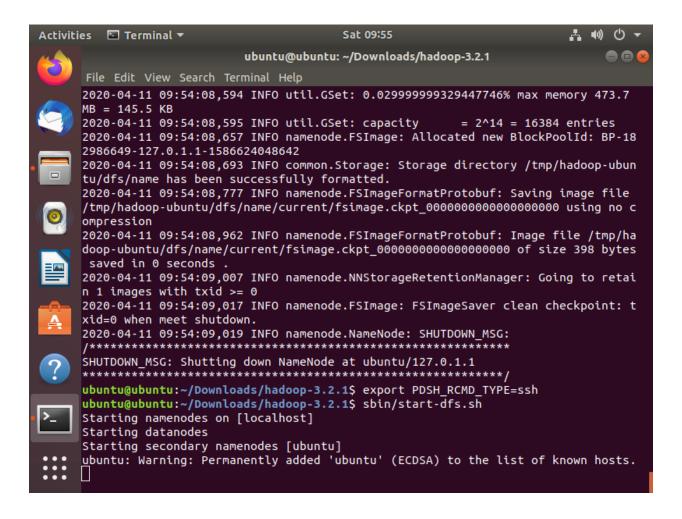
Installing ubuntu in VirtualBox

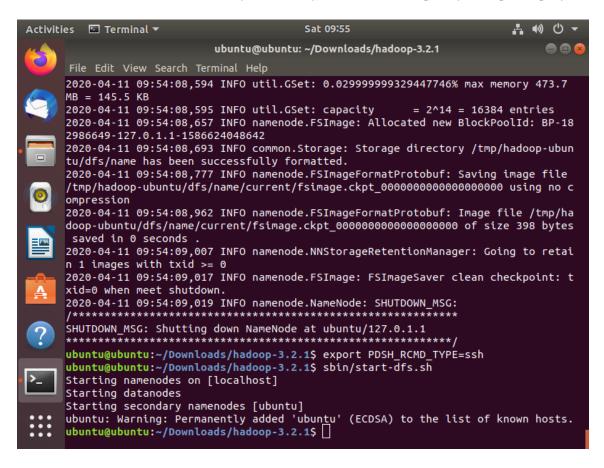
According to the provided manuals, I Installed the Hadoop as following screenshot

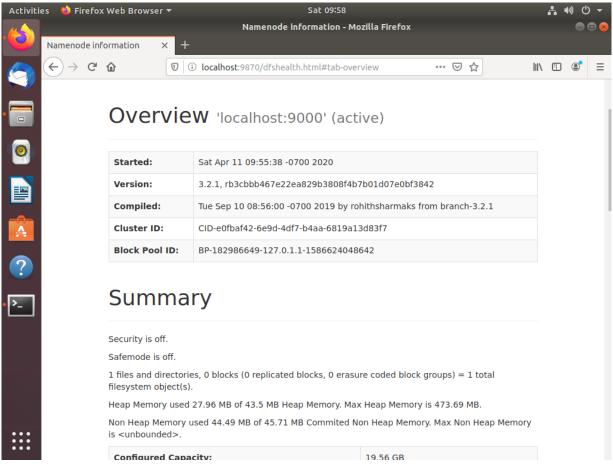


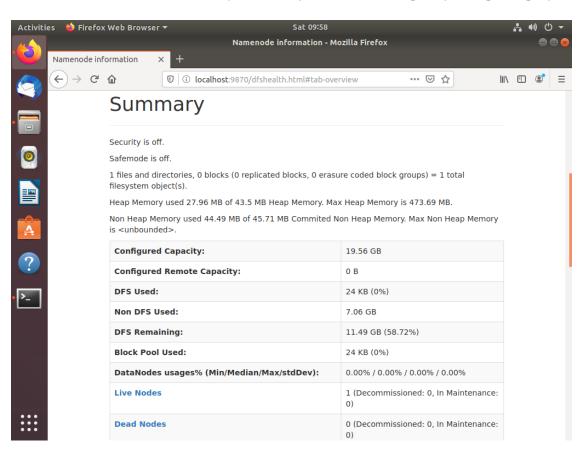
Then the Hadoop setup,
downloaded from the cloudera
archive, for more resource
efficiency I added another kind of
virtual machine like VMware, a

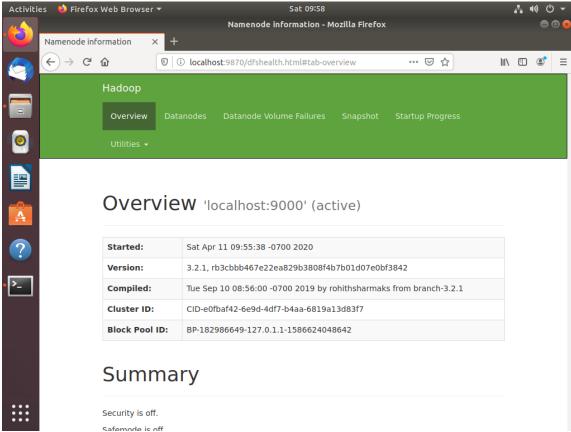
light-weight version to use it with ubuntu installed and setting up Hadoop cluster single node.











Note: Screenshots shows some fractions at the cloudera and Hadoop setup installation.

Setting up a Four Node Hadoop Cluster using AWS.

The idea was clear on what to do, and I would do it as was shown in the 7 parts with videos, the problem was that at the moment I don't have a credit card to sign up and creating a AWS account.