

<b>Ex. No: 1</b>	<b>Basic Linux Commands</b>
<b>16.7.2024</b>	

### **Aim:**

To execute various basic Linux commands

### **Program:**

1. Display information about current directory - `ls`
2. Display the current working directory - `pwd`
3. Create a new directory - `mkdir SEMESTER-7`
4. Navigate between different folders - `cd SEMESTER-6`
5. Remove empty directories find . -type d -empty -delete
6. Copy files from one directory to same directory - `cp matmul.cpp ./copied.cpp`
7. Copy files from one directory to another directory - `cp matmul.cpp ../SEMESTER-7`
8. Rename a filename to another name - `mv SEMESTER-7 semester-7`
9. Move a file from one directory to another - `mv copied.cpp trial/`
10. Delete individual files from a directory - `rm -rf matmul.cpp`
11. Delete an unempty directory - `rm -rf trial/`
12. Get basic information about the OS - `lscpu`
13. Find a file in the directory - `find *.cpp`
14. Create empty files - `touch empty.txt`
15. Display file contents on terminal - `cat empty.txt`
16. Clear Terminal - `clear`
17. Display the processes in terminal - `ps -A`
18. Access manual for all Linux commands - `help`
19. Search for a specific string in an output - `ls -l | grep "s/. *w://p"`
20. 20. Display active processes on the terminal - `ps`
21. 21. Download files from the internet - `wget <url>`
22. Create or update passwords for existing users - `passwd`
23. View exact location of any tool/software installed - `which bash`
24. Check the details of the file system - `df -Th`
25. Check the lines, word count and characters in a file using different options :  
For characters - `wc -c empty.txt`  
For words - `wc -w empty.txt`  
For lines - `wc -l empty.txt`

Output:

```
rheaubuntu@LAPTOP-RB9PEJ x + v - □ ×  
rheaubuntu@LAPTOP-RB9PEJTU:~$ ls  
SEMESTER-6  
rheaubuntu@LAPTOP-RB9PEJTU:~$ pwd  
/home/rheaubuntu  
rheaubuntu@LAPTOP-RB9PEJTU:~$ mkdir SEMESTER-7  
rheaubuntu@LAPTOP-RB9PEJTU:~$ ls  
SEMESTER-6 SEMESTER-7  
rheaubuntu@LAPTOP-RB9PEJTU:~$ cd SEMESTER-6  
rheaubuntu@LAPTOP-RB9PEJTU:~/SEMESTER-6$ cd ../SEMESTER-7  
rheaubuntu@LAPTOP-RB9PEJTU:~/SEMESTER-7$ ls  
rheaubuntu@LAPTOP-RB9PEJTU:~/SEMESTER-7$ find . -type d -empty  
.  
rheaubuntu@LAPTOP-RB9PEJTU:~/SEMESTER-7$ find . -type d -empty -delete  
e  
rheaubuntu@LAPTOP-RB9PEJTU:~/SEMESTER-7$ cd ../SEMESTER-6  
rheaubuntu@LAPTOP-RB9PEJTU:~/SEMESTER-6$ ls  
a.exe      matmul.cpp  matmul_output.txt  pointer.cpp  rand_mat.exe  
hello.cpp  matmul.exe  openmp             rand_mat.cpp  
rheaubuntu@LAPTOP-RB9PEJTU:~/SEMESTER-6$ cp matmul.cpp ./  
cp: 'matmul.cpp' and './matmul.cpp' are the same file  
rheaubuntu@LAPTOP-RB9PEJTU:~/SEMESTER-6$ cp matmul.cpp ./copied.cpp  
rheaubuntu@LAPTOP-RB9PEJTU:~/SEMESTER-6$ ls  
a.exe      hello.cpp  matmul.exe  openmp      rand_mat.cpp  
copied.cpp matmul.cpp matmul_output.txt  pointer.cpp  rand_mat.exe  
rheaubuntu@LAPTOP-RB9PEJTU:~/SEMESTER-6$ cp matmul.cpp ../SEMESTER-7  
rheaubuntu@LAPTOP-RB9PEJTU:~/SEMESTER-6$ cd ../SEMESTER-7  
rheaubuntu@LAPTOP-RB9PEJTU:~/SEMESTER-7$ ls  
matmul.cpp  
rheaubuntu@LAPTOP-RB9PEJTU:~/SEMESTER-7$ |
```

```
rheaubuntu@LAPTOP-RB9PEJ x + v - □ ×  
rheaubuntu@LAPTOP-RB9PEJTU:~/SEMESTER-7$ cd ..  
rheaubuntu@LAPTOP-RB9PEJTU:~$ mv SEMESTER-7 semester-7  
rheaubuntu@LAPTOP-RB9PEJTU:~$ ls  
SEMESTER-6 semester-7  
rheaubuntu@LAPTOP-RB9PEJTU:~$ |
```

```
rheaubuntu@LAPTOP-RB9PEJ x + v - □ ×  
rheaubuntu@LAPTOP-RB9PEJTU:~/semester-7$ ls  
matmul.cpp  
rheaubuntu@LAPTOP-RB9PEJTU:~/semester-7$ rm -rf matmul.cpp  
rheaubuntu@LAPTOP-RB9PEJTU:~/semester-7$ ls  
rheaubuntu@LAPTOP-RB9PEJTU:~/semester-7$ |
```

```
rheaubuntu@LAPTOP-RB9PEJ x + v - □ ×  
rheaubuntu@LAPTOP-RB9PEJTU:~/SEMESTER-6$ ls  
a.exe      matmul.cpp  matmul_output.txt  pointer.cpp  rand_mat.exe  
hello.cpp  matmul.exe  openmp             rand_mat.cpp  trial  
rheaubuntu@LAPTOP-RB9PEJTU:~/SEMESTER-6$ rm -rf trial/  
rheaubuntu@LAPTOP-RB9PEJTU:~/SEMESTER-6$ ls  
a.exe      matmul.cpp  matmul_output.txt  pointer.cpp  rand_mat.exe  
hello.cpp  matmul.exe  openmp             rand_mat.cpp  
rheaubuntu@LAPTOP-RB9PEJTU:~/SEMESTER-6$ |
```

```
rheaubuntu@LAPTOP-RB9PEJ  X  +  v  -  □  X

rheaubuntu@LAPTOP-RB9PEJTU:~/SEMESTER-6$ lscpu
Architecture:          x86_64
CPU op-mode(s):        32-bit, 64-bit
Address sizes:          39 bits physical, 48 bits virtual
Byte Order:             Little Endian
CPU(s):                 8
On-line CPU(s) list:   0-7
Vendor ID:              GenuineIntel
Model name:             Intel(R) Core(TM) i5-9300H CPU @ 2.40GHz
CPU family:             6
Model:                  158
Thread(s) per core:    2
Core(s) per socket:    4
Socket(s):              1
Stepping:               10
BogoMIPS:               4800.01
Flags:                  fpu vme de pse tsc msr pae mce cx8 apic sep
                        mtrr pge mca cmov pat pse36 clflush mmx fxsr
                        sse sse2 ss ht syscall nx pdpe1gb rdtscp lm
                        constant_tsc arch_perfmon rep_good nopl xto
                        pology cpuid pni pclmulqdq vmx ssse3 fma cx1
                        6 pdcm pcid sse4_1 sse4_2 movbe popcnt aes x
                        save avx f16c rdrand hypervisor lahf_lm abm
                        3dnowprefetch invpcid_single pti ssbd ibrs i
                        bpb stibp tpr_shadow vnmi ept vpid ept_ad fs
                        gsbase bml avx2 smep bmi2 erms invpcid rdse
                        ed adx smap clflushopt xsaveopt xsavec xgetb
                        v1 xsaves md_clear flush_lld arch_capabilities
Virtualization features:
  Virtualization:       VT-x
  Hypervisor vendor:    Microsoft
  Virtualization type:  full
Caches (sum of all):
  L1d:                  128 KiB (4 instances)
  L1i:                  128 KiB (4 instances)
  L2:                   1 MiB (4 instances)
  L3:                   8 MiB (1 instance)
Vulnerabilities:
  Gather data sampling:  Unknown: Dependent on hypervisor status
  Itlb multihit:        KVM: Mitigation: VMX disabled
  L1tf:                 Mitigation; PTE Inversion; VMX conditional c
                        ache flushes, SMT vulnerable
  Mds:                  Mitigation; Clear CPU buffers; SMT Host stat
                        e unknown
  Meltdown:             Mitigation; PTI
  Mmio stale data:      Mitigation; Clear CPU buffers; SMT Host stat
                        e unknown
  Retbleed:             Mitigation; IBRS
```

```
rheaubuntu@LAPTOP-RB9PEJ X + v - □ X
rheaubuntu@LAPTOP-RB9PEJTU:~/SEMESTER-6$ ls
a.exe      matmul.cpp  matmul_output.txt  pointer.cpp  rand_mat.exe
hello.cpp  matmul.exe  openmp             rand_mat.cpp
rheaubuntu@LAPTOP-RB9PEJTU:~/SEMESTER-6$ find *.cpp
hello.cpp
matmul.cpp
pointer.cpp
rand_mat.cpp
rheaubuntu@LAPTOP-RB9PEJTU:~/SEMESTER-6$ |
```

```
rheaubuntu@LAPTOP-RB9PEJ X + v - □ X
rheaubuntu@LAPTOP-RB9PEJTU:~/semester-7$ touch empty.txt
rheaubuntu@LAPTOP-RB9PEJTU:~/semester-7$ ls
empty.txt
rheaubuntu@LAPTOP-RB9PEJTU:~/semester-7$ |
```

```
rheaubuntu@LAPTOP-RB9PEJ X + v - □ X
rheaubuntu@LAPTOP-RB9PEJTU:~/semester-7$ nano empty.txt
rheaubuntu@LAPTOP-RB9PEJTU:~/semester-7$ cat empty.txt
Hello
Empty File!!!!
rheaubuntu@LAPTOP-RB9PEJTU:~/semester-7$ clear|
```

```
rheaubuntu@LAPTOP-RB9PEJ X + v - □ X
rheaubuntu@LAPTOP-RB9PEJTU:~/semester-7$ cd ..
rheaubuntu@LAPTOP-RB9PEJTU:~$ ps -A
  PID TTY          TIME CMD
    1 ?           00:00:09 systemd
    2 ?           00:00:00 init-systemd(Ub
    7 ?           00:00:00 init
   40 ?           00:00:00 systemd-journal
   63 ?           00:00:00 systemd-udevd
   80 ?           00:00:00 snapfuse
   81 ?           00:00:00 snapfuse
   82 ?           00:00:00 snapfuse
   87 ?           00:00:00 snapfuse
   91 ?           00:00:00 snapfuse
  102 ?           00:00:02 snapfuse
  103 ?           00:00:00 snapfuse
  104 ?           00:00:01 snapfuse
  111 ?           00:00:00 systemd-resolve
  135 ?           00:00:00 cron
  139 ?           00:00:00 dbus-daemon
  155 ?           00:00:00 networkd-dispat
```

```
rheaubuntu@LAPTOP-RB9PEJ:~$ help
GNU bash, version 5.1.16(1)-release (x86_64-pc-linux-gnu)
These shell commands are defined internally. Type 'help' to see this list.
Type 'help name' to find out more about the function 'name'.
Use 'info bash' to find out more about the shell in general.
Use 'man -k' or 'info' to find out more about commands not in this list.

A star (*) next to a name means that the command is disabled.

job_spec [&]
(( expression ))
. filename [arguments]
:
[ arg... ]
[[ expression ]]
alias [-p] [name=value] ... ]
bg [job_spec ...]
bind [-lpsvPSVX] [-m keymap] [-f filename] [-q name] [-u name] [->
break [n]
builtin [shell-builtin [arg ...]]
caller [expr]
case WORD in [PATTERN [| PATTERN]...] COMMANDS ;;)... esac
cd [-L|[-P [-e]] [-@]] [dir]
command [-pVv] command [arg ...]
compgen [-abcdefgjkuv] [-o option] [-A action] [-G globpat] [-W >
complete [-abcdefgjkuv] [-pr] [-DEI] [-o option] [-A action] [-G>
comptopt [-o|+o option] [-DEI] [name ...]
continue [n]
coproc [NAME] command [redirections]
history [-c] [-d offset] [n] or history -anrw [filename] or hist>
if COMMANDS; then COMMANDS; [ elif COMMANDS; then COMMANDS; ]...>
jobs [-lnprs] [jobspec ...] or jobs -x command [args]
kill [-s sigspec | -n signum | -sigspec] pid | jobspec ... or ki>
let arg [arg ...]
local [option] name[=value] ...
logout [n]
mapfile [-d delim] [-n count] [-O origin] [-s count] [-t] [-u fd>
popd [-n] [+N | -N]
printf [-v var] format [arguments]
pushd [-n] [+N | -N | dir]
pwd [-LP]
read [-ers] [-a array] [-d delim] [-i text] [-n nchars] [-N ncha>
readarray [-d delim] [-n count] [-O origin] [-s count] [-t] [-u >
readonly [-aAf] [name[=value] ...] or readonly -p
return [n]
select NAME [in WORDS ... ;] do COMMANDS; done
set [-abefhkmnptuvxBCHP] [-o option-name] [--] [arg ...]
shift [n]
shopt [-pqsu] [-o] [optname ...]
```

```
rheaubuntu@LAPTOP-RB9PEJ:~/SEMESTER-6$ ls
a.exe      matmul.cpp  matmul_output.txt  pointer.cpp  rand_mat.exe
hello.cpp  matmul.exe  openmp             rand_mat.cpp
rheaubuntu@LAPTOP-RB9PEJ:~/SEMESTER-6$ ls -l | grep "s/.*/p"
rheaubuntu@LAPTOP-RB9PEJ:~/SEMESTER-6$ |
```

```
rheaubuntu@LAPTOP-RB9PEJ:~/SEMESTER-6$ ps
  PID TTY          TIME CMD
  377 pts/0        00:00:00 bash
 9579 pts/0        00:00:00 ps
rheaubuntu@LAPTOP-RB9PEJ:~/SEMESTER-6$ |
```

```
rheaubuntu@LAPTOP-RB9PEJ x + v
rheaubuntu@LAPTOP-RB9PEJTU:~/SEMESTER-6$ wget https://wordpress.org/latest.tar.gz
--2024-07-16 16:41:32-- https://wordpress.org/latest.tar.gz
Resolving wordpress.org (wordpress.org)... 198.143.164.252
Connecting to wordpress.org (wordpress.org)|198.143.164.252|:443... connected.
HTTP request sent, awaiting response... 200 OK
Length: 24696391 (24M) [application/octet-stream]
Saving to: 'latest.tar.gz'

latest.tar.gz      100%[=====>]  23.55M  7.35MB/s   in 5.1s

2024-07-16 16:41:39 (4.63 MB/s) - 'latest.tar.gz' saved [24696391/24696391]

rheaubuntu@LAPTOP-RB9PEJTU:~/SEMESTER-6$ |
```

```
rheaubuntu@LAPTOP-RB9PEJ x + v
rheaubuntu@LAPTOP-RB9PEJTU:~/SEMESTER-6$ passwd|
```

```
rheaubuntu@LAPTOP-RB9PEJ x + v
rheaubuntu@LAPTOP-RB9PEJTU:/$ which bash
/usr/bin/bash
rheaubuntu@LAPTOP-RB9PEJTU:/$ |
```

```
rheaubuntu@LAPTOP-RB9PEJ:~$ df -Th
Filesystem      Type      Size  Used Avail Use% Mounted on
none            tmpfs     1.9G  4.0K  1.9G   1% /mnt/wsl
drivers         9p        238G  217G   22G  92% /usr/lib/wsl/drivers
none           tmpfs     1.9G    0  1.9G   0% /usr/lib/modules
none           overlay   1.9G    0  1.9G   0% /usr/lib/modules/5.15.153.1-microsoft-standard-WSL2
/dev/sdc        ext4     1007G  3.7G  952G   1% /
none           tmpfs     1.9G   96K  1.9G   1% /mnt/wslg
none           overlay   1.9G    0  1.9G   0% /usr/lib/wsl/lib
rootfs         rootfs    1.9G  2.1M  1.9G   1% /init
none           tmpfs     1.9G  844K  1.9G   1% /run
none           tmpfs     1.9G    0  1.9G   0% /run/lock
none           tmpfs     1.9G    0  1.9G   0% /run/shm
tmpfs          tmpfs     4.0M    0  4.0M   0% /sys/fs/cgroup
none           overlay   1.9G   76K  1.9G   1% /mnt/wslg/versions.txt
none           overlay   1.9G   76K  1.9G   1% /mnt/wslg/doc
C:\             9p        238G  217G   22G  92% /mnt/c
D:\            9p        932G   13G  919G   2% /mnt/d
snapfuse       fuse.snapfuse 64M   64M    0 100% /snap/core20/2264
snapfuse       fuse.snapfuse 75M   75M    0 100% /snap/core22/1033
snapfuse       fuse.snapfuse 128K  128K    0 100% /snap/bare/5
snapfuse       fuse.snapfuse 75M   75M    0 100% /snap/core22/1380
snapfuse       fuse.snapfuse 92M   92M    0 100% /snap/gtk-common-themes/1535
snapfuse       fuse.snapfuse 39M   39M    0 100% /snap/snapd/21465
snapfuse       fuse.snapfuse 131M  131M    0 100% /snap/ubuntu-desktop-installer/1284
snapfuse       fuse.snapfuse 132M  132M    0 100% /snap/ubuntu-desktop-installer/1286
snapfuse       fuse.snapfuse 39M   39M    0 100% /snap/snapd/21759
snapfuse       fuse.snapfuse 64M   64M    0 100% /snap/core20/2318
rheaubuntu@LAPTOP-RB9PEJ:~$
```

```
rheaubuntu@LAPTOP-RB9PEJ:~/semester-7$ cat empty.txt
My name is Rhea.
rheaubuntu@LAPTOP-RB9PEJ:~/semester-7$ wc -c empty.txt
17 empty.txt
rheaubuntu@LAPTOP-RB9PEJ:~/semester-7$ wc -w empty.txt
4 empty.txt
rheaubuntu@LAPTOP-RB9PEJ:~/semester-7$ wc -l empty.txt
1 empty.txt
rheaubuntu@LAPTOP-RB9PEJ:~/semester-7$
```

**Result:**

Successfully implemented various basic linux commands



<b>Ex. No: 2</b>	<b>Install and Configure Hadoop</b>
<b>23.7.2024</b>	

### **Aim:**

To install and configure Hadoop in Ubuntu

### **Program:**

#### 1. Install Java:

- Check if Java is installed: `java -version`
- If Java is not installed, download and install it:  
`sudo apt update`  
`sudo apt install openjdk-8-jdk`

#### 2. Download Hadoop:

- Go to the Apache Hadoop releases page and download the binary distribution (e.g., Hadoop 3.3.6):

```
wget https://downloads.apache.org/hadoop/common/hadoop-3.3.0/hadoop-3.4.0.tar.gz
```

- Extract the downloaded file:  
`tar -xzvf hadoop-3.3.6.tar.gz`

#### 3. Move Hadoop to the desired directory:

```
sudo mv hadoop-3.3.6 /opt/hadoop
```

#### 4. Set environment variables:

- Open the `~/.bashrc` file for editing: `nano ~/.bashrc`
- Add the following lines at the end of the file:  
`export HADOOP_HOME=/opt/hadoop` `export`  
`PATH=$PATH:$HADOOP_HOME/bin:$HADOOP_HOME/sbin`
- Save and exit, then refresh your terminal

#### 5. Configure Hadoop: `source ~/.bashrc`

- Navigate to the Hadoop configuration directory:  
`cd $HADOOP_HOME/etc/hadoop/`
- Edit `hadoop-env.sh` to set the Java home path:  
`nano hadoop-env.sh`
- Add or modify the line to include your Java installation path:  
`export JAVA_HOME=/usr/lib/jvm/java-11-openjdk-amd64`

#### 6. Format the NameNode: `hadoop namenode -format`

#### 7. Start YARN: `start-yarn.sh`

#### 8. Start all Hadoop services: `start-all.sh`

#### 9. Check running processes: `jps`

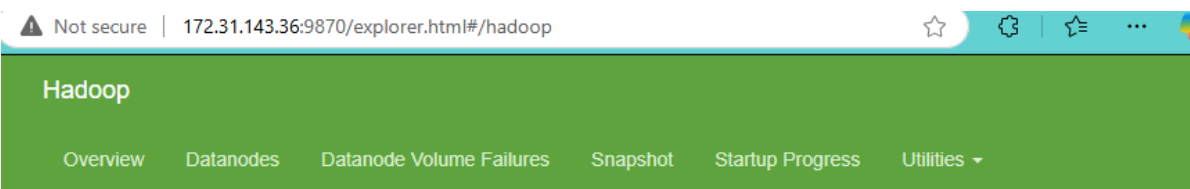
#### 10. Access Hadoop web interface: <http://localhost:9870>

#### 11. Stop all Hadoop services when needed: `stop-all.sh`

### **Output:**

```
hadoop@LAPTOP-RB9PEJTU:~$ hadoop version
Hadoop 3.3.6
Source code repository https://github.com/apache/hadoop.git -r 1be78238728da9266a4f88195058f08fd012bf9c
Compiled by ubuntu on 2023-06-18T08:22Z
Compiled on platform linux-x86_64
Compiled with protoc 3.7.1
From source with checksum 5652179ad55f76cb287d9c633bb53bbd
This command was run using /home/hadoop/hadoop/share/hadoop/common/hadoop-common-3.3.6.jar
hadoop@LAPTOP-RB9PEJTU:~$
```

```
hadoop@LAPTOP-RB9PEJTU:~$ start-all.sh
WARNING: Attempting to start all Apache Hadoop daemons as hadoop in 10 seconds.
WARNING: This is not a recommended production deployment configuration.
WARNING: Use CTRL-C to abort.
Starting namenodes on [localhost]
Starting datanodes
Starting secondary namenodes [LAPTOP-RB9PEJTU]
Starting resourcemanager
Starting nodemanagers
```



## Browse Directory

Search:

Show 

25

 entries

<input type="checkbox"/>	Permission	Owner	Group	Size	Last Modified	Replication	Block Size	Name	
<input type="checkbox"/>	<a href="#">drwxrwxr-x</a>	<a href="#">hadoop</a>	<a href="#">supergroup</a>	0 B	Nov 03 17:18	<a href="#">0</a>	0 B	<a href="#">hadoop</a>	

Showing 1 to 1 of 1 entries

Previous

1

Next

Hadoop, 2023.

### Result:

Successfully installed and configured Hadoop in Ubuntu

<b>Ex. No: 3</b>	<b>Implementing MapReduce</b>
<b>30.7.2024</b>	

### **Aim:**

To implement a simple map-reduce code for the wordcount problem in Hadoop.

### **Program:**

#### **WordCount.java**

```
import java.io.IOException;
import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.*;
import org.apache.hadoop.mapreduce.*;
import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;

public class WordCount {

    public static class TokenizerMapper extends Mapper<LongWritable,
Text, Text, IntWritable> {

        private final static IntWritable one = new IntWritable(1);
        private Text word = new Text();

        public void map(LongWritable key, Text value, Context context)
throws IOException, InterruptedException {

            String[] words = value.toString().split("\\s+");
            for (String str : words) {
                word.set(str);
                context.write(word, one);
            }
        }
    }
}
```

```
}
```

```
    public static class IntSumReducer extends Reducer<Text, IntWritable,  
Text, IntWritable> {
```

```
        private IntWritable result = new IntWritable();
```

```
        public void reduce(Text key, Iterable<IntWritable> values,  
Context context) throws IOException, InterruptedException {
```

```
            int sum = 0;
```

```
            for (IntWritable val : values) {
```

```
                sum += val.get();
```

```
            }
```

```
            result.set(sum);
```

```
            context.write(key, result);
```

```
        }
```

```
    }
```

```
    public static void main(String[] args) throws Exception {
```

```
        Configuration conf = new Configuration();
```

```
        Job job = Job.getInstance(conf, "word count");
```

```
        job.setJarByClass(WordCount.class);
```

```
        job.setMapperClass(TokenizerMapper.class);
```

```
        job.setCombinerClass(IntSumReducer.class);
```

```
        job.setReducerClass(IntSumReducer.class);
```

```
        job.setOutputKeyClass(Text.class);
```

```
        job.setOutputValueClass(IntWritable.class);
```

```
        FileInputFormat.addInputPath(job, new Path(args[0]));
```

```
        FileOutputFormat.setOutputPath(job, new Path(args[1]));
```

```
        System.exit(job.waitForCompletion(true) ? 0 : 1);
```

```
    }
```

```
}
```

Output:

```
hadoop@LAPTOP-RB9PEJ7U:~/hadoop/wordcount_java$ hdfs dfs -cat /hadoop/hadoop/wordcount_java/output/part-r-00000
Hadoop 1
Hello 3
Rhea 1
World 1
```

Browse Directory

Show 

25

 entries 

Search:

<input type="checkbox"/>	Permission	Owner	Group	Size	Last Modified	Replication	Block Size	Name	
<input type="checkbox"/>	<a href="#">-rw-r--r--</a>	<a href="#">hadoop</a>	<a href="#">supergroup</a>	0 B	Oct 20 12:01	<a href="#">1</a>	128 MB	<a href="#">_SUCCESS</a>	
<input type="checkbox"/>	<a href="#">-rw-r--r--</a>	<a href="#">hadoop</a>	<a href="#">supergroup</a>	32 B	Oct 20 12:01	<a href="#">1</a>	128 MB	<a href="#">part-r-00000</a>	

Showing 1 to 2 of 2 entries

Previous

1

Next

Hadoop, 2023.

Result:

Successfully implemented map reduce for wordcount problem in Hadoop

<b>Ex. No: 4</b>	<b>Implementing MapReduce 2</b>
<b>6.8.2024</b>	

**Aim:**

1.Implement map reduce for NCDC weather dataset using Hadoop – find the max and min temperature

2.Implement Apriori algorithm using map reduce paradigm

**Program 1:**

**Temperature.java**

```
import java.io.IOException;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.io.*;
import org.apache.hadoop.mapreduce.*;
import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;

public class Temperatures {

    public static class TempMapper extends Mapper<LongWritable, Text,
    Text, IntWritable> {

        public void map(LongWritable key, Text value, Context context)
        throws IOException, InterruptedException {

            String line = value.toString();
            String[] fields = line.split("\t");

            if (fields.length >= 8) { // Ensure there are enough fields
                try {
                    // Extract the minimum and maximum temperature fields
```

```

        int minTemp = Integer.parseInt(fields[6].trim()); //
MLY-TMIN-NORMAL

        int maxTemp = Integer.parseInt(fields[7].trim()); //
MLY-TMAX-NORMAL

        // Write the min and max temperatures to the context
        context.write(new Text("Min Temperature"), new
IntWritable(minTemp));

        context.write(new Text("Max Temperature"), new
IntWritable(maxTemp));

        } catch (NumberFormatException e) {
            // Ignore invalid data
        }
    }
}

}

}

}

    public static class TempReducer extends Reducer<Text, IntWritable,
Text, IntWritable> {

        public void reduce(Text key, Iterable<IntWritable> values,
Context context) throws IOException, InterruptedException {

            int extremeTemp = key.toString().equals("Min Temperature") ?
Integer.MAX_VALUE : Integer.MIN_VALUE;

            for (IntWritable value : values) {
                int temp = value.get();
                if (key.toString().equals("Min Temperature")) {
                    if (temp < extremeTemp) {
                        extremeTemp = temp;
                    }
                } else { // Max Temperature
                    if (temp > extremeTemp) {

```

```

        extremeTemp = temp;
    }
}
}
context.write(key, new IntWritable(extremeTemp));
}
}

public static void main(String[] args) throws Exception {
    Configuration conf = new Configuration();
    Job job = Job.getInstance(conf, "min and max temperatures");
    job.setJarByClass(Temperatures.class);
    job.setMapperClass(TempMapper.class);
    job.setReducerClass(TempReducer.class);
    job.setOutputKeyClass(Text.class);
    job.setOutputValueClass(IntWritable.class);
    FileInputFormat.addInputPath(job, new Path(args[0]));
    FileOutputFormat.setOutputPath(job, new Path(args[1]));
    System.exit(job.waitForCompletion(true) ? 0 : 1);
}
}

```

## **Program 2:**

### **Mapper.py**

```

#!/usr/bin/env python3

import sys

from itertools import combinations

def generate_combinations(item_list, length):
    return list(combinations(item_list, length))

# Input comes from standard input (stdin)

```



```

for line in sys.stdin:
    line = line.strip()
    items = line.split()
    for length in range(1, len(items) + 1):
        for combination in generate_combinations(items, length):
            print(f"{'.'.join(combination)}\t1")

```

### **Reducer.py**

```

#!/usr/bin/env python3
import sys
current_itemset = None
current_count = 0 #
Input comes from standard input (stdin)
for line in sys.stdin:
    line = line.strip()
    itemset, count = line.split('\t', 1)
    count = int(count)
    if current_itemset == itemset:
        current_count += count
    else:
        if current_itemset:
            print(f"{current_itemset}\t{current_count}")
            current_count = count
            current_itemset = itemset
        if current_itemset == itemset:
            print(f"{current_itemset}\t{current_count}")

```

### **Output 1 : NCDC Dataset**

```
hadoop@LAPTOP-RB9PEJTU:~/hadoopdata/mapreduce_code$ hdfs dfs -cat /temperature/output/part-r-00000
Max Temperature 793
Min Temperature -43
```

## Browse Directory

Show  entries
Search:

<input type="checkbox"/>	Permission	Owner	Group	Size	Last Modified	Replication	Block Size	Name	
<input type="checkbox"/>	<a href="#">-rw-r--r--</a>	<a href="#">hadoop</a>	<a href="#">supergroup</a>	0 B	Sep 10 14:47	<a href="#">1</a>	128 MB	<a href="#">_SUCCESS</a>	
<input type="checkbox"/>	<a href="#">-rw-r--r--</a>	<a href="#">hadoop</a>	<a href="#">supergroup</a>	40 B	Sep 10 14:47	<a href="#">1</a>	128 MB	<a href="#">part-r-00000</a>	

Showing 1 to 2 of 2 entries

Hadoop, 2023.

## Output 2 : Apriori Algorithm

```
File Input Format Counters
  Bytes Read=67
File Output Format Counters
  Bytes Written=94
2024-11-03 17:20:55,898 INFO streaming.StreamJob: Output directory: /hadoop/hadoop/apriori_python/output
hadoop@LAPTOP-RB9PEJTU:~/hadoop/apriori_python$ hdfs dfs -cat /hadoop/hadoop/apriori_python/output/part-00000
apple 3
apple,banana 2
apple,banana,orange 1
apple,orange 1
banana 4
banana,orange 3
orange 3
hadoop@LAPTOP-RB9PEJTU:~/hadoop/apriori_python$ |
```

# Browse Directory

/hadoop/hadoop/apriori\_python/output

Go!

Show

25

entries

Search:

	Permission	Owner	Group	Size	Last Modified	Replication	Block Size	Name	
<input type="checkbox"/>	<a href="#">-rw-r--r--</a>	<a href="#">hadoop</a>	<a href="#">supergroup</a>	0 B	Nov 03 17:20	<a href="#">1</a>	128 MB	<a href="#">_SUCCESS</a>	
<input type="checkbox"/>	<a href="#">-rw-r--r--</a>	<a href="#">hadoop</a>	<a href="#">supergroup</a>	94 B	Nov 03 17:20	<a href="#">1</a>	128 MB	<a href="#">part-00000</a>	

Showing 1 to 2 of 2 entries

Previous

1

Next

Hadoop, 2023.

## Result:

Successfully implemented map reduce for NCDC Dataset and executed apriori algorithm using map reduce paradigm.

<b>Ex. No: 5</b>	<b>Spark and PySpark</b>
<b>13.8.2024</b>	

### **Aim:**

Install spark and pyspark. Run a spark shell and test the installation. Run the wordcount program that you did using Hadoop using pyspark. Use the movielens dataset and try to find out for each movie, how are the ratings distributed.

### **Program**

#### 1. Download and Install Spark

- a. `wget https://www.apache.org/dyn/closer.lua/spark/spark-3.5.0/spark-3.5.0-bin-hadoop3.tgz`
- b. `tar -xzf spark-3.5.0-bin-hadoop3.tgz`
- c. `mv spark-3.5.0-bin-hadoop3 spark`

#### 2. Set Environment Variables

- a. `nano ~/.bashrc`
- b. `export SPARK_HOME=~/.spark`
- c. `export PATH=$PATH:$SPARK_HOME/bin`
- d. `export PYTHONPATH=$SPARK_HOME/python:$PYTHONPATH`
- e. `export PYSARK_PYTHON=python3`
- f. `source ~/.bashrc`

#### 3. Running Programs : `spark-submit file.py`

### **Wordcount :**

```
from pyspark import SparkContext
```

```
# Initialize SparkContext
```

```
sc = SparkContext("local", "Word Count")
```

```

# Read input file
input_file = "/home/snucse/Desktop/wordcount_input.txt"
text_file = sc.textFile(input_file)

# Count words
counts = text_file.flatMap(lambda line: line.split(" ")) \
                    .map(lambda word: (word, 1)) \
                    .reduceByKey(lambda a, b: a + b)

# Collect and print results
for word, count in counts.collect():
    print(f"{word}: {count}")

# Stop the SparkContext
sc.stop()

```

### **Program Movie-lens :**

```

from pyspark.sql import SparkSession
from pyspark.sql import functions as F

# Initialize Spark session
spark = SparkSession.builder \
    .appName("MovieLens Ratings Distribution") \
    .getOrCreate()

# Load the MovieLens dataset (change the path accordingly)
ratings = spark.read.csv("/home/snucse/Desktop/movielens.csv",
    header=True, inferSchema=True)

# Show the original DataFrame and column names
ratings.show()

```

```

print(ratings.columns)

# Clean up column names (if needed)
ratings = ratings.toDF(*[c.strip() for c in ratings.columns])

# Calculate the count of ratings and average rating for each movieId
rating_distribution = ratings.groupBy("movieId").agg(
    F.count("rating").alias("rating_count"),
    F.avg("rating").alias("average_rating")
)

# Show the results
rating_distribution.show()

# Stop the Spark session
spark.stop()

```

### Output Wordcount :

```

24/10/08 09:45:40 INFO TaskSchedulerImpl: Killing all running tasks in stage 1: Stage finished
24/10/08 09:45:40 INFO DAGScheduler: Job 0 finished: collect at /home/snucse/Desktop/wordcount.py:16, took 0.894432 s
My: 1
name: 1
is: 2
Rhea.: 1
Rhea: 1
sleeping.: 1
24/10/08 09:45:40 INFO SparkContext: SparkContext is stopping with exitCode 0.
24/10/08 09:45:40 INFO ShutdownHookManager: Shutdown hook called
24/10/08 09:45:40 INFO ShutdownHookManager: Shutdown hook called

```

### Output Movie-lens:

```

24/10/08 10:18:46 INFO DAGScheduler: Job 2 finished: showString at NativeMethodAccessorImpl.java:0, took 0.042936 s
24/10/08 10:18:46 INFO CodeGenerator: Code generated in 6.975829 ms
+-----+
|userId|movieId|rating|timestamp|
+-----+
| 196| 242| 3|881250949|
| 186| 302| 3|891717742|
| 22| 377| 1|878887116|
| 244| 51| 2|880606923|
| 166| 346| 1|886397596|
| 298| 474| 4|884182806|
| 115| 265| 2|881171488|
| 253| 465| 5|891628467|
| 305| 451| 3|886324817|
| 6| 86| 3|883603013|
| 62| 257| 2|879372434|
| 286| 1014| 5|879781125|
| 200| 222| 5|876042340|
| 210| 40| 3|891035994|
| 224| 29| 3|888104457|
| 303| 785| 3|879485318|
| 122| 387| 5|879270459|
| 194| 274| 2|879539794|
| 291| 1042| 4|874834944|
| 234| 1184| 2|892079237|
+-----+
only showing top 20 rows

['userId', 'movieId', 'rating', 'timestamp']
24/10/08 10:18:46 INFO BlockManagerInfo: Removed broadcast_1_piece0 on 10.23.22.124:35821 in memory (size: 6.4 KiB, free: 366.2 MiB)
24/10/08 10:18:46 INFO DAGScheduler: Job 4 finished: showString at NativeMethodAccessorImpl.java:0, took 0.063932 s
24/10/08 10:18:46 INFO CodeGenerator: Code generated in 3.72489 ms
+-----+
|movieId|rating_count| average_rating|
+-----+
| 496| 231| 4.121212121212121|
| 471| 221| 3.6108597285067874|
| 463| 71| 3.859154929577465|
| 148| 128| 3.203125|
| 1342| 2| 2.5|
| 833| 49| 3.204081632653061|
| 1088| 13| 2.230769230769231|
| 1591| 6| 3.1666666666666665|
| 1238| 8| 3.125|
| 1580| 1| 1.0|
| 1645| 1| 4.0|
| 392| 68| 3.5441176470588234|
| 623| 39| 2.923076923076923|
| 540| 43| 2.511627906976744|
| 858| 3| 1.0|
| 737| 59| 2.983050847457627|
| 243| 132| 2.4393939393939394|
| 1025| 44| 2.9318181818181817|
| 1084| 21| 3.857142857142857|
| 1127| 11| 2.909090909090909|
+-----+
only showing top 20 rows

24/10/08 10:18:46 INFO SparkContext: SparkContext is stopping with exitCode 0.
24/10/08 10:18:46 INFO SparkUI: Stopped Spark web UI at http://10.23.22.124:4040
24/10/08 10:18:46 INFO MapOutputTrackerMasterEndpoint: MapOutputTrackerMasterEndpoint stopped!
24/10/08 10:18:46 INFO MemoryStore: MemoryStore cleared
24/10/08 10:18:46 INFO BlockManager: BlockManager stopped
24/10/08 10:18:46 INFO BlockManagerMaster: BlockManagerMaster stopped
24/10/08 10:18:46 INFO OutputCommitCoordinator$OutputCommitCoordinatorEndpoint: OutputCommitCoordinator stopped!
24/10/08 10:18:46 INFO SparkContext: Successfully stopped SparkContext
24/10/08 10:18:47 INFO ShutdownHookManager: Shutdown hook called
24/10/08 10:18:47 INFO ShutdownHookManager: Deleting directory /tmp/spark-1c7edb2d-80f9-49ef-83a8-dfbee853b850

```

## Result:

Successfully installed Spark and Pyspark, and implemented the Wordcount problem and distribution of movie reviews using them.

<b>Ex. No: 6</b>	<b>PySpark 2</b>
<b>20.8.2024</b>	

### **Aim:**

1. Use the friends\_test dataset. Col1 is ID, Col2 is name, Col3 is Age and Col4 is num of friends. Understand mapvalues function of RDD in spark and find the average number of friends for each unique age present in the dataset.
2. Use the temp\_csv dataset. Column headers are present in the dataset. Understand filter operations and filter out only the "TMIN" values from the "desc" column. With the resultant data (RDD), find the following:
  - a. Minimum temperature (overall)
  - b. Minimum temperature for every ItemID
  - c. Minimum emperature for every StationID
3. Use the same dataset, filter only "TMAX" column and find the maximum temperatures just like the ones mentioned above.

### **Program:**

#### **Friends\_test\_analysis.py**

```

from pyspark import SparkConf, SparkContext

from pyspark.sql

import * from pyspark.sql.functions
import * from pyspark.sql.types import *


spark = SparkSession.builder.appName("friends
test").config("spark.memory.offHeap.e

df = spark.read.csv('friends_test.csv',header=False)
df.explain()
spark.stop()

conf = SparkConf().setAppName("Basicapp").setMaster("local[*]")
sc = SparkContext(conf=conf)

rdd = sc.textFile("friends_test.csv")
rdd.first()

rdd_split = rdd.map(lambda line: line.split(","))
for row in rdd_split.take(5):

```



```
print(row)
```

### **Temp\_analysis.py**

```
from pyspark import SparkConf, SparkContext
```

```
# Initialize Spark Context
```

```
conf = SparkConf().setAppName("TempDataset").setMaster("local[*]")
```

```
sc = SparkContext(conf=conf)
```

```
rdd = sc.textFile("temp.csv")
```

```
# Split data and remove header
```

```
rdd_header = rdd.first()
```

```
rdd_data = rdd.filter(lambda row: row != rdd_header).map(lambda row:  
row.split(","))
```

```
# Filter for TMIN and compute minimum temperatures
```

```
rdd_TMIN_filter = rdd_data.filter(lambda row: row[2] == "TMIN")
```

```
rdd_min_overall = rdd_TMIN_filter.map(lambda x: int(x[3])).reduce(lambda  
a, b: a if a < b else b)
```

```
print("Minimum temperature overall:", rdd_min_overall)
```

```
rdd_min_itemID = rdd_TMIN_filter.map(lambda x: (x[0],  
int(x[3]))).reduceByKey(lambda a, b: a if a < b else b)
```

```
print("Minimum temperature by ItemID:", rdd_min_itemID.collect())
```

```
rdd_min_stationID = rdd_TMIN_filter.map(lambda x: (x[1],  
int(x[3]))).reduceByKey(lambda a, b: a if a < b else b)
```

```
print("Minimum temperature by StationID:", rdd_min_stationID.collect())
```

```
# Filter for TMAX and compute maximum temperatures
```

```
rdd_TMAX_filter = rdd_data.filter(lambda row: row[2] == "TMAX")
```

```
rdd_max_overall = rdd_TMAX_filter.map(lambda x: int(x[3])).reduce(lambda
a, b: a if a > b else b)
```

```
print("Maximum temperature overall:", rdd_max_overall)
```

```
rdd_max_itemID = rdd_TMAX_filter.map(lambda x: (x[0],
int(x[3]))).reduceByKey(lambda a, b: a if a > b else b)
```

```
print("Maximum temperature by ItemID:", rdd_max_itemID.collect())
```

```
rdd_max_stationID = rdd_TMAX_filter.map(lambda x: (x[1],
int(x[3]))).reduceByKey(lambda a, b: a if a > b else b)
```

```
print("Maximum temperature by StationID:", rdd_max_stationID.collect())
```

```
sc.stop()
```

## Output : Friends\_test\_analysis.py (part a)

```
hadoop@LAPTOP-RB9PEJTU:~/hadoop/pyspark_friends$ spark-submit friends_test_analysis.py
24/11/03 18:03:06 WARN Utils: Your hostname, LAPTOP-RB9PEJTU resolves to a loopback address: 127.0.1.1; using 10.255.255.254 instead (on interface lo)
24/11/03 18:03:06 WARN Utils: Set SPARK_LOCAL_IP if you need to bind to another address
Using Spark's default log4j profile: org/apache/spark/log4j-defaults.properties
24/11/03 18:03:07 INFO SparkContext: Running Spark version 3.2.0
```

```
24/11/03 18:03:11 INFO TaskSchedulerImpl: Removed TaskSet 6.0, whose tasks have all completed, from pool
24/11/03 18:03:11 INFO DAGScheduler: ResultStage 6 (collect at /home/hadoop/hadoop/pyspark_friends/friends_test_analysis.py:17) finished in 0.068 s
24/11/03 18:03:11 INFO DAGScheduler: Job 2 is finished. Cancelling potential speculative or zombie tasks for this job
24/11/03 18:03:11 INFO TaskSchedulerImpl: Killing all running tasks in stage 6: Stage finished
24/11/03 18:03:11 INFO DAGScheduler: Job 2 finished: collect at /home/hadoop/hadoop/pyspark_friends/friends_test_analysis.py:17, took 0.156635 s
Age: 18, Average Number of Friends: 343.38
Age: 19, Average Number of Friends: 213.27
Age: 20, Average Number of Friends: 165.00
Age: 21, Average Number of Friends: 350.88
Age: 22, Average Number of Friends: 206.43
Age: 23, Average Number of Friends: 246.30
Age: 24, Average Number of Friends: 233.80
Age: 25, Average Number of Friends: 197.45
Age: 26, Average Number of Friends: 242.06
Age: 27, Average Number of Friends: 228.12
Age: 28, Average Number of Friends: 209.10
Age: 29, Average Number of Friends: 215.92
Age: 30, Average Number of Friends: 235.82
Age: 31, Average Number of Friends: 267.25
Age: 32, Average Number of Friends: 207.91
Age: 33, Average Number of Friends: 325.33
Age: 34, Average Number of Friends: 245.50
Age: 35, Average Number of Friends: 211.62
Age: 36, Average Number of Friends: 246.60
Age: 37, Average Number of Friends: 249.33
Age: 38, Average Number of Friends: 193.53
Age: 39, Average Number of Friends: 169.29
Age: 40, Average Number of Friends: 250.82
Age: 41, Average Number of Friends: 268.56
Age: 42, Average Number of Friends: 303.50
Age: 43, Average Number of Friends: 230.57
Age: 44, Average Number of Friends: 282.17
Age: 45, Average Number of Friends: 309.54
Age: 46, Average Number of Friends: 223.69
Age: 47, Average Number of Friends: 233.22
Age: 48, Average Number of Friends: 281.40
Age: 49, Average Number of Friends: 184.67
Age: 50, Average Number of Friends: 254.60
Age: 51, Average Number of Friends: 302.14
```

```

Age: 52, Average Number of Friends: 340.64
Age: 53, Average Number of Friends: 222.86
Age: 54, Average Number of Friends: 278.08
Age: 55, Average Number of Friends: 295.54
Age: 56, Average Number of Friends: 306.67
Age: 57, Average Number of Friends: 258.83
Age: 58, Average Number of Friends: 116.55
Age: 59, Average Number of Friends: 220.00
Age: 60, Average Number of Friends: 202.71
Age: 61, Average Number of Friends: 256.22
Age: 62, Average Number of Friends: 220.77
Age: 63, Average Number of Friends: 384.00
Age: 64, Average Number of Friends: 281.33
Age: 65, Average Number of Friends: 298.20
Age: 66, Average Number of Friends: 276.44
Age: 67, Average Number of Friends: 214.62
Age: 68, Average Number of Friends: 269.60
Age: 69, Average Number of Friends: 235.20
24/11/03 18:03:11 INFO SparkUI: Stopped Spark web UI at http://10.255.255.254:4040
24/11/03 18:03:11 INFO MapOutputTrackerMasterEndpoint: MapOutputTrackerMasterEndpoint stopped!
24/11/03 18:03:11 INFO MemoryStore: MemoryStore cleared
24/11/03 18:03:11 INFO BlockManager: BlockManager stopped
24/11/03 18:03:11 INFO BlockManagerMaster: BlockManagerMaster stopped
24/11/03 18:03:11 INFO OutputCommitCoordinator$OutputCommitCoordinatorEndpoint: OutputCommitCoordinator stopped!
24/11/03 18:03:11 INFO SparkContext: Successfully stopped SparkContext
24/11/03 18:03:12 INFO ShutdownHookManager: Shutdown hook called
24/11/03 18:03:12 INFO ShutdownHookManager: Deleting directory /tmp/spark-bb2dc12a-5411-48e4-b790-a529bdf63ffc
24/11/03 18:03:12 INFO ShutdownHookManager: Deleting directory /tmp/spark-47d3c8d8-4f6d-468d-bb3c-029f245e1e42/pyspark-dda7de57-8078-4f52-be24-e258d9a9e689
24/11/03 18:03:12 INFO ShutdownHookManager: Deleting directory /tmp/spark-47d3c8d8-4f6d-468d-bb3c-029f245e1e42

```

## Output : Test\_analysis.py (part b and c)

```

hadoop@LAPTOP-RB9PEJTU:~/hadoop/pyspark_friends$ spark-submit temp_analysis.py
24/11/03 18:03:32 WARN Utils: Your hostname, LAPTOP-RB9PEJTU resolves to a loopback address: 127.0.0.1; use --hostname if you need to bind to another address
Using Spark's default log4j profile: org/apache/spark/log4j-defaults.properties
24/11/03 18:03:33 INFO SparkContext: Running Spark version 3.2.0
24/11/03 18:03:33 WARN NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java versions instead
24/11/03 18:03:33 INFO ResourceUtils: ==============================================================

```

```
Minimum temperature overall: -148
```

```
Minimum temperature by ItemID: [('ITE00100554', -148), ('EZE00100082', -135)]
```

```

Minimum temperature by StationID: [('18000102', -13),
('18000114', -35), ('18000115', -23), ('18000116',
18000127', 15), ('18000128', 33), ('18000130', 3),
0211', -102), ('18000212', -78), ('18000213', -42),
80000228', -43), ('18000303', -30), ('18000305', -47),
('18000315', -52), ('18000317', -52), ('18000318',
('18000331', 11), ('18000401', 50), ('18000402',
18000420', 125), ('18000421', 153), ('18000422', 15),
('18000502', 139), ('18000503', 154), ('18000504',
('18000519', 105), ('18000526', 160), ('18000531',
), ('18000611', 104), ('18000614', 126), ('18000616',
('18000624', 115), ('18000625', 100), ('18000626',
('18000714', 144), ('18000715', 124), ('18000716',
0), ('18000802', 154), ('18000803', 183), ('18000800',
181), ('18000816', 198), ('18000818', 185), ('18000819',
146), ('18000829', 149), ('18000830', 154), ('18000831',
4', 104), ('18000917', 110), ('18000918', 118), ('18000919',
1', 110), ('18001003', 94), ('18001008', 104), ('18001009',
52), ('18001026', 50), ('18001027', 45), ('18001030',

```

```
Maximum temperature overall: 323
```

```
Maximum temperature by ItemID: [('ITE00100554', 323), ('EZE00100082', 323)]
```

```
24/11/03 18:03:37 INFO DAGScheduler: Job 6 finished
Maximum temperature by StationID: [('18000102', 114), ('18000115', 54), ('18000116', 56), ('18000128', 79), ('18000130', 66), ('18000212', 13), ('18000213', 13), ('18000217', 7), ('18000303', 54), ('18000305', 79), ('18000308', 38), ('18000318', 91), ('18000320', 116), ('18000402', 128), ('18000403', 128), ('18000422', 277), ('18000423', 250), ('18000503', 240), ('18000504', 238), ('18000507', 231), ('18000526', 263), ('18000531', 254), ('18000602', 216), ('18000614', 216), ('18000616', 216), ('18000617', 216), ('18000625', 254), ('18000626', 263), ('18000703', 266), ('18000715', 266), ('18000716', 248), ('18000723', 281), ('18000803', 281), ('18000805', 309), ('18000818', 323), ('18000819', 323), ('18000830', 241), ('18000830', 241), ('18000903', 216), ('18000918', 246), ('18000920', 225), ('18001003', 200), ('18001008', 184), ('18001012', 175), ('18001026', 129), ('18001027', 134), ('18001030', 150), ('18001113', 113), ('18001116', 90), ('18001118', 93)]
```

### Result:

Successfully utilized map value functions for various tasks in pyspark.

Ex. No: 7	Hadoop and Docker
27.8.2024	

### Aim:

Set up a simple Hadoop environment using Docker container, including atleast one NameNode and one DataNode. Ensure the containers are properly configured to interact with each other. After the setup, verify the Hadoop cluster is operational by running a simple HDFS file operation (e.g, uploading a file to HDFS)

### Output:

```
hadoop@LAPTOP-RB9PEJTV:~$ docker network create hadoop-net
c2c3697ab95ac230b5adbc9f905420439251c7df0f024ad2a30145dc8803369
hadoop@LAPTOP-RB9PEJTV:~$ docker run -d --name namenode \
--network hadoop-net \
-e CLUSTER_NAME=test-cluster \
-e CORE_CONF_fs_defaultFS=hdfs://namenode:9000 \
bde2020/hadoop-namenode
Unable to find image 'bde2020/hadoop-namenode:latest' locally
latest: Pulling from bde2020/hadoop-namenode
3192219afd04: Pull complete
7127a1d8cced: Pull complete
883a89599900: Pull complete
77920a3e82af: Pull complete
92329e81aec4: Pull complete
f373218fec59: Pull complete
aa53513fe997: Pull complete
8b1800105b98: Pull complete
c3a84a3e49c8: Pull complete
a65640a64a76: Pull complete
a29cc756d786: Pull complete
abf352b16046: Pull complete
dddd5a449e99: Pull complete
Digest: sha256:fd74110805132d646cf6f12635efc0919e1fb2ac5bd376c5366272fc261301e
Status: Downloaded newer image for bde2020/hadoop-namenode:latest
d4aa4e4288c5f71009642a02e01a2a3a1a1cca455c525581c0d98c4c162447dd
```

```
hadoop@LAPTOP-RB9PEJTV:~$ docker run -d --name datanode \
--network hadoop-net \
-e CORE_CONF_fs_defaultFS=hdfs://namenode:9000 \
-e SERVICE_PRECONDITION="namenode:9000" \
bde2020/hadoop-datanode
Unable to find image 'bde2020/hadoop-datanode:latest' locally
latest: Pulling from bde2020/hadoop-datanode
3192219afd04: Already exists
7127a1d8cced: Already exists
883a89599900: Already exists
77920a3e82af: Already exists
92329e81aec4: Already exists
f373218fec59: Already exists
aa53513fe997: Already exists
8b1800105b98: Already exists
c3a84a3e49c8: Already exists
a65640a64a76: Already exists
4bf0ae3d5cc8: Pull complete
b91d0b0b68c8: Pull complete
5e185246c615: Pull complete
Digest: sha256:35f899bcbe9f983825a8a3bdc135ed0e8e0eaf3b58f9b08bf257b5e86bae3b47
Status: Downloaded newer image for bde2020/hadoop-datanode:latest
ae8eec62ba572dd934f228a09179876d0d8aa6295568bd220b31179a4797cbb6
```

```
hadoop@LAPTOP-RB9PEJTU:~$ docker exec -it namenode bash
root@d4aa4e4288c5:/# /usr/local/hadoop/bin/hdfs dfs -mkdir /test
bash: /usr/local/hadoop/bin/hdfs: No such file or directory
root@d4aa4e4288c5:/# hdfs dfs -mkdir /test
root@d4aa4e4288c5:/# echo "Hello Hadoop!" > hello.txt
root@d4aa4e4288c5:/# hdfs dfs -put hello.txt /test
2024-10-28 13:13:01,509 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostTrusted = false, remoteHostTrusted = false
root@d4aa4e4288c5:/# hdfs dfs -ls /test
Found 1 items
-rw-r--r--  3 root supergroup          14 2024-10-28 13:13 /test/hello.txt
root@d4aa4e4288c5:/# hdfs dfs -cat /test/hello.txt
2024-10-28 13:18:00,956 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostTrusted = false, remoteHostTrusted = false
Hello Hadoop!
root@d4aa4e4288c5:/# |
```

## Result:

Successfully setup Hadoop environment and ran a Hadoop code in Docker

<b>Ex. No: 8</b>	<b>Public and Private Keys</b>
<b>3.9.2024</b>	

### **Aim:**

To secure an EC2 instance using SSH Keys and Network Access Control.

### **Program:**

#### 1. Generate Private and Public Key Pairs

o Linux/Mac: ■ Create a new directory for key pairs: `mkdir key-pair-labs && cd key-pair-labs`

■ Generate a private key: `openssl genrsa-out snu-privatekey.pem 2048`

■ Generate a public key from the private key: `openssl rsa-in snu-privatekey.pem-pubout-out snu-publickey.pem`

■ Set permissions on the private key: `chmod 400 snu-privatekey.pem`

■ Copy the public key: `cat snu-publickey.pem`

o Paste the public key into the AWS console.

#### 2. Launch an Ubuntu EC2 Instance:

o Launch a new EC2 instance using the public key.

o Choose an appropriate instance type and AMI.

o Allocate an Elastic IP address and attach it to the instance.

#### 3. Login to the Instance: o Use the private key generated in step 1 to SSH into the instance.

o Linux / Mac : `ssh-i snu-privatekey.pem ubuntu@`

#### 4. Edit Security Group: o Open the Security Group settings for the instance.

o Add an inbound rule to allow ICMP traffic (ping) from anywhere:

■ Type: Custom TCP Rule

■ Protocol: ICMP

■ Port Range: All

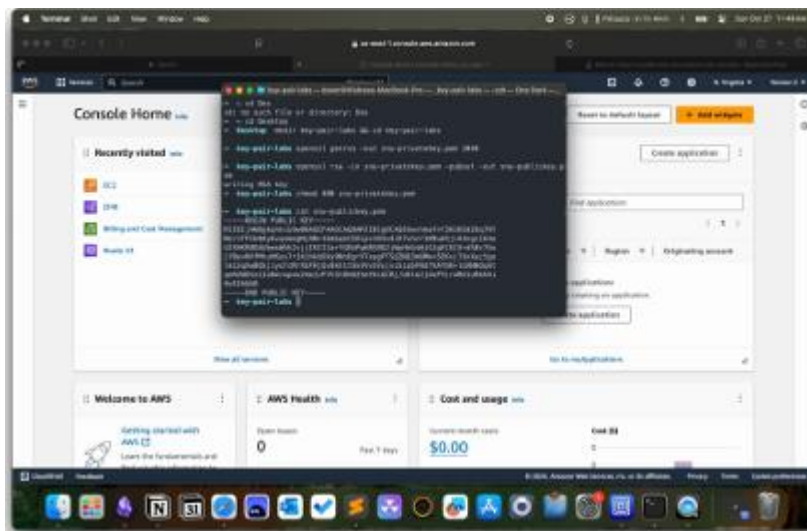
■ Source: 0.0.0.0/0

#### 5. Ping the Instance:

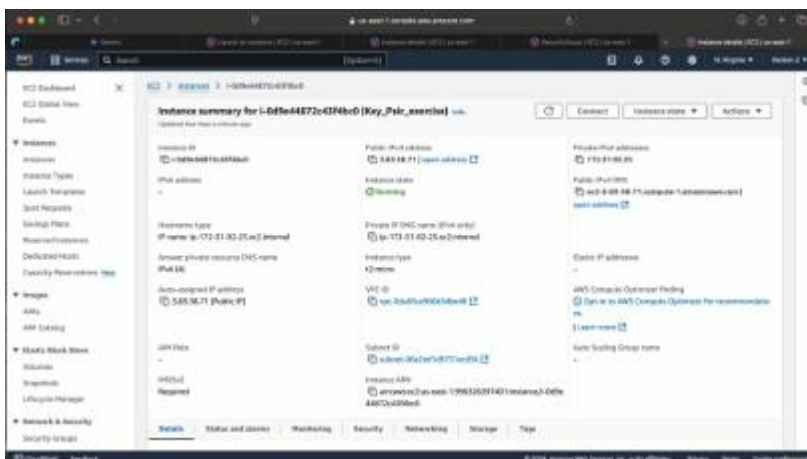
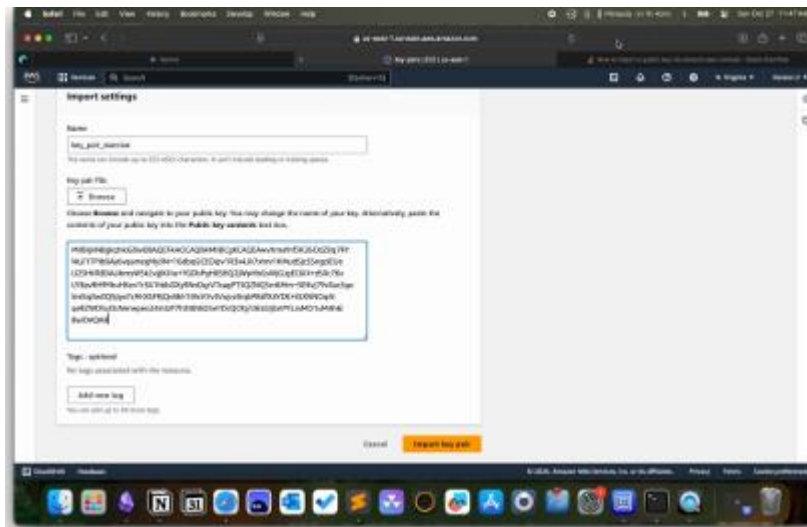
o Ping the public IP address of the instance to verify connectivity.

o You should be able to ping the instance successfully.

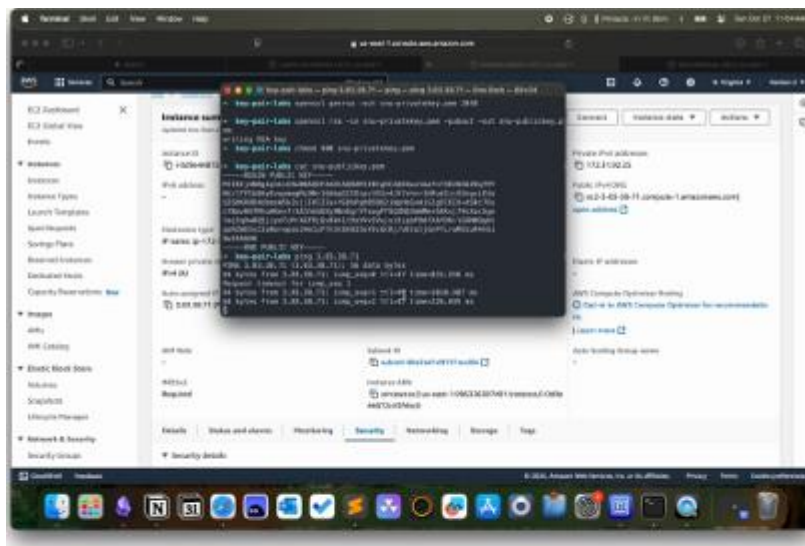
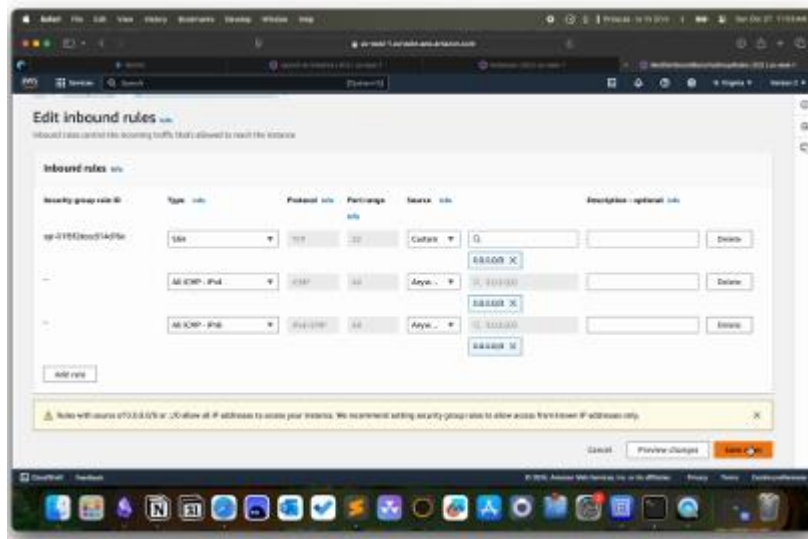
6. Edit Default NACL:
  - Open the Network Access Control List (NACL) settings for the instance.
  - Edit the default NACL to block ICMP traffic.
  - Add an egress rule to deny ICMP traffic:
    - Type: Custom TCP Rule
    - Protocol: ICMP
    - Port Range: All
    - Destination: 0.0.0.0/0
7. Ping the Instance (Again):
  - Ping the public IP address of the instance again.
  - You should not be able to ping the instance this time, as ICMP traffic is now blocked











**Result:**

Successfully created public and private keys and utilized it to secure EC2 instances

<b>Ex. No: 9</b>	<b>EC2</b>
<b>10.8.2024</b>	

### **Aim:**

To create an EC2 instance and Deploy a web server on it using Apache2

### **Program:**

#### 1. Create a Key Pair:

- o Use the AWS Management Console to create a new key pair.
- o Download the private key file (.pem) and save it securely.

#### 2. Launch EC2 Instance:

- o Ubuntu AMI, o Choose an instance type based on your requirements.
- o Configure security groups to allow SSH and HTTP traffic.
- o Launch the instance and provide the private key file during the launch process.

#### 3. Connect to Instance:

- o Use the SSH client to connect to the instance using the public IP address and private key file.

#### 4. Update Package Lists:

- o Run the following command to update the package lists : `sudo apt update`

#### 5. Install Apache2: `sudo apt install apache2`

#### 6. Check Apache Status: `sudo service apache2 status`

#### 7. Test Apache:

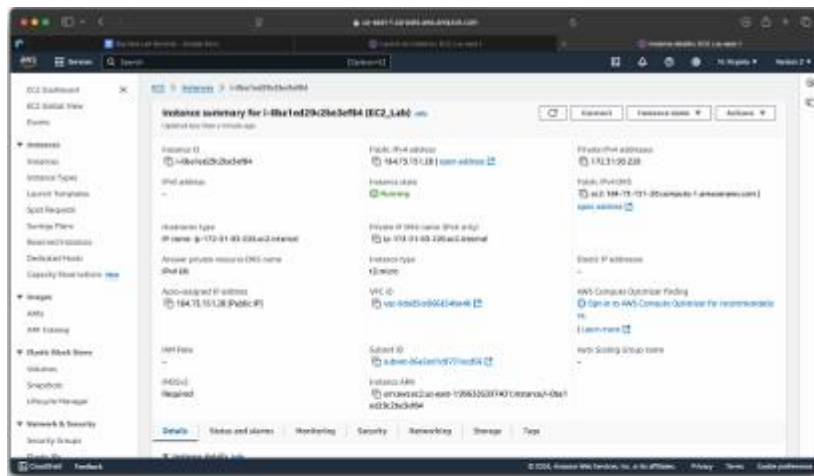
- o Open a web browser and enter the public IP address of the instance.
- o You should see the default Apache2 welcome page.

#### 8. Change the Contents of the Webpage:

- o `cd /var/www/html/` o `sudo chmod 777 index.html`
- o `echo "website" > index.html`

#### 9. Test the website to show new content

### **Output:**



```
key-pair-labs — ubuntu@ip-172-31-93-220: ~ — ssh -i snu-privatekey.pem ubuntu@1...
+ Desktop mkdir key-pair-labs && cd key-pair-labs
+ key-pair-labs openssl genrsa -out snu-privatekey.pem 2048
+ key-pair-labs openssl rsa -in snu-privatekey.pem -pubout -out snu-publickey.pem
writing RSA key
+ key-pair-labs chmod 400 snu-privatekey.pem
+ key-pair-labs cat snu-publickey.pem
-----BEGIN PUBLIC KEY-----
MIIBIjANBgkqhkiG9w0BAQEFAAQCAQ8AMIIBBgKCAQEA5V9u25EVk5nNfs0qMo1
shWlxxwJP9A1MGYOH1QJdInUdk1t4U1NR57e+9YlRwaYyYoj1Men11JcVdHGzTl
R3BbGmPSKhr71lawYKwb3IoCP3s1NPbg/Mq88PJzEdKPPJUP9xisnncNKlNkoTwo
tnpXETWUXdJ/8pWvwaGVQsUWHjN8YHbL3/oRoVndhMps9BGn4cqRV1DMpQ5QwkkI
+IapDy4sJ3xp24fha0pnbGzDYJ5IkarqLvt4TVISuNqSuxwJ0qoHGZhkmold8Tcy
I873xrMgcAmxwRfQ3FAhD39dRroatZPvvhDL0TnJEFVee1RpVrs5Q0yFCRva08Y
uwIDAQAB
-----END PUBLIC KEY-----
+ key-pair-labs ssh -i snu-privatekey.pem ubuntu@184.73.151.28
The authenticity of host '184.73.151.28 (184.73.151.28)' can't be established.
ED25519 key fingerprint is SHA256:0447y6SkAs3URJwfwJ4XN8w9P1G0vUeMPPMA040dgfnY.
This key is not known by any other names.
Are you sure you want to continue connecting (yes/no/[fingerprint])? yes
```

```
key-pair-labs — ubuntu@ip-172-31-93-220: ~ — ssh -i snu-privatekey.pem ubuntu@1...
ubuntu@ip-172-31-93-220:~$ sudo apt update
Hit:1 http://us-east-1.ec2.archive.ubuntu.com/ubuntu noble InRelease
Get:2 http://us-east-1.ec2.archive.ubuntu.com/ubuntu noble-updates InRelease [12
6 kB]
Get:3 http://us-east-1.ec2.archive.ubuntu.com/ubuntu noble-backports InRelease [
126 kB]
Get:4 http://security.ubuntu.com/ubuntu noble-security InRelease [126 kB]
Get:5 http://us-east-1.ec2.archive.ubuntu.com/ubuntu noble/universe amd64 Packag
es [15.0 MB]
Get:6 http://us-east-1.ec2.archive.ubuntu.com/ubuntu noble/universe Translation-
en [5982 kB]
Get:7 http://security.ubuntu.com/ubuntu noble-security/main amd64 Packages [433
kB]
Get:8 http://us-east-1.ec2.archive.ubuntu.com/ubuntu noble/universe amd64 Compon
ents [3871 kB]
Get:9 http://us-east-1.ec2.archive.ubuntu.com/ubuntu noble/universe amd64 c-n-f
Metadata [301 kB]
Get:10 http://us-east-1.ec2.archive.ubuntu.com/ubuntu noble/multiverse amd64 Pac
kages [269 kB]
Get:11 http://us-east-1.ec2.archive.ubuntu.com/ubuntu noble/multiverse Translati
on-en [118 kB]
Get:12 http://us-east-1.ec2.archive.ubuntu.com/ubuntu noble/multiverse amd64 Com
ponents [35.0 kB]
Get:13 http://us-east-1.ec2.archive.ubuntu.com/ubuntu noble/multiverse amd64 c-n
```

```
key-pair-labs — ubuntu@ip-172-31-93-220: ~ — ssh -i snu-privatekey.pem ubuntu@1...
ubuntu@ip-172-31-93-220:~$ sudo apt install apache2
Reading package lists... Done
Building dependency tree... Done
Reading state information... Done
The following additional packages will be installed:
  apache2-bin apache2-data apache2-utils libapr1t64 libaprutil1-dbd-sqlite3
  libaprutil1-ldap libaprutil1t64 liblua5.4-0 ssl-cert
Suggested packages:
  apache2-doc apache2-suexec-pristine | apache2-suexec-custom www-browser
The following NEW packages will be installed:
  apache2 apache2-bin apache2-data apache2-utils libapr1t64
  libaprutil1-dbd-sqlite3 libaprutil1-ldap libaprutil1t64 liblua5.4-0 ssl-cert
0 upgraded, 10 newly installed, 0 to remove and 28 not upgraded.
Need to get 2084 kB of archives.
After this operation, 8094 kB of additional disk space will be used.
Do you want to continue? [Y/n] Y
Get:1 http://us-east-1.ec2.archive.ubuntu.com/ubuntu noble-updates/main amd64 libapr1t64 amd64 1.7.2-3.1ubuntu0.1 [108 kB]
Get:2 http://us-east-1.ec2.archive.ubuntu.com/ubuntu noble/main amd64 libaprutil1t64 amd64 1.6.3-1.1ubuntu7 [91.9 kB]
Get:3 http://us-east-1.ec2.archive.ubuntu.com/ubuntu noble/main amd64 libaprutil1-dbd-sqlite3 amd64 1.6.3-1.1ubuntu7 [11.2 kB]
Get:4 http://us-east-1.ec2.archive.ubuntu.com/ubuntu noble/main amd64 libaprutil1-ldap amd64 1.6.3-1.1ubuntu7 [9116 B]
```



```

key-pair-labs — ubuntu@ip-172-31-93-220: ~ — ssh -i smu-privatekey.pem ubuntu@1...
No services need to be restarted.

No containers need to be restarted.

No user sessions are running outdated binaries.

No VM guests are running outdated hypervisor (qemu) binaries on this host.
ubuntu@ip-172-31-93-220:~$ sudo service apache2 status
● apache2.service - The Apache HTTP Server
   Loaded: loaded (/usr/lib/systemd/system/apache2.service; enabled; preset:
   Active: active (running) since Sun 2024-10-27 07:58:20 UTC; 14s ago
     Docs: https://httpd.apache.org/docs/2.4/
   Main PID: 2060 (apache2)
      Tasks: 55 (limit: 1130)
    Memory: 5.4M (peak: 5.7M)
       CPU: 32ms
    CGroup: /system.slice/apache2.service
            └─2060 /usr/sbin/apache2 -k start
              └─2063 /usr/sbin/apache2 -k start
                └─2064 /usr/sbin/apache2 -k start

Oct 27 07:58:20 ip-172-31-93-220 systemd[1]: Starting apache2.service - The Apache HTTP Server
Oct 27 07:58:20 ip-172-31-93-220 systemd[1]: Started apache2.service - The Apache HTTP Server
lines 1-15/15 (END)

```

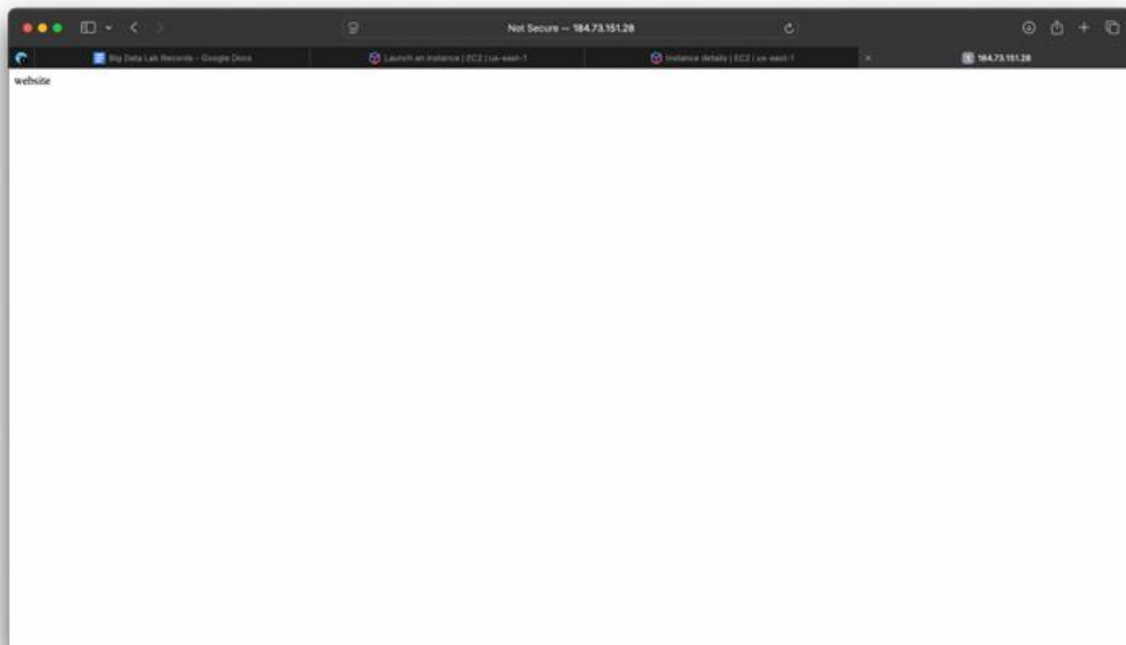




```
key-pair-labs — ubuntu@ip-172-31-93-220: /var/www/html — ssh -i snu-privatekey.p...
21 of these updates are standard security updates.
To see these additional updates run: apt list --upgradable

Enable ESM Apps to receive additional future security updates.
See https://ubuntu.com/esm or run: sudo pro status

Last login: Sun Oct 27 07:57:17 2024 from 49.43.251.85
ubuntu@ip-172-31-93-220:~$ cd /var/www/html/
ubuntu@ip-172-31-93-220:/var/www/html$ sudo chmod 777 index.html
Command 'suo' not found, did you mean:
  command 'su' from deb util-linux (2.39.3-9ubuntu6.1)
  command 'sur' from deb subtle (0.11.3224-xi-2.2build5)
  command 'sup' from deb sup (20100519-3)
  command 'sui' from deb hxttools (20231101-1)
  command 'sudo' from deb sudo (1.9.14p2-1ubuntu1)
  command 'sudo' from deb sudo-ldap (1.9.14p2-1ubuntu1)
  command 'sumo' from deb sumo (1.18.0+dfsg-3build2)
  command 'zuo' from deb zuo (1.9-1)
  command 'sum' from deb coreutils (9.4-2ubuntu2)
Try: sudo apt install <deb name>
ubuntu@ip-172-31-93-220:/var/www/html$ sudo chmod 777 index.html
ubuntu@ip-172-31-93-220:/var/www/html$ echo "website" > index.html
ubuntu@ip-172-31-93-220:/var/www/html$
```



## Result:

Successfully deployed an EC2 instance and tested using Apache Webservices.

<b>Ex. No: 10</b>	<b>Route53</b>
<b>1.10.2024</b>	

### **Aim:**

To set up a web server on an AWS EC2 instance and configure a domain name using GoDaddy and Route 53

### **Program:**

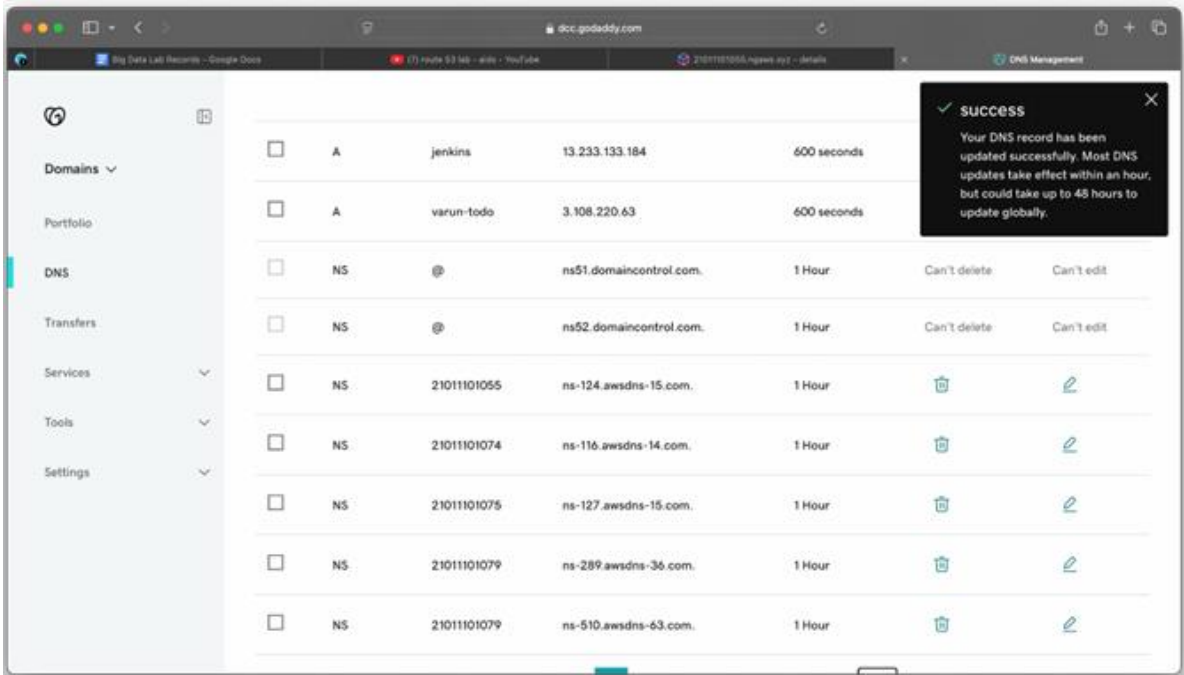
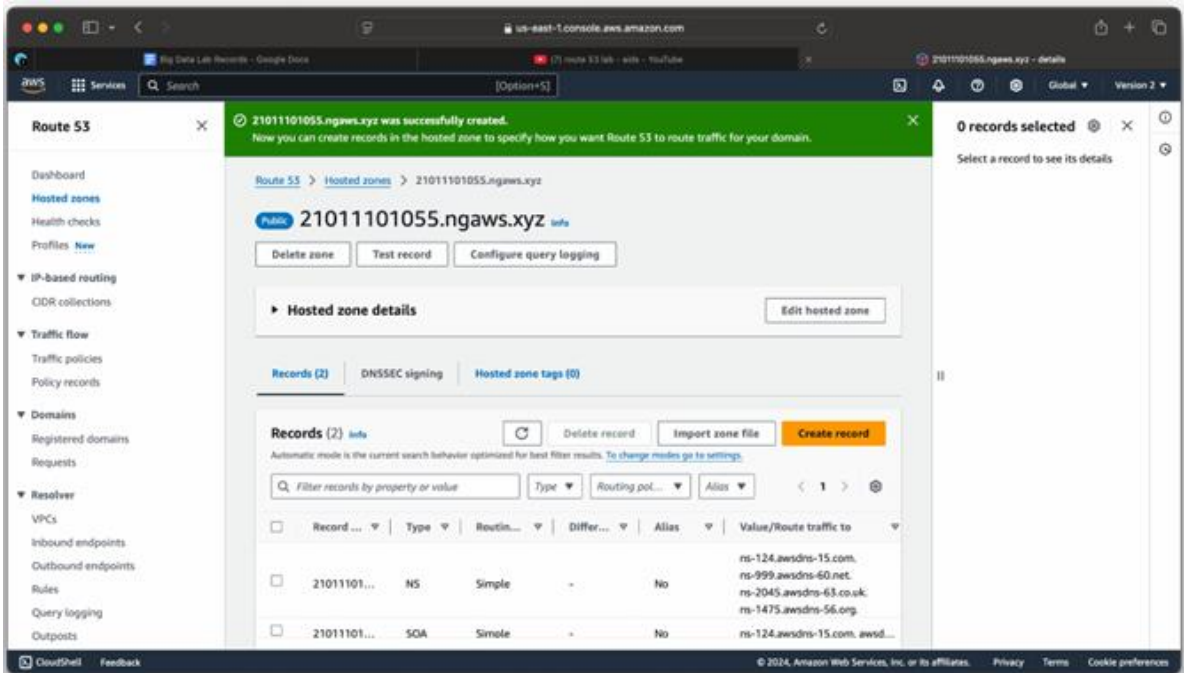
1. Login to GoDaddy and Create a subdomain:
  - o Use your 12-digit registration number.
  - o Example: 21100101001.ngaws.xyz
2. Create a Hosted Zone in Route 53:
  - o Navigate to the Route 53 service in the AWS console.
  - o Create a new Hosted Zone for your subdomain.
3. Get the Name Server information:
  - o Retrieve the name server information from the Route 53 dashboard.
  - o Example: ns-1234.awsdns-12.org, [invalid URL removed]
4. Update the NS record in GoDaddy:
  - o Log in to the GoDaddy portal.
  - o Navigate to the DNS settings for your subdomain.
  - o Update the NS records with the name servers obtained from Route 53.
5. Create an EC2 instance:
  - o Launch an EC2 instance with an elastic public IP address.
6. Install Apache:
  - o Connect to the EC2 instance using SSH.
  - o Update the package lists: `sudo apt update`
  - o Install Apache: `sudo apt install apache2`
7. Configure Apache:
  - o Edit the Apache configuration file: `sudo nano /etc/apache2/sites-available/000-default.conf`
  - o Modify the DocumentRoot and ServerName directives to point to your desired web content directory and domain name.

- o Save the configuration file and restart Apache: `sudo systemctl restart apache2`

8. Create an A record in Route 53:

- o Navigate to your Hosted Zone in Route 53.
- o Create a new A record with the following details:
  - Name: @ (for the root domain)
  - Value: The public IP address of your EC2 instance
  - TTL: 3600 (1 hour)

**Output:**



us-east-1.console.aws.amazon.com

EC2 Dashboard

EC2 Global View

Events

Instances

Instance Types

Launch Templates

Spot Requests

Savings Plans

Reserved Instances

Dedicated Hosts

Capacity Reservations

Images

AMIs

AMI Catalog

Elastic Block Store

Volumes

Snapshots

Lifecycle Manager

Network & Security

Security Groups

Instance summary for i-09a90d79a-0d79acc628dbd

Updated less than a minute ago

Public IPv4 address copied

Instance ID: i-09a90d79acc628dbd

IPv6 address: -

Hostname type: IP name: ip-172-31-45-111.ec2.internal

Answer private resource DNS name: IPv4 (A)

Auto-assigned IP address: 54.162.240.174 [Public IP]

IAM Role: -

IMDSv2: Required

Instance state: Running

Private IP DNS name (IPv4 only): ip-172-31-45-111.ec2.internal

Instance type: t2.micro

VPC ID: vpc-06a85ce966b34be48

Subnet ID: subnet-08c82bc8826320fef

Instance ARN: arn:aws:ec2:us-east-1:996326397401:instance/i-09a90d79acc628dbd

Private IPv4 addresses: 172.31.45.111

Public IPv4 DNS: ec2-54-162-240-174.compute-1.amazonaws.com

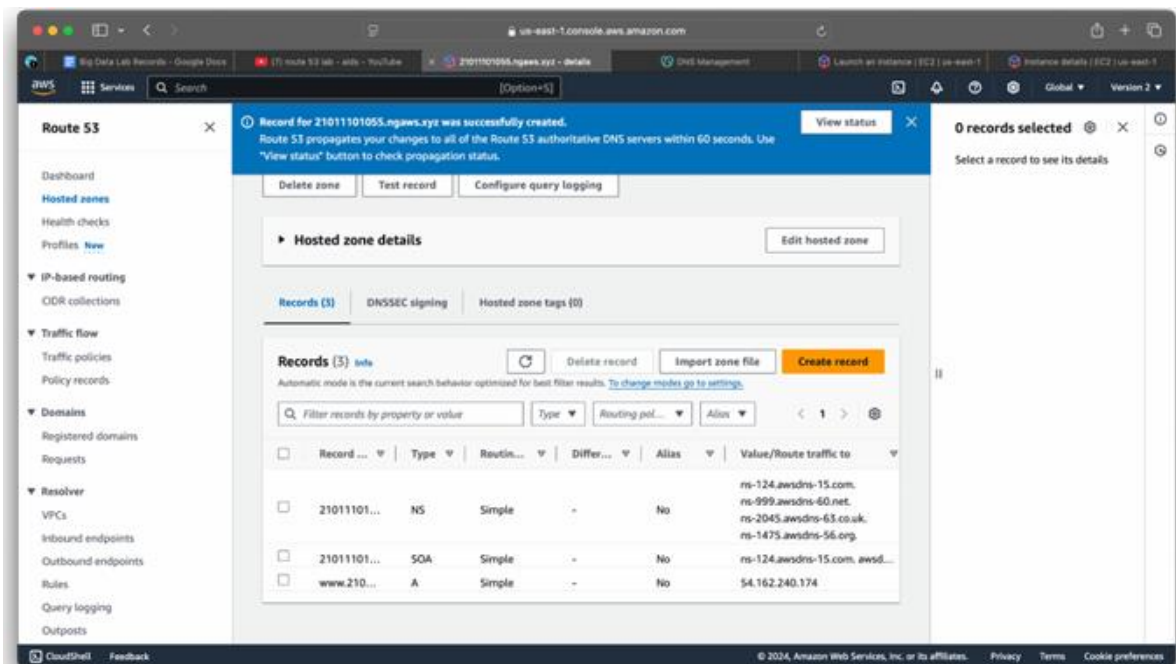
Elastic IP addresses: -

AWS Compute Optimizer finding: Opt-in to AWS Compute Optimizer for recommendations. Learn more

Auto Scaling Group name: -

Details | Status and alarms | Monitoring | Security | Networking | Storage | Tags

© 2024 Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences



```
daver — daver@Vishwas-MacBook-Pro — — -zsh — One Dark — 80x24

[+] ~ links www.21011101055.ngaws.xyz
zsh: command not found: links
[+] ~ nslookup www.21011101055.ngaws.xyz
Server:      172.16.0.2
Address:     172.16.0.2#53

Non-authoritative answer:
Name:   www.21011101055.ngaws.xyz
Address: 54.162.240.174

[+] ~
```

## Result:

Successfully implemented a Route53 DNS Lookup on AWS.

<b>Ex. No: 11</b>	<b>IAM</b>
<b>8.10.2024</b>	

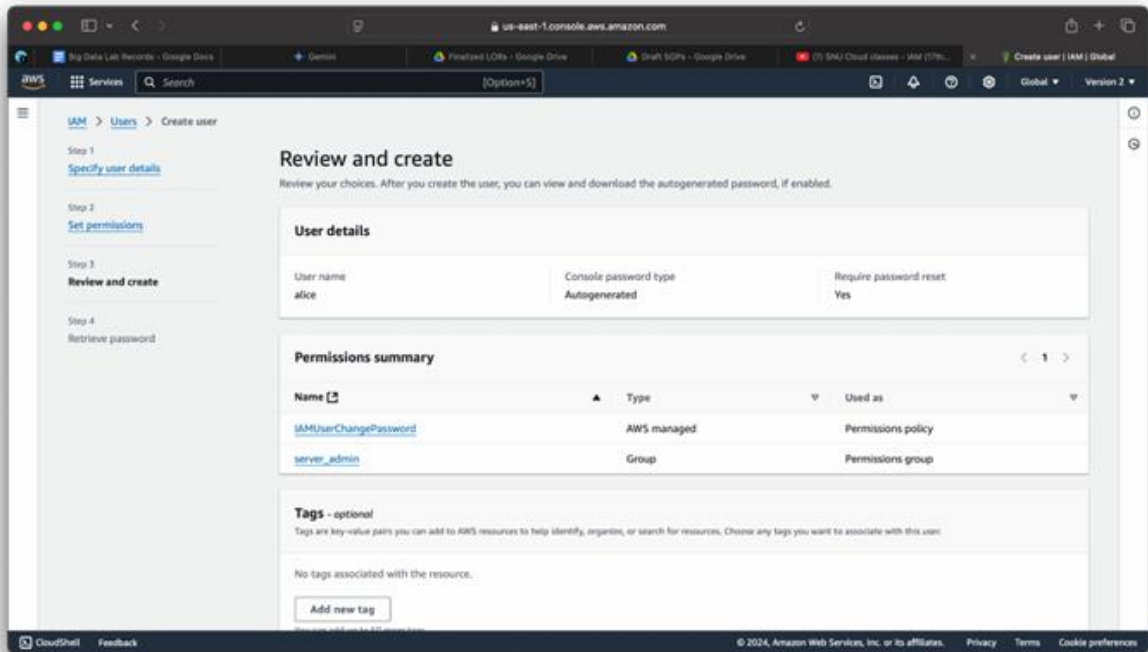
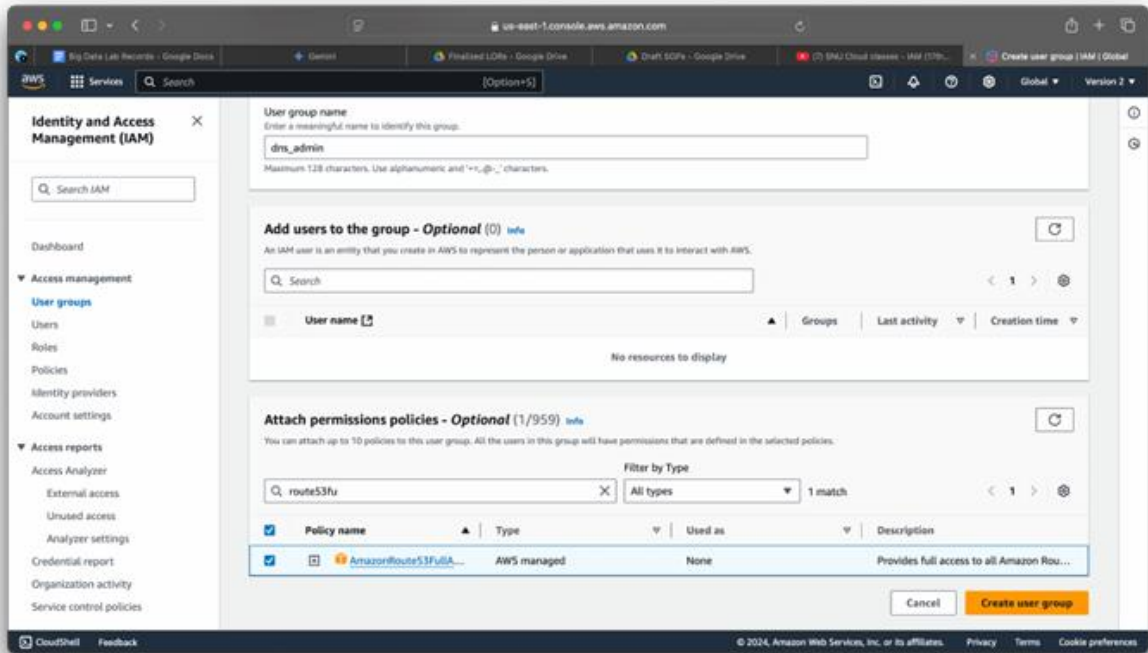
**Aim:**

To implement IAM user access in AWS Console.

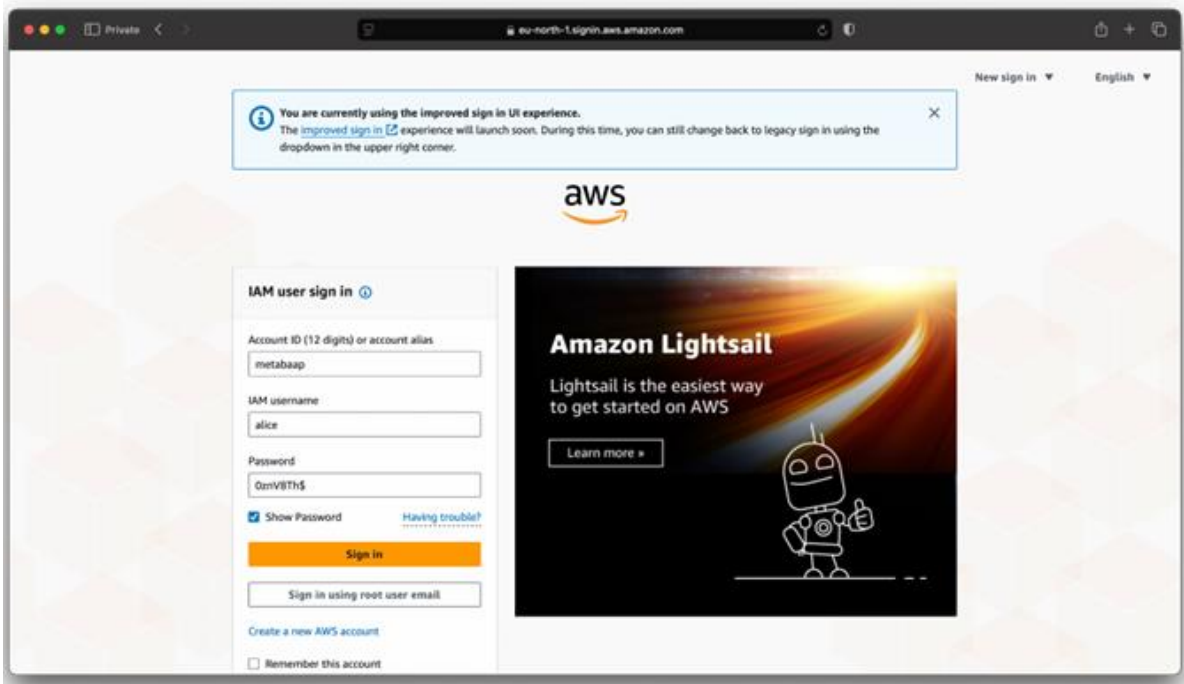
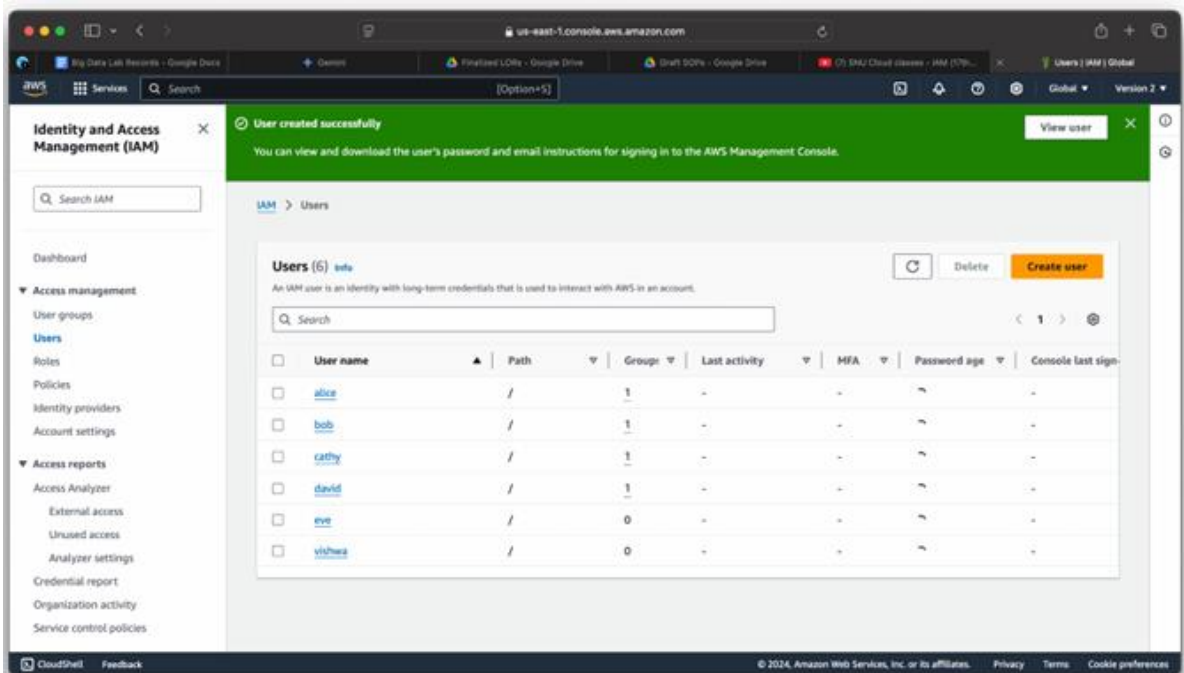
**Program:**

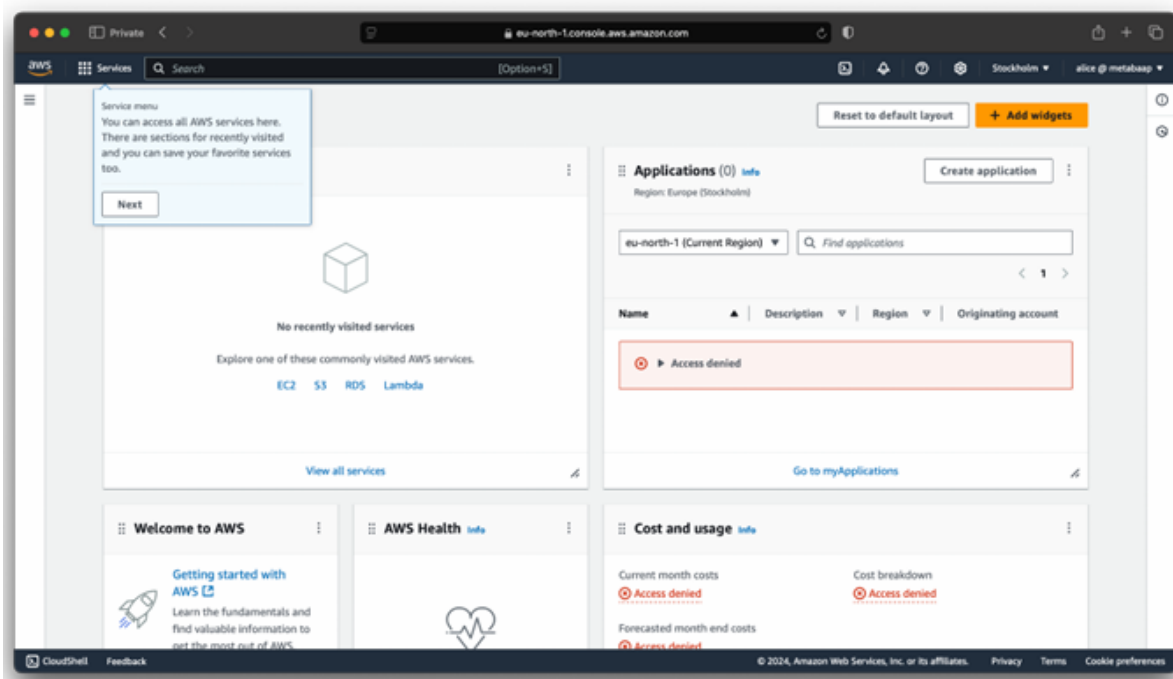
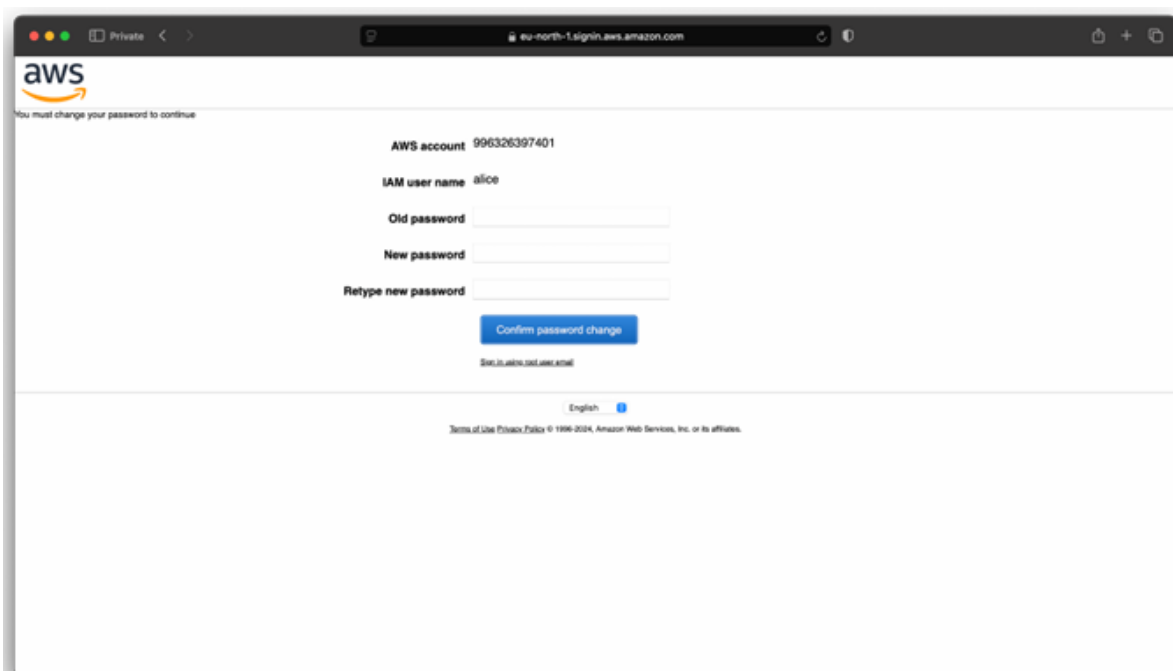
1. Create User Groups server\_admin and dns\_admin. Give them full access to EC2 Service and Route53 services.
2. Create users, alice, bob, cathay and david. Set passwords for them
3. Add alice and bob to user group server\_adin and cathy and david to user group dns\_admin.
4. Create user eve and give him Billing access.
5. Create user hadoop and give full (admin) access to the services.
6. Create an alias name for your account.
7. Use the alias URL to login to your account instead of the account ID.

**Output:**









## Result:

Successfully implemented IAM labs in AWS