# Practical 1 :
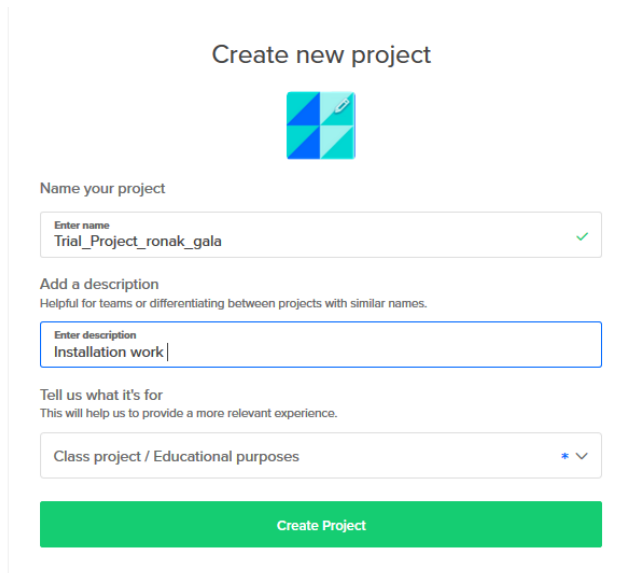
Aim : Install, configure and run Hadoop and HDFS ad explore HDFS.

## Phase 1 : Creating Droplet
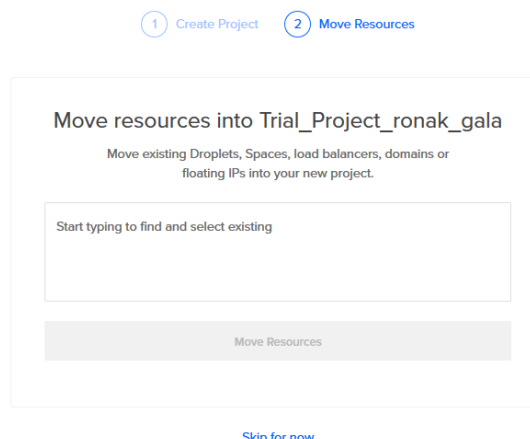
Step 1 : Sign up on https://www.digitalocean.com/
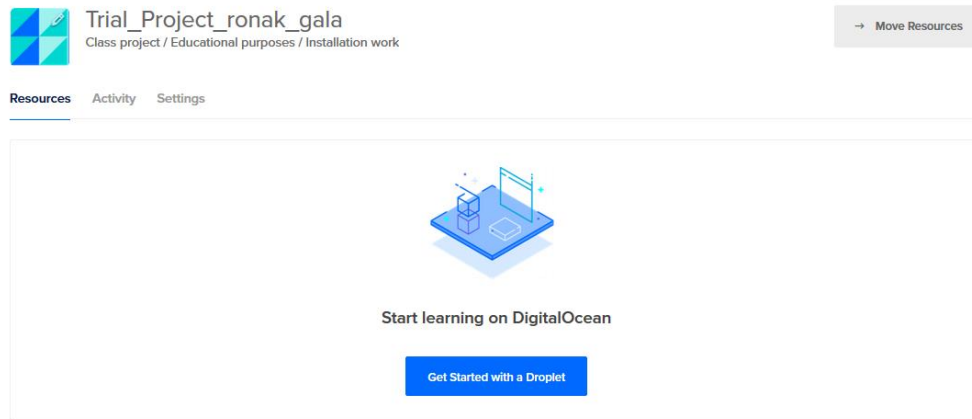
Step 2 : Create a new project.

Step 3 : You can skip the move resources page.

Step 4 : Click ' Get Started with a Droplet'.



Step 5 : Select the following configurations :-

Choose an Image – Ubuntu

Choose datacentre region – Bangalore

Authentication – OTP (you can also select SSH but you might need to set it up)
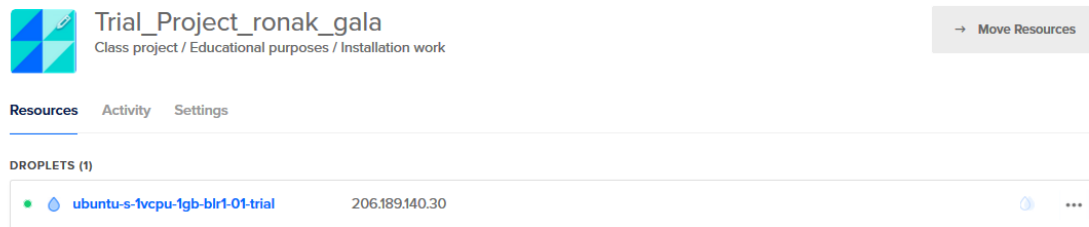
How many droplets – 1

Choose a new hostname or go with the provided one

Then select Create Droplet.

Phase 2: Installing JDK

Step 6 : Click on the droplet after it has been created.



Step 7 : Click on console.



A new window will open with the console.

Step 8 : Login in to your registered email id and look for a mail with username and password of your new droplet. Insert the credentials and change the password as per your requirement.

Note : For password the cursor doesn't move, i.e. it won't show the password length just like ubuntu and other Linux flavours.

Step 9 : Check for updates

Command : sudo apt update

Step 10 : Install OpenJDK

Command : sudo apt install default-jdk

It will ask for permission as it will take up 727mb of disk, so allow it accordingly.

Step 11 : Check for the java version.

Command : java -version

```
root@ubuntu-s-1vcpu-1gb-blr1-01-trial:~# java -version
openjdk version "11.0.6" 2020-01-14
OpenJDK Runtime Environment (build 11.0.6+10-post-Ubuntu-1ubuntu118.04.1)
OpenJDK 64-Bit Server VM (build 11.0.6+10-post-Ubuntu-1ubuntu118.04.1, mixed mod
e, sharing)
root@ubuntu-s-1vcpu-1gb-blr1-01-trial:~# _
```

Phase 3: Installing Hadoop

Step 12 : Go to https://hadoop.apache.org/releases.html to get the stable release.

Step 13 : Navigate to the binary you wish to install. In this guide we will be installing 3.1.3



Step 14 : Click on binary, from your desired version's row.

Step 15 : Copy the http url.

Step 16 : On the droplet's console we will use 'wget' to fetch it.

Command    :    wget    https://downloads.apache.org/hadoop/common/hadoop-3.1.3/hadoop-3.1.3.tar.gz


Step 17 : We'll use the tar command with the -x flag to extract, -z to uncompress,  -v for verbose output, and -f to specify that we're extracting from a file.

Command : tar -xzvf hadoop–3.1.3.tar.gz


Step 18 : We'll move the extracted files into /usr/local

Command : sudo mv hadoop-3.1.3 /usr/local/Hadoop


Phase 4: Configuring Hadoop's Java Home

Step 19 : We'll find java's default path

Command : readlink -f /usr/bin/java | sed "s:bin/java::"

Output : /usr/lib/jvm/java-11-openjdk-amd64/


Step 20 : You can copy this output to set Hadoop's Java home to this specific version, which ensures that if the default Java changes, this value will not. Alternatively, you can use the readlink command dynamically in the file so that Hadoop will automatically use whatever Java version is set as the system default.

To begin, open hadoop-env.sh:

Command : sudo nano /usr/local/hadoop/etc/hadoop/hadoop-env.sh

Step 21 : Move the cursor till the last line and then to set the value to Java_Home write this command:

Option 1 : Set a Static Value

Command : #export JAVA_HOME=${JAVA_HOME} ( this line is a comment)

export JAVA_HOME=/usr/lib/jvm/java-11-openjdk-amd64/


Option 2 :Use Readlink to Set the Value Dynamically

Command : #export JAVA_HOME=${JAVA_HOME} (this line is a comment)

export JAVA_HOME=$(readlink -f /usr/bin/java | sed "s:bin/java::")
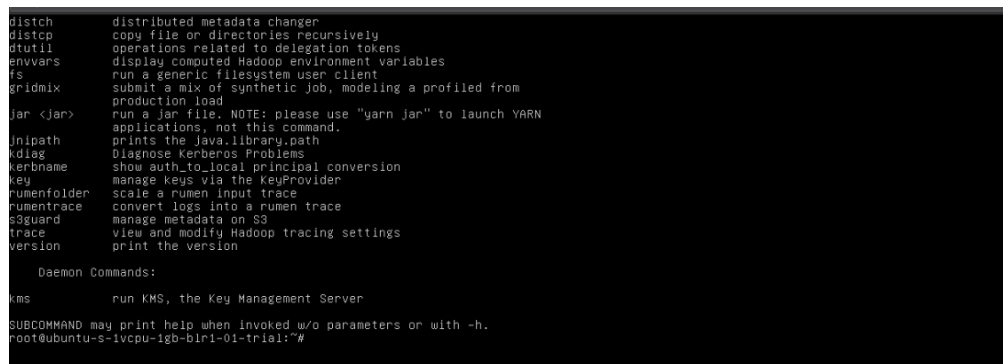

Press ctrl + X to exit file

It will ask whether you want to save the modified buffer, answer with YES

It will ask for a file name to write press ENTER


Phase 5: Running Hadoop

Step 22: We will run Hadoop

Command : /usr/local/hadoop/bin/hadoop



The help means we've successfully configured Hadoop to run in stand-alone mode.

Step 23 : We'll ensure that it is functioning properly by running the example MapReduce program it ships with. To do so, create a directory called input in our home directory and copy Hadoop's configuration files into it to use those files as our data.

Command 1:  mkdir ~/input

Command 2: cp /usr/local/hadoop/etc/hadoop/*.xml ~/input

```
root@ubuntu-s-1vcpu-1gb-blr1-01-trial:~# mkdir ~/input
root@ubuntu-s-1vcpu-1gb-blr1-01-trial:~# cp /usr/local/hadoop/etc/hadoop/*.xml ~
/input
root@ubuntu-s-1vcpu-1gb-blr1-01-trial:~#
```

Step 24 : We can use the following command to run the MapReduce hadoop-mapreduce-examples program, a Java archive with several options. We'll invoke its grep program, one of the many examples included in hadoop-mapreduce-examples, followed by the input directory, input and the output directory grep_example. The MapReduce grep program will count the matches of a literal word or regular expression. Finally, we'll supply the regular expression allowed[.]* to find occurrences of the word allowed within or at the end of a declarative sentence. The expression is case-sensitive, so we wouldn't find the word if it were capitalized at the beginning of a sentence:

Command            :            /usr/local/hadoop/bin/hadoop            jar /usr/local/hadoop/share/hadoop/mapreduce/hadoop-mapreduce-examples-3.1.3.jar grep ~/input ~/grep_example 'allowed[.]*'

Note: If the output directory already exists, the program will fail, and rather than seeing the summary, the output will look something like:

Output

. . .

    at java.base/java.lang.reflect.Method.invoke(Method.java:564)

    at org.apache.hadoop.util.RunJar.run(RunJar.java:244)

    at org.apache.hadoop.util.RunJar.main(RunJar.java:158)

Step 25 : Results are stored in the output directory and can be checked by running cat on the output directory:

Command : cat ~/grep_example/*

```
root@ubuntu-s-1vcpu-1gb-blr1-01-trial:~# cat ~/grep_example/*
21      allowed.
1       allowed
root@ubuntu-s-1vcpu-1gb-blr1-01-trial:~# _
```

The MapReduce task found 21 occurrences of the word allowed followed by a period and one occurrence where it was not. Running the example program has verified that our stand-alone installation is working properly and that non-privileged users on the system can run Hadoop for exploration or debugging.