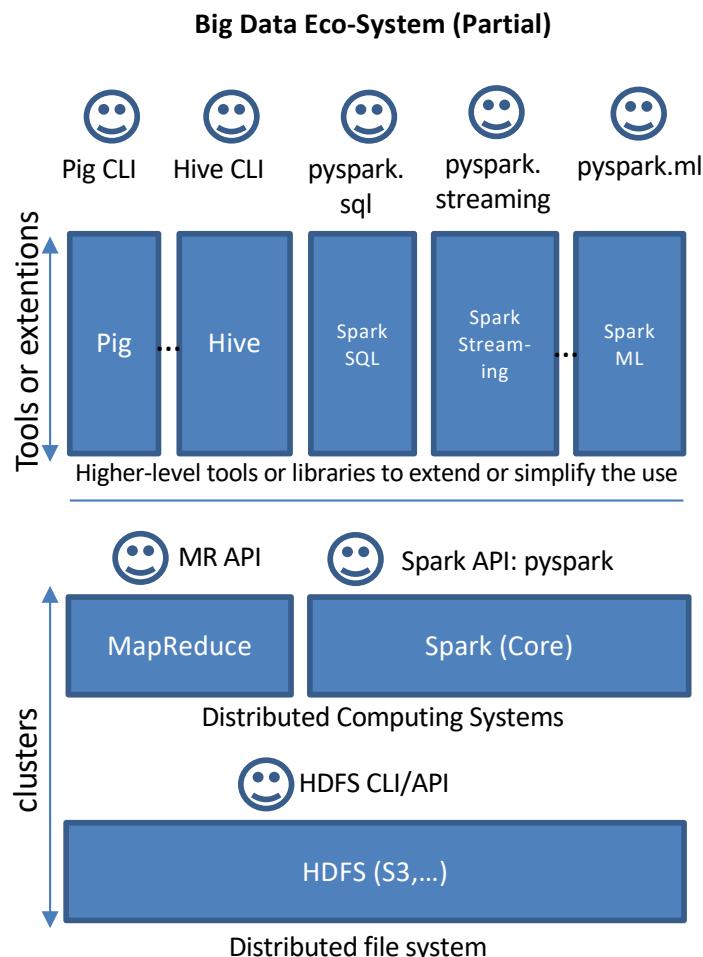
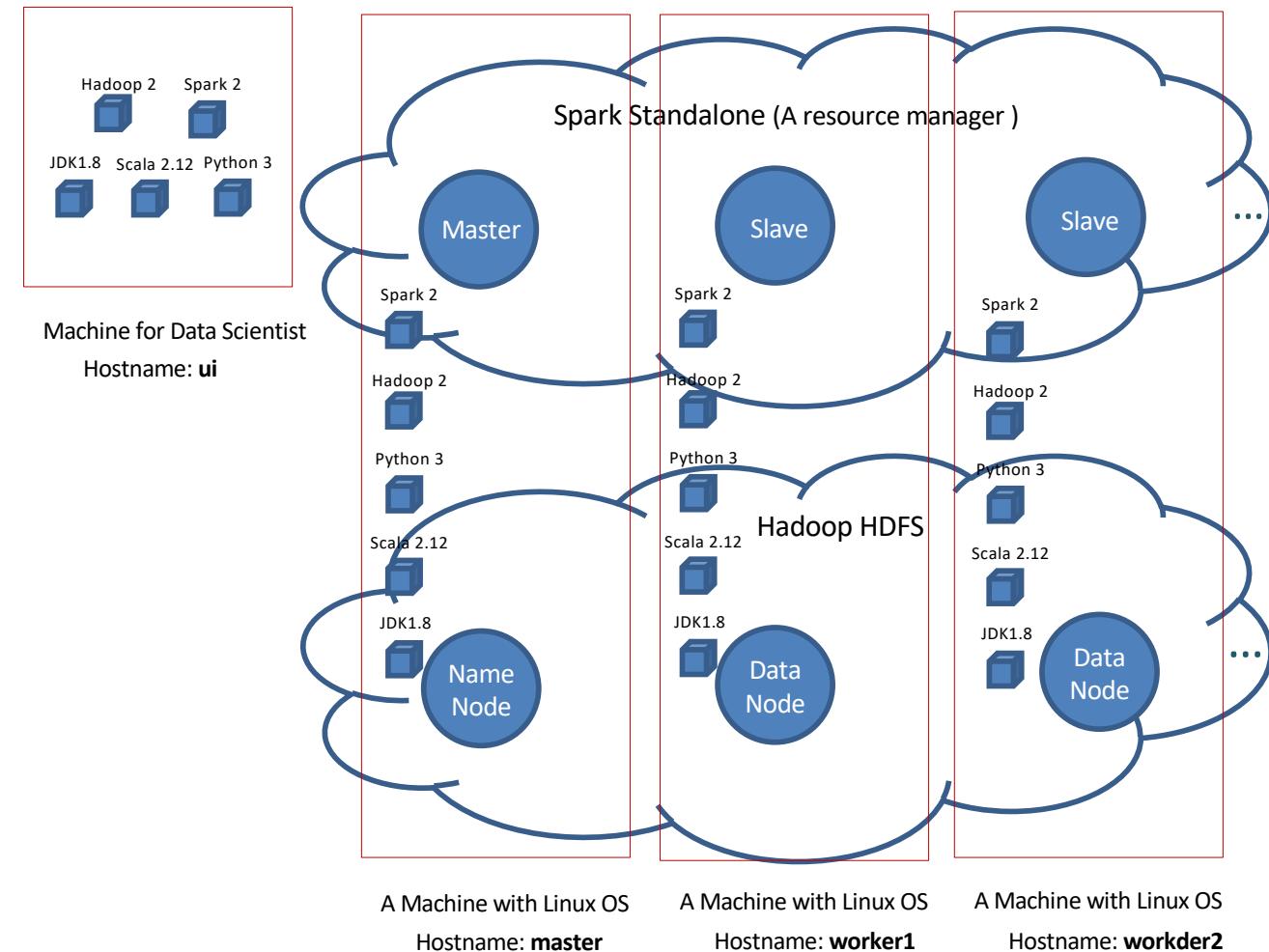


Production Environment Setup

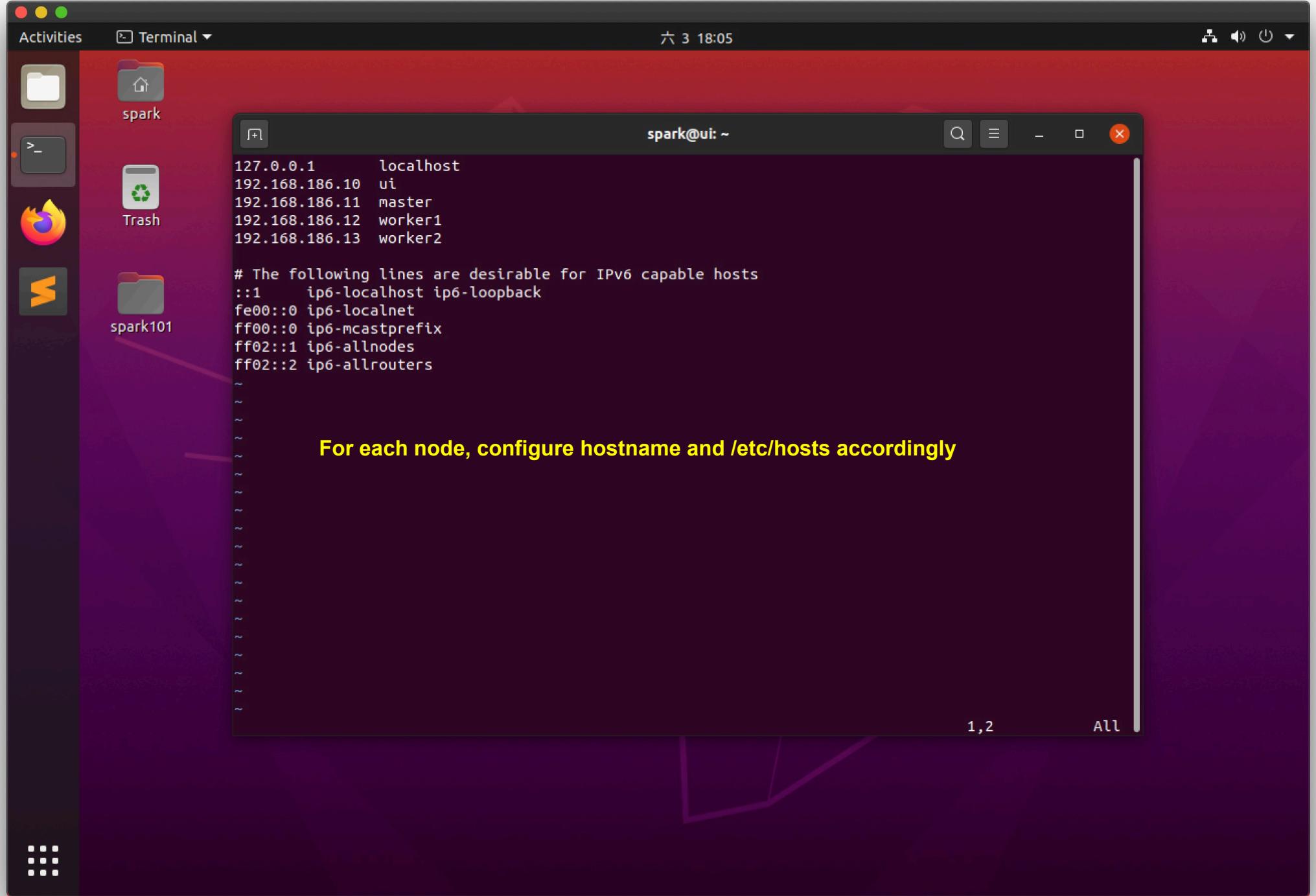


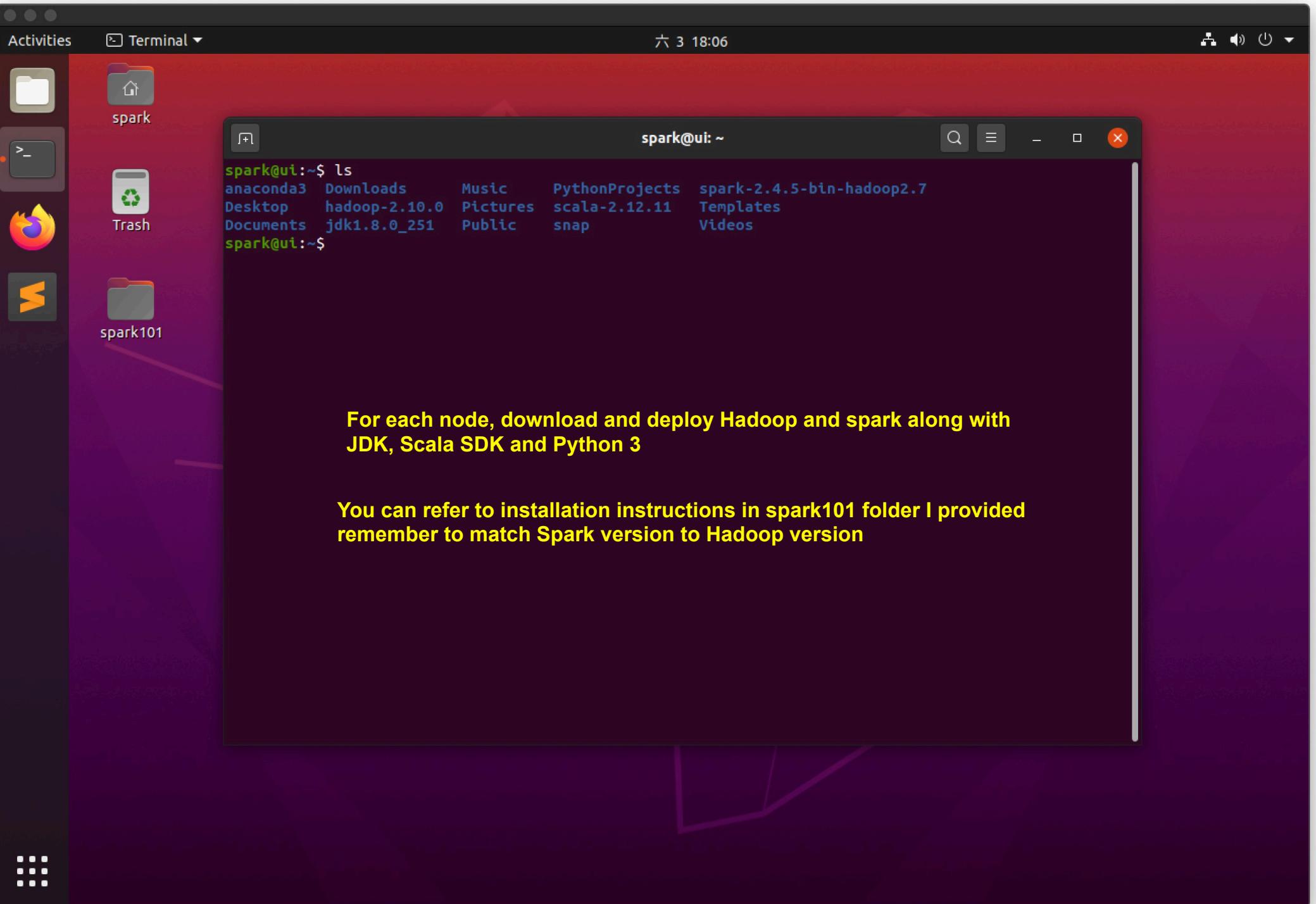
Data Scientists use Python, Scala API, Hadoop API to access



Note that:

1. Make sure each machine in Spark/HDFS cluster has static IP address
2. Make sure the /etc/hosts has records for all nodes and synchronized on each node
3. The firewall on each node need to allow access to required ports
(can turn off the firewall temporarily during the setup and put it on after the test is done)





Activities Terminal ▾ 六 3 18:06

spark@ui: ~

```
# sleep 10; alert
alias alert='notify-send --urgency=low -i "$( [ $? = 0 ] && echo terminal || echo error)" "$(history|tail -n1|sed -e '\''$s/^\\s*[0-9]\\+\\s*//;s/[;&|]\\s*alert$//\\''")"

# Alias definitions.
# You may want to put all your additions into a separate file like
# ~/.bash_aliases, instead of adding them here directly.
# See /usr/share/doc/bash-doc/examples in the bash-doc package.

if [ -f ~/.bash_aliases ]; then
    . ~/.bash_aliases
fi

# enable programmable completion features (you don't need to enable
# this, if it's already enabled in /etc/bash.bashrc and /etc/profile
# sources /etc/bash.bashrc).
if ! shopt -oq posix; then
    if [ -f /usr/share/bash-completion/bash_completion ]; then
        . /usr/share/bash-completion/bash_completion
    elif [ -f /etc/bash_completion ]; then
        . /etc/bash_completion
    fi
fi

export JAVA_HOME=/home/spark/jdk1.8.0_251
export PATH=$JAVA_HOME/bin:$PATH

export SCALA_HOME=/home/spark/scala-2.12.11
export PATH=$SCALA_HOME/bin:$PATH

export ANACONDA_HOME=/home/spark/anaconda3
export PATH=$ANACONDA_HOME/bin:$PATH

export HADOOP_HOME=/home/spark/hadoop-2.10.0
export PATH=$HADOOP_HOME/bin:$PATH

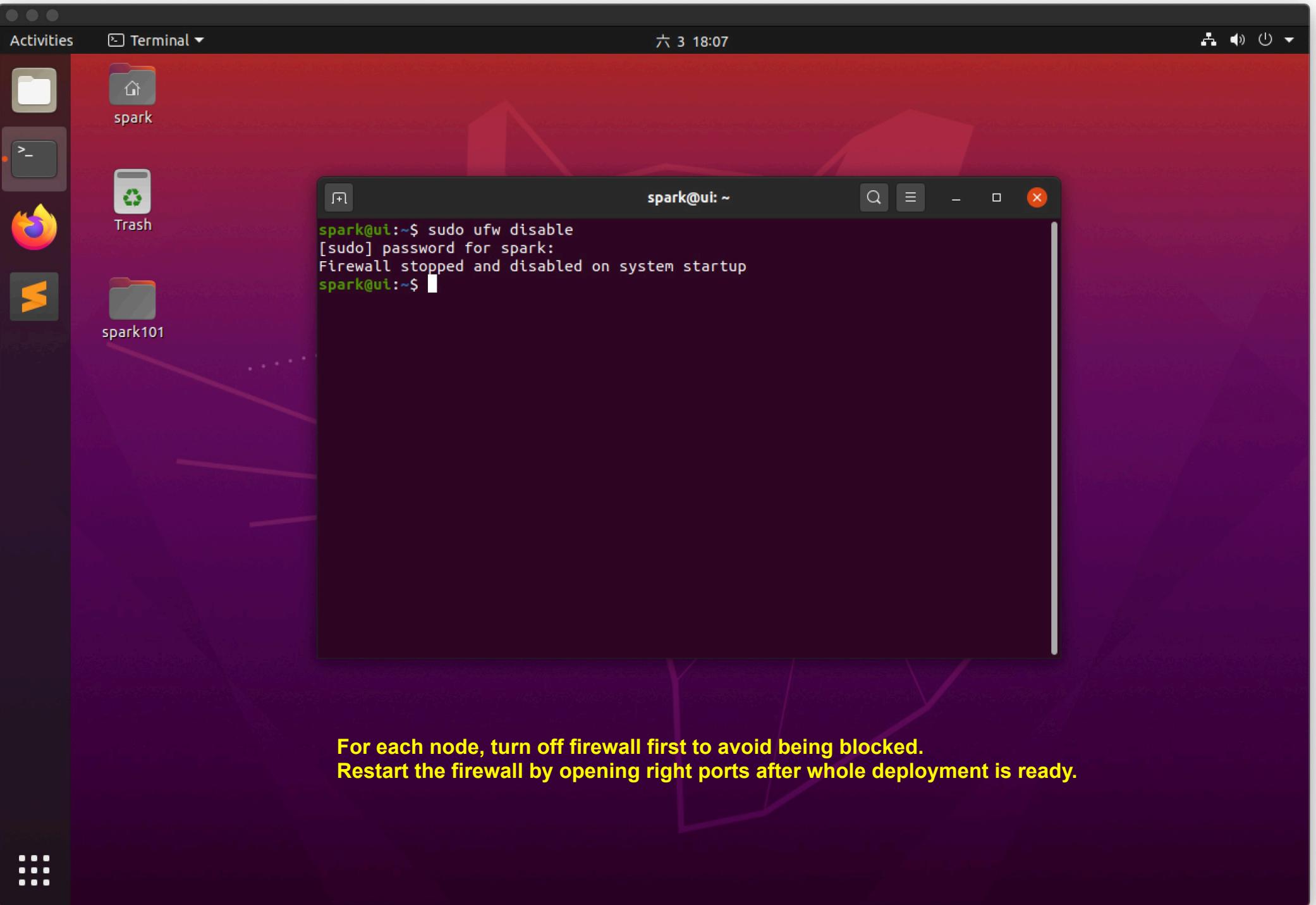
export SPARK_HOME=/home/spark/spark-2.4.5-bin-hadoop2.7
export PATH=$SPARK_HOME/bin:$PATH
export PYSPARK_PYTHON=python

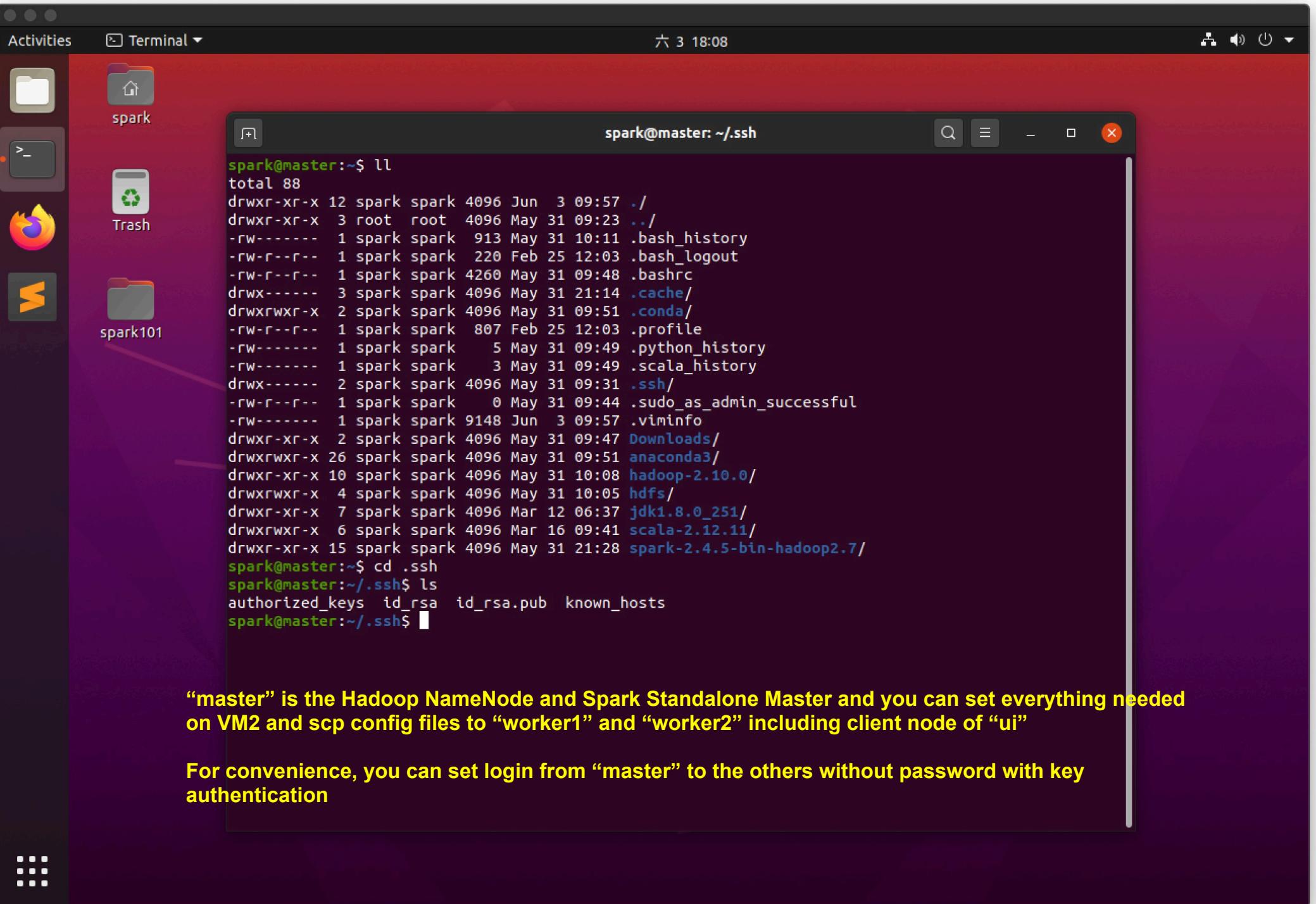
# use ipython
export PYSPARK_DRIVER_PYTHON=ipython
```

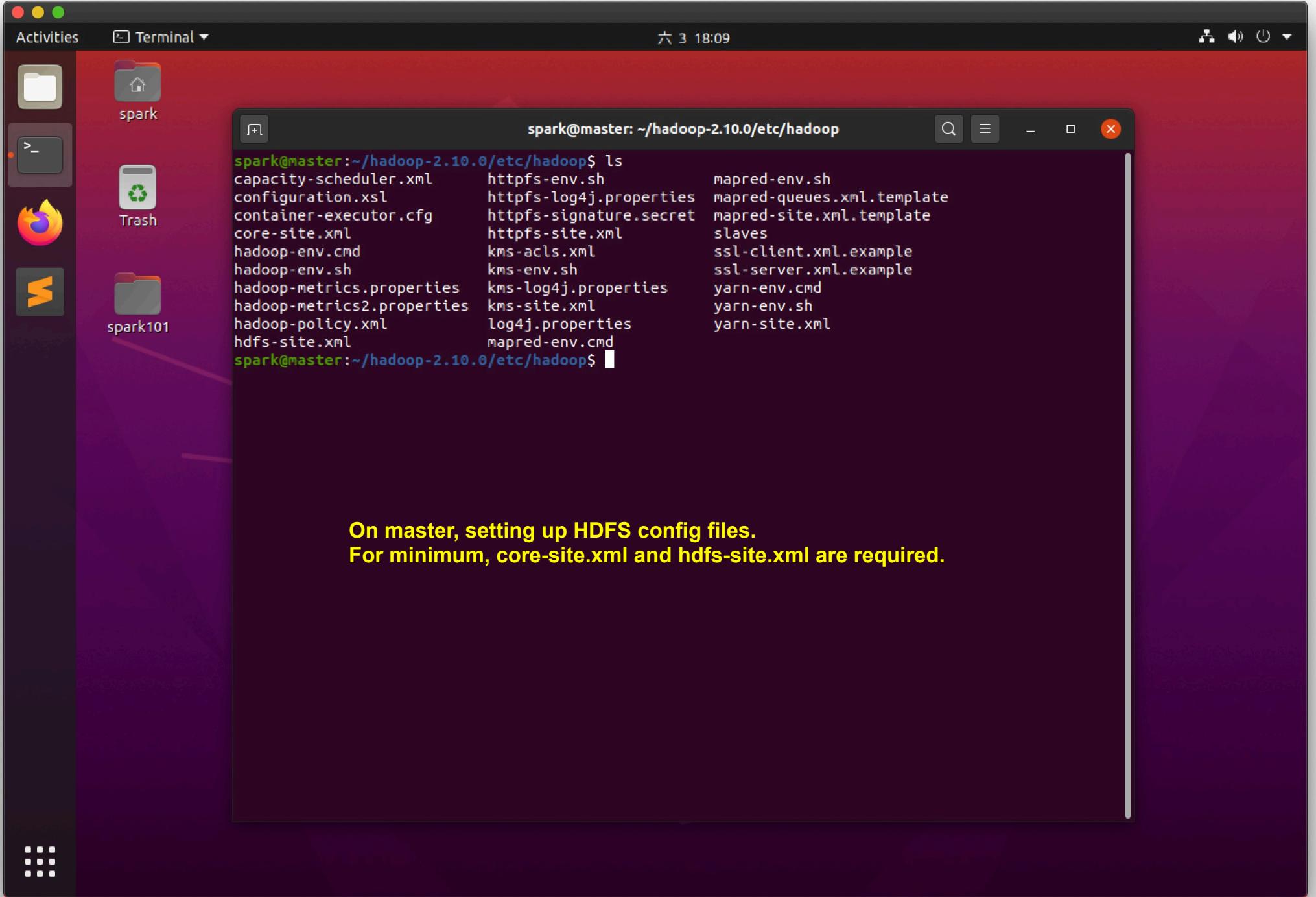
115 × 43

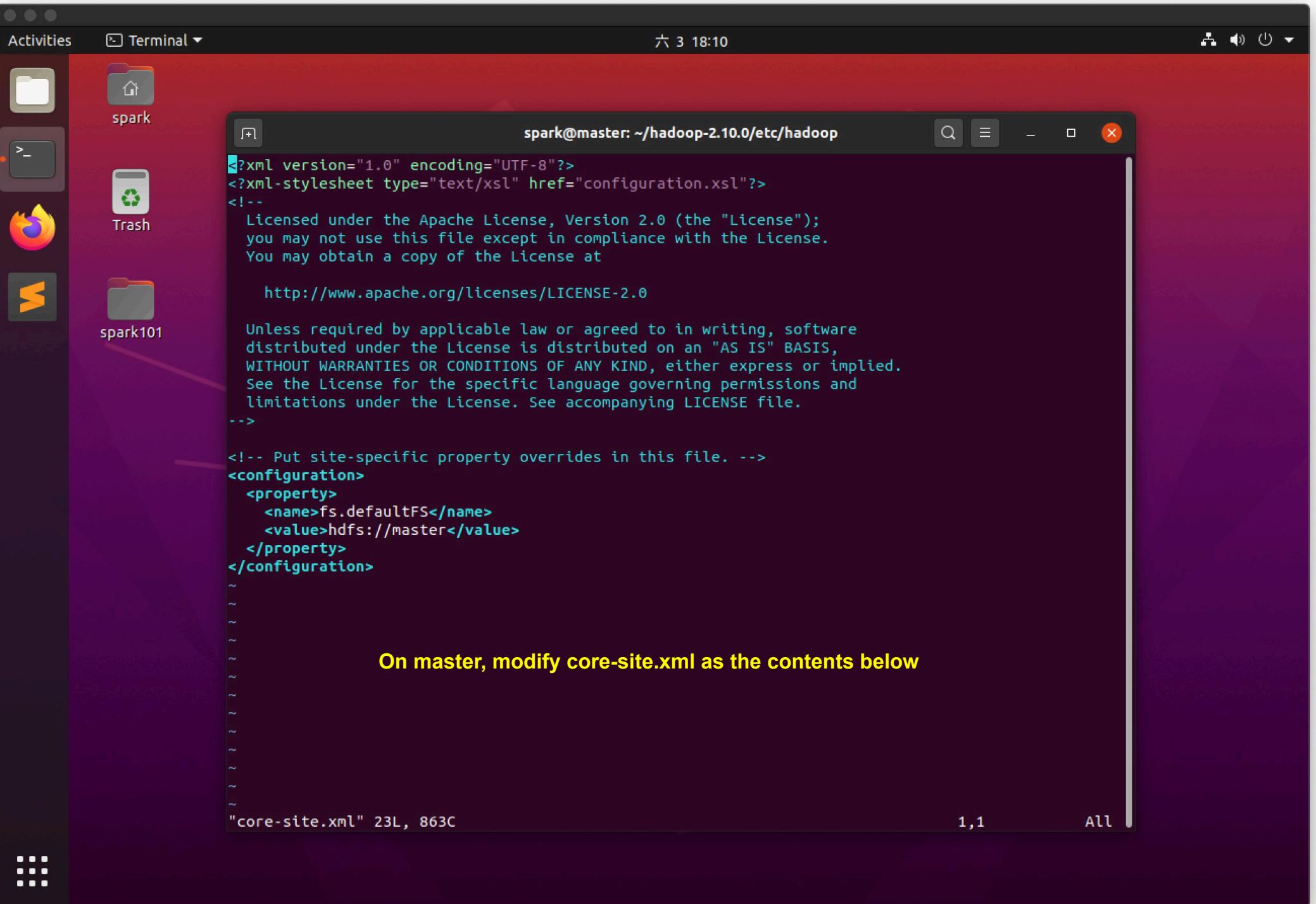
For each node, set relevant environmental variables in `~/.bashrc`

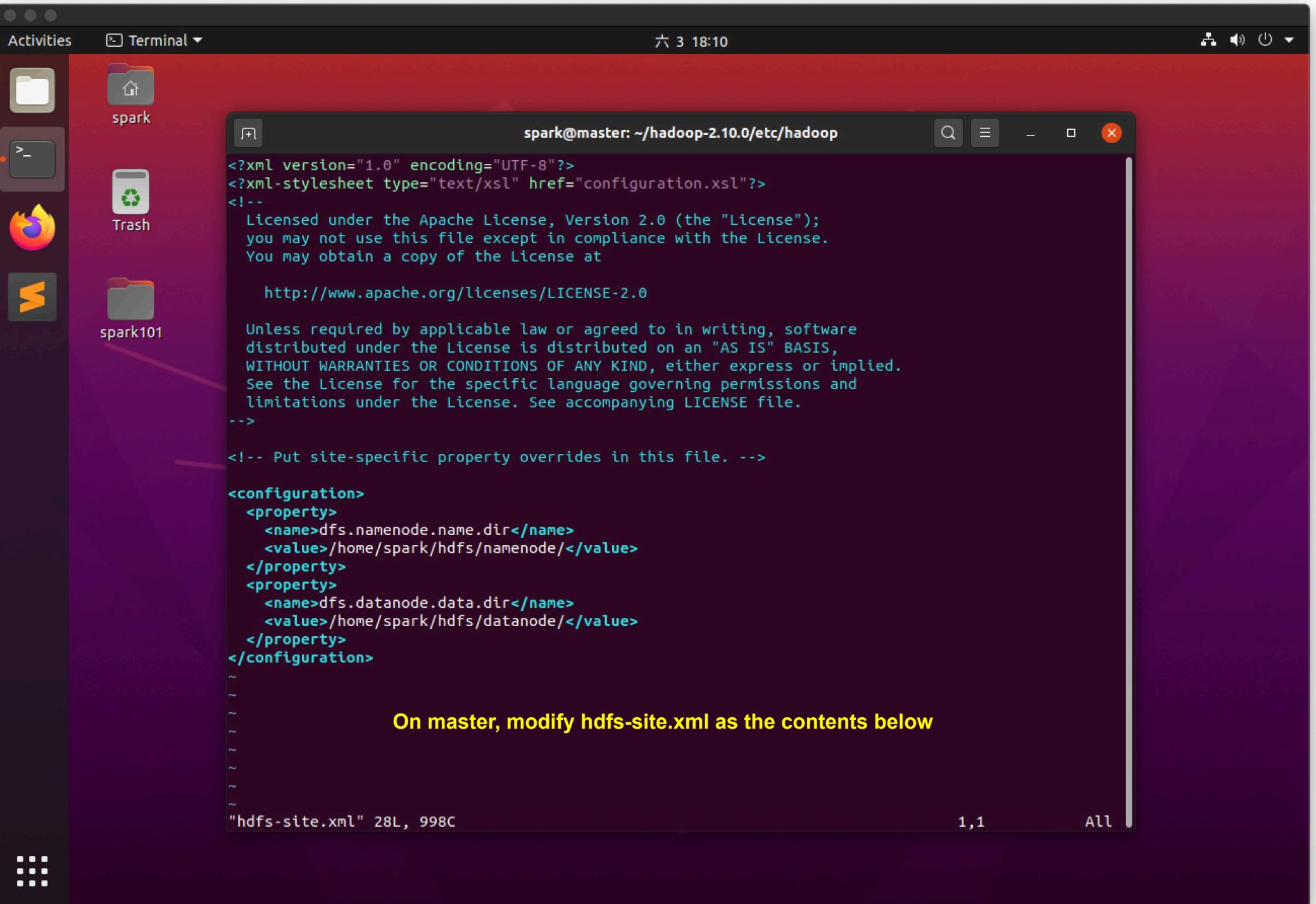
136,1 Bot

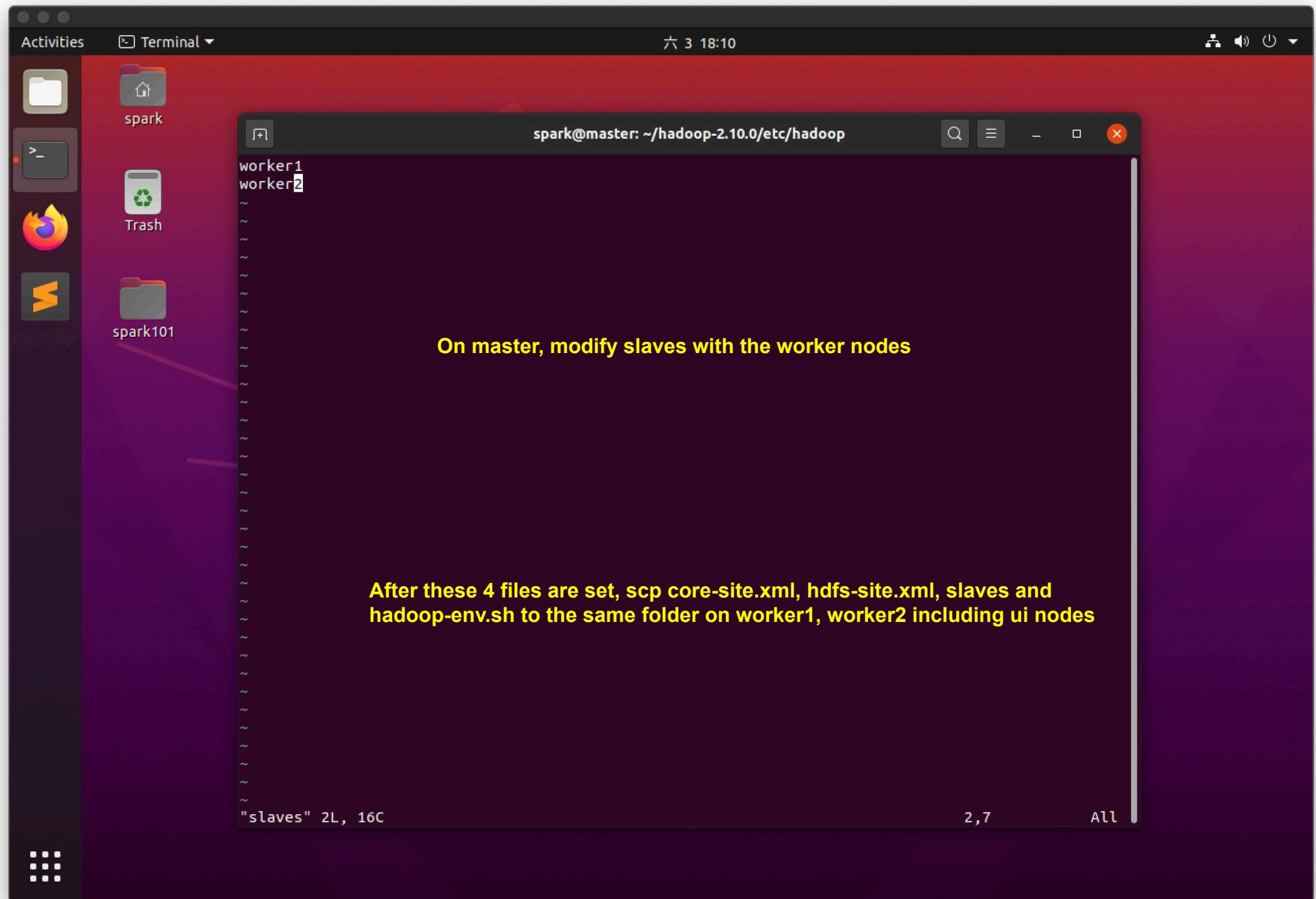


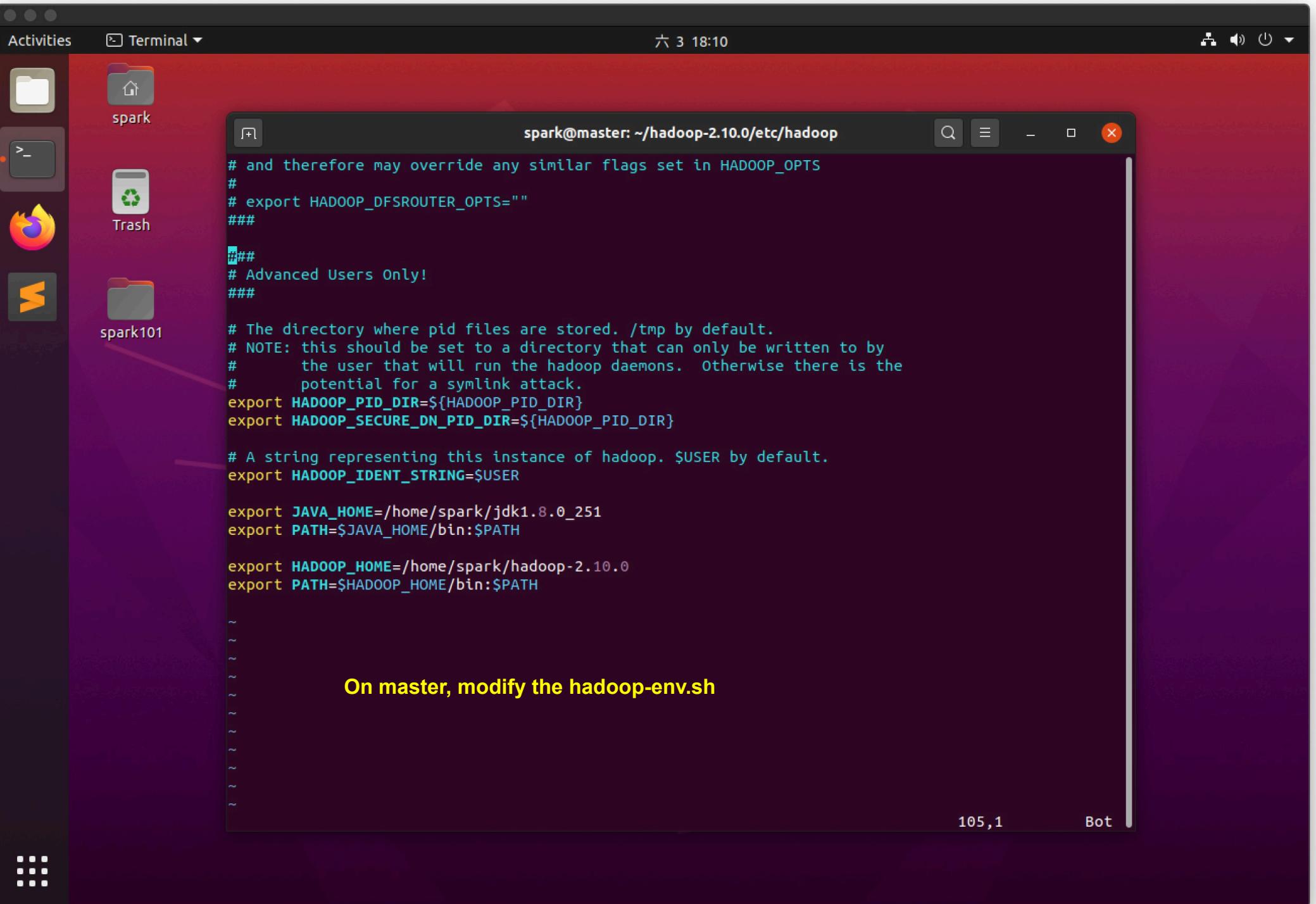


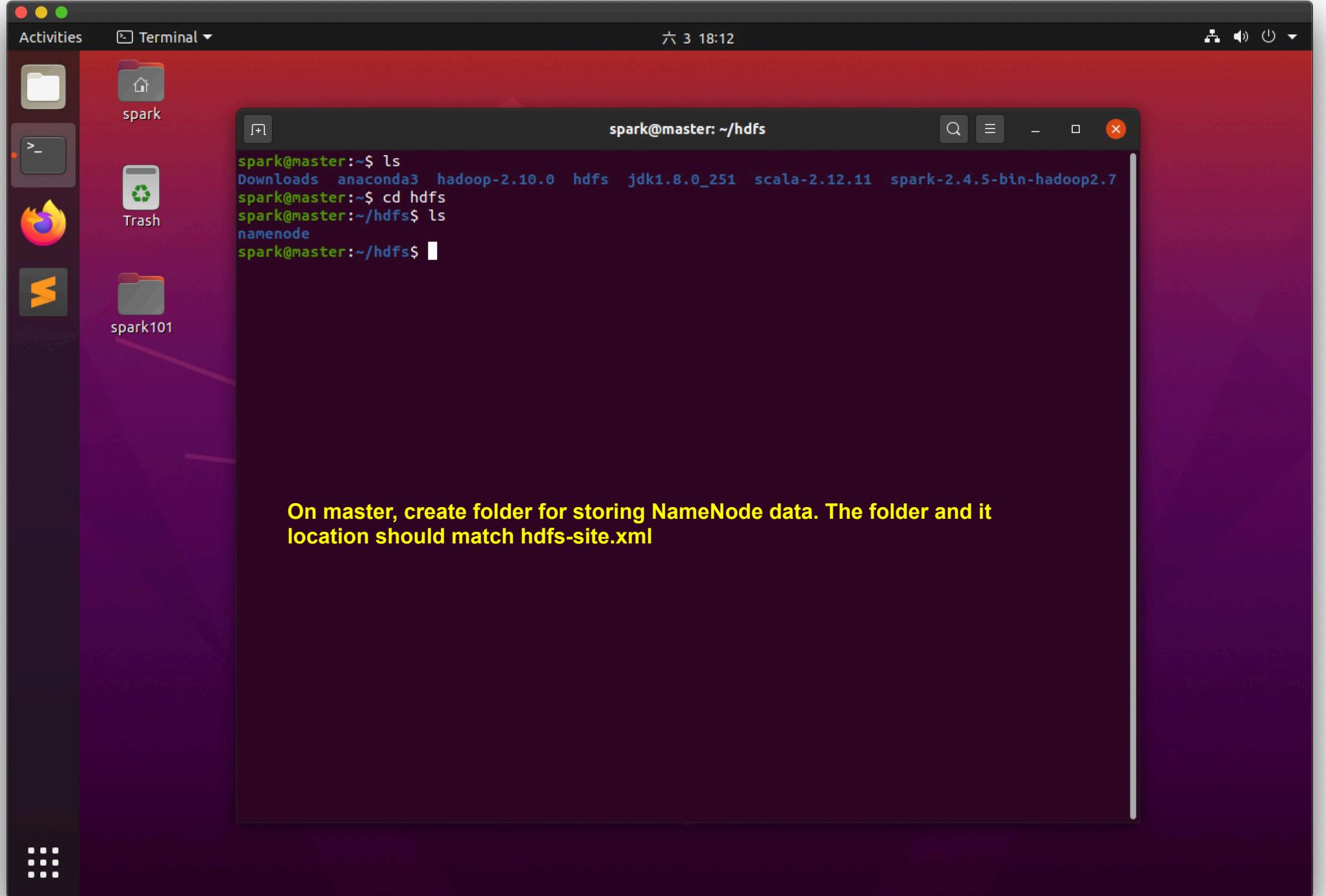


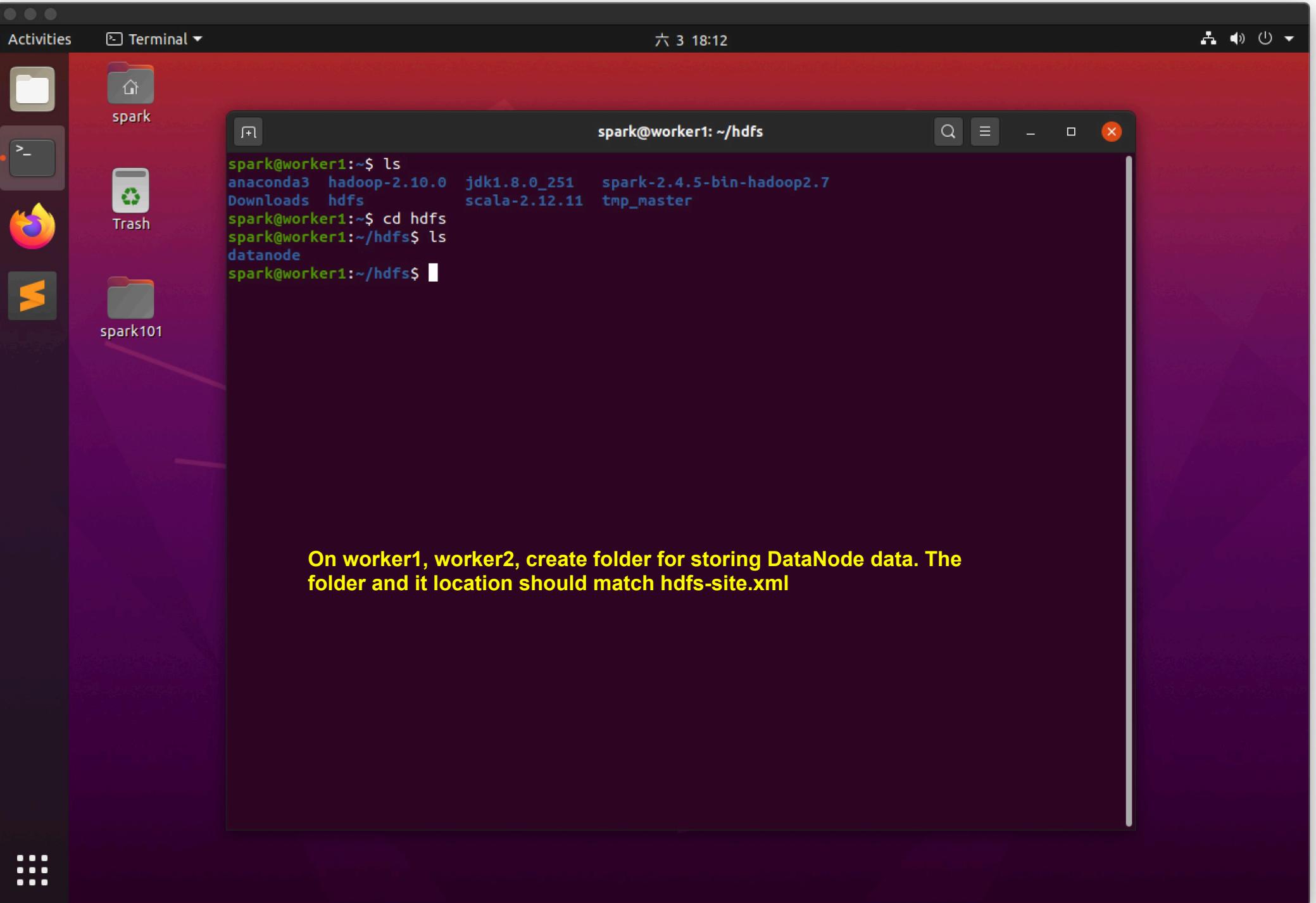


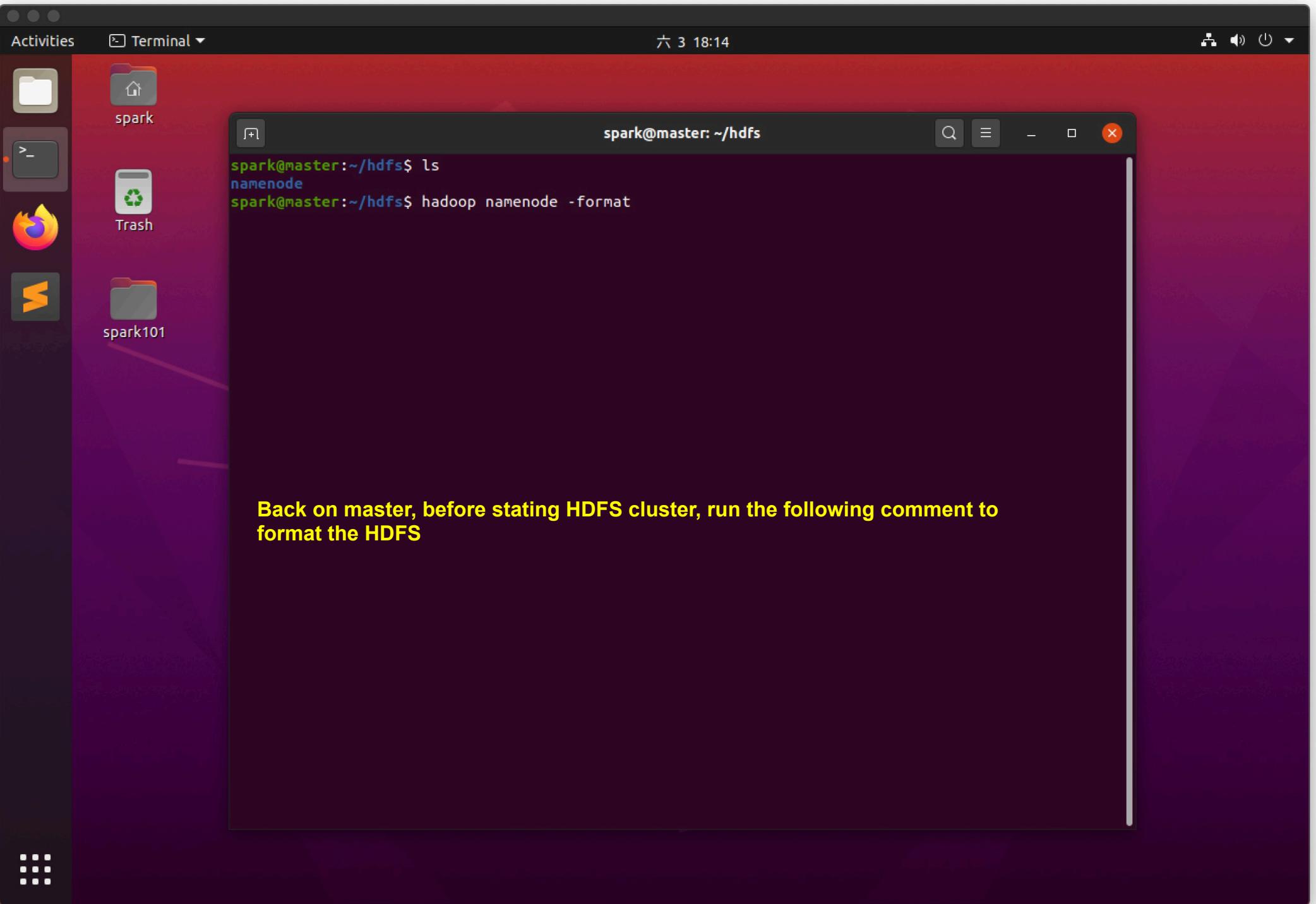


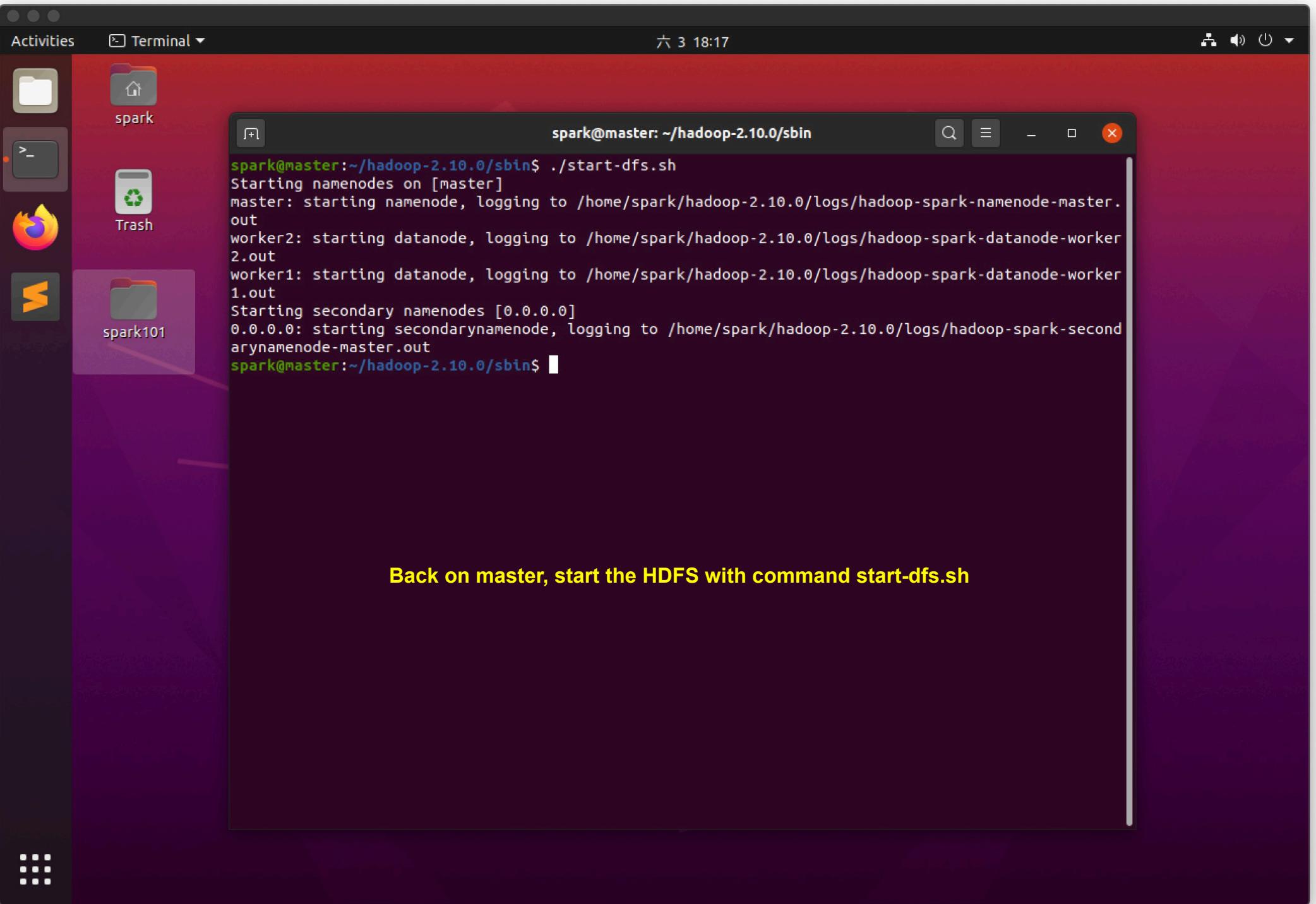


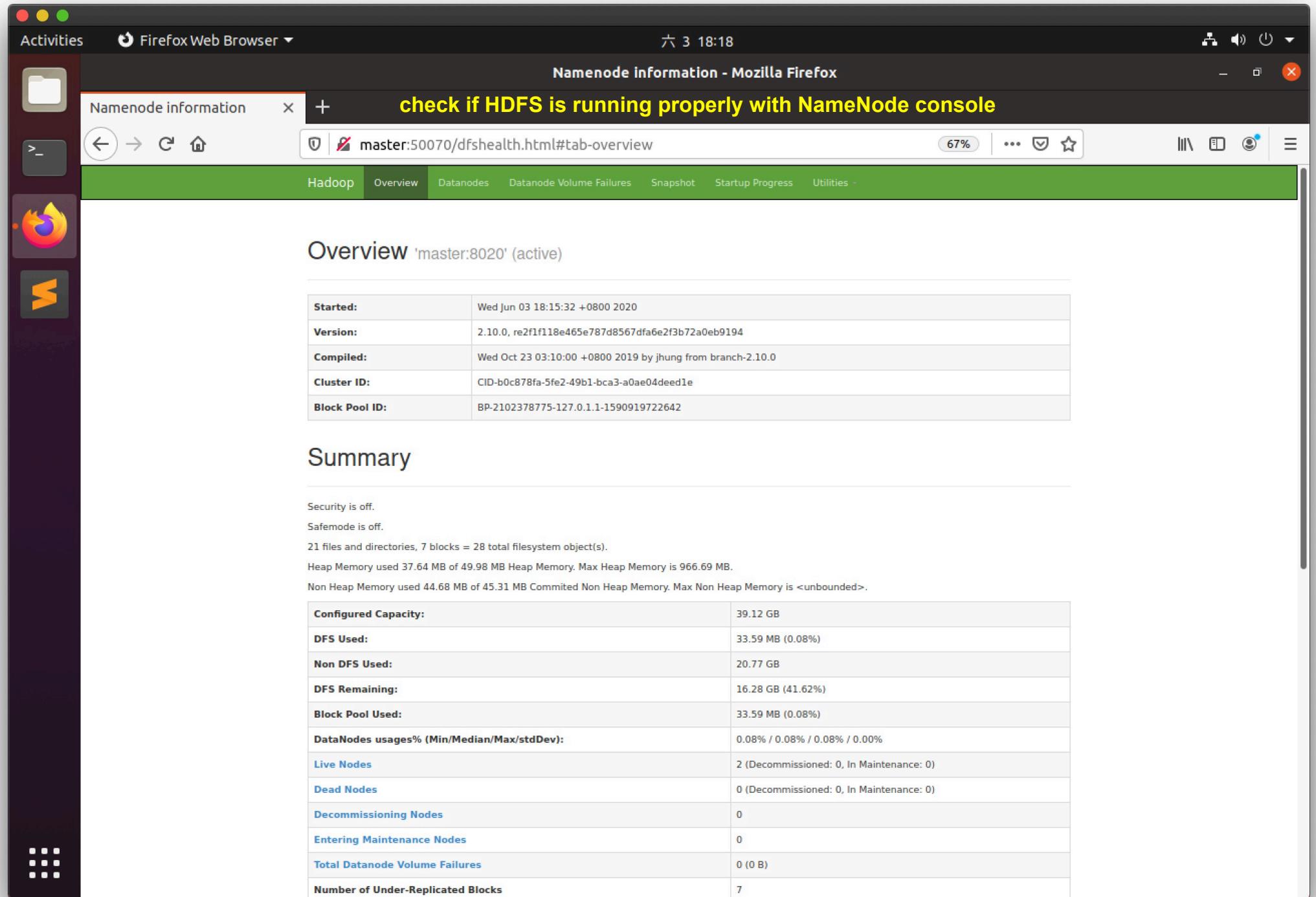


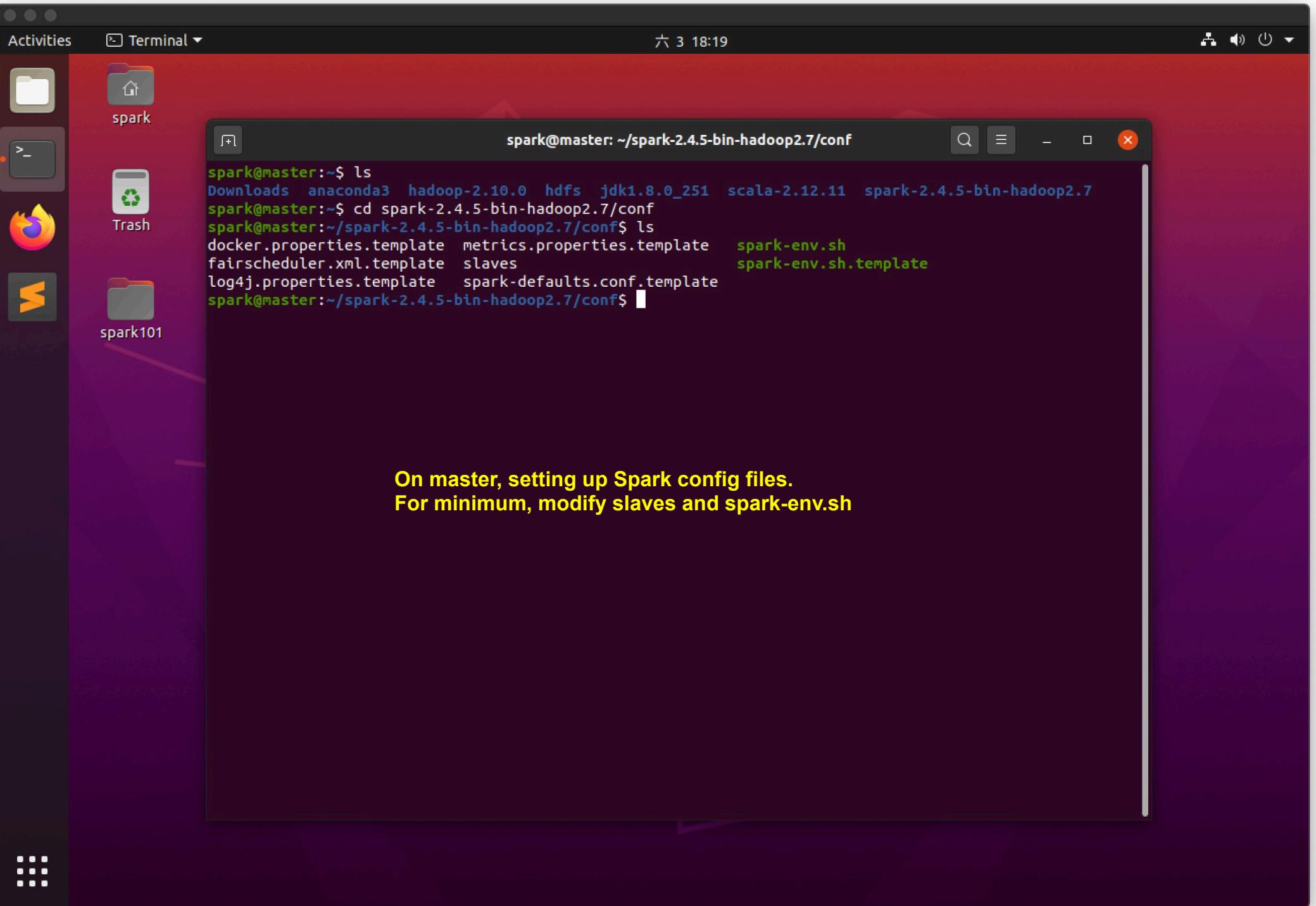












Activities

Terminal

六 3 18:19

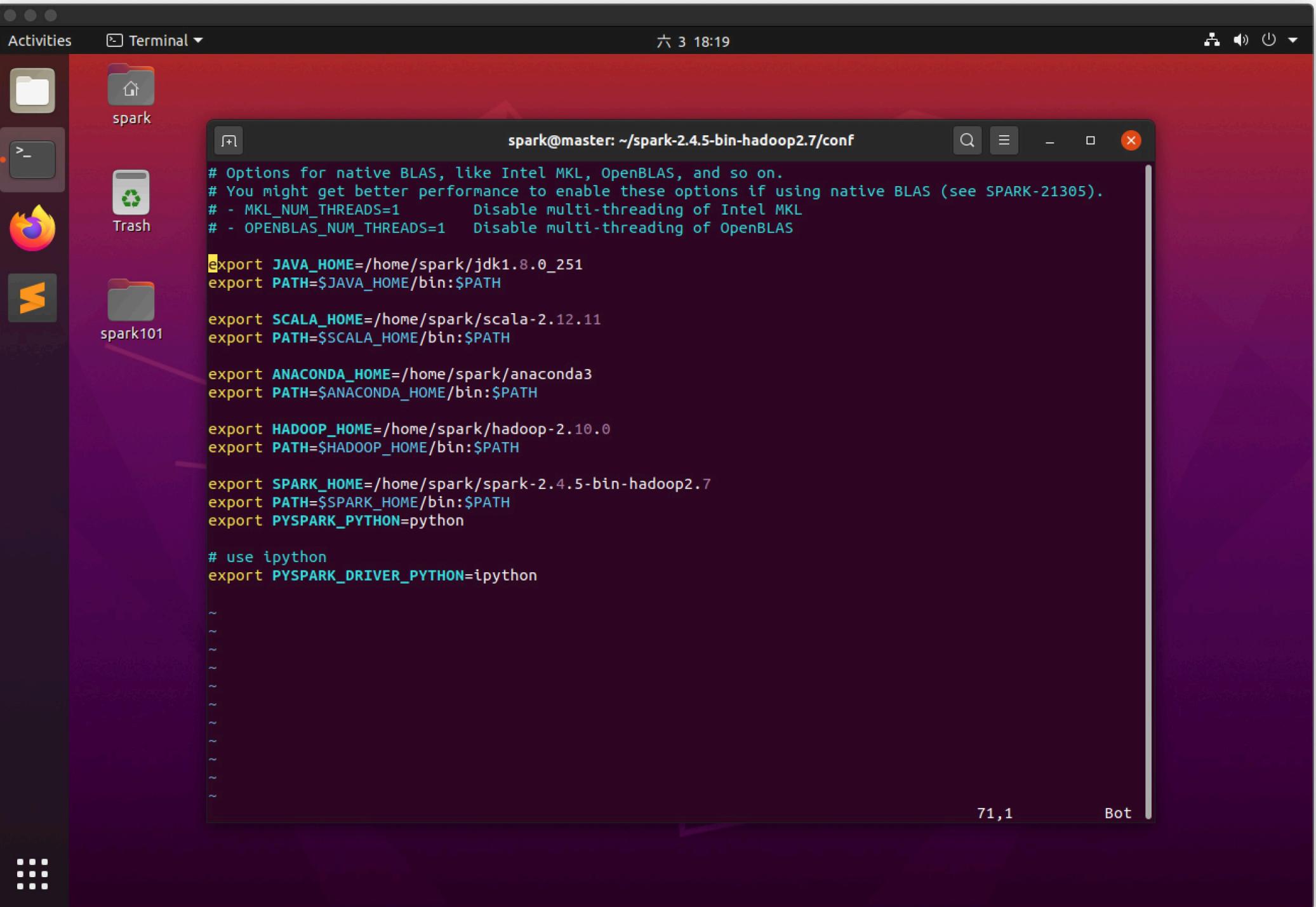


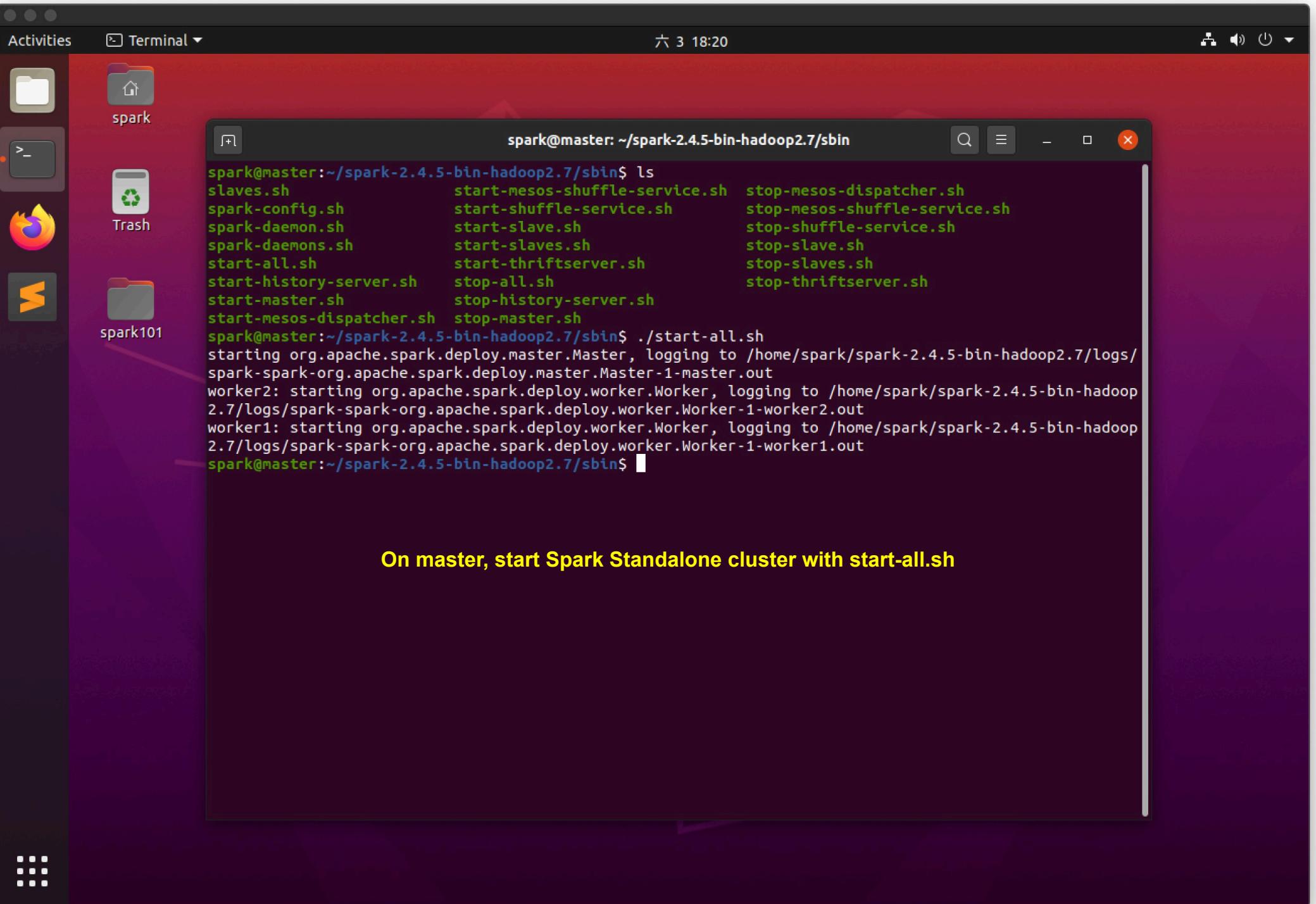
spark

Trash

spark101

spark@master: ~/spark-2.4.5-bin-hadoop2.7/conf





Check the Spark Standalone console to see if the cluster is running properly

六 3 18:21

Activities Firefox Web Browser

Spark Master at spark://master:7077 - Mozilla Firefox

Namenode information X Spark Master at spark:///r X +

master:8080

Apache Spark 2.4.5

Spark Master at spark://master:7077

URL: spark://master:7077

Alive Workers: 2

Cores in use: 2 Total, 0 Used

Memory in use: 2.0 GB Total, 0.0 B Used

Applications: 0 Running, 0 Completed

Drivers: 0 Running, 0 Completed

Status: ALIVE

▼ Workers (2)

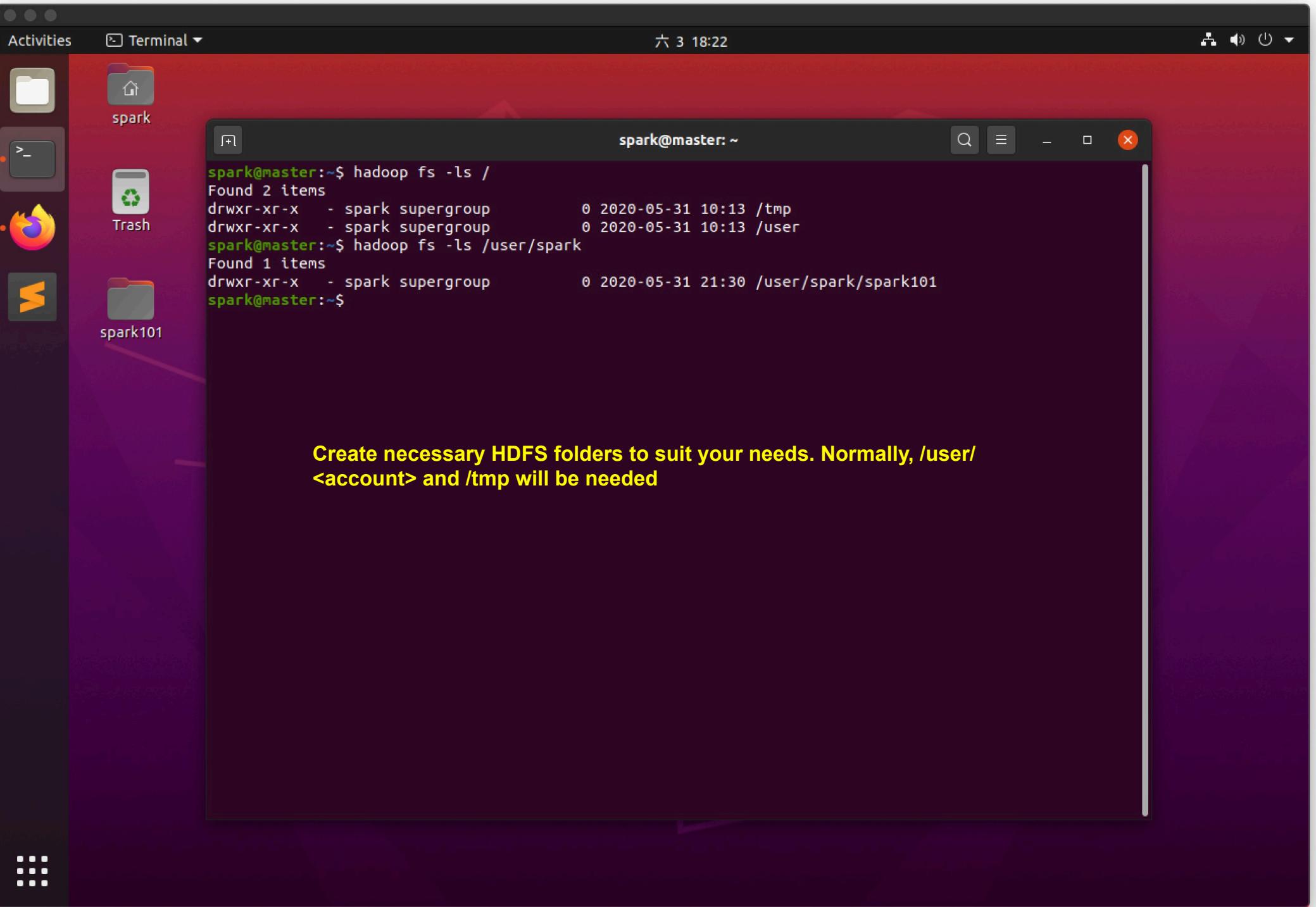
Worker Id	Address	State	Cores	Memory
worker-20200603102030-192.168.186.12-42649	192.168.186.12:42649	ALIVE	1 (0 Used)	1024.0 MB (0.0 B Used)
worker-20200603102030-192.168.186.13-34561	192.168.186.13:34561	ALIVE	1 (0 Used)	1024.0 MB (0.0 B Used)

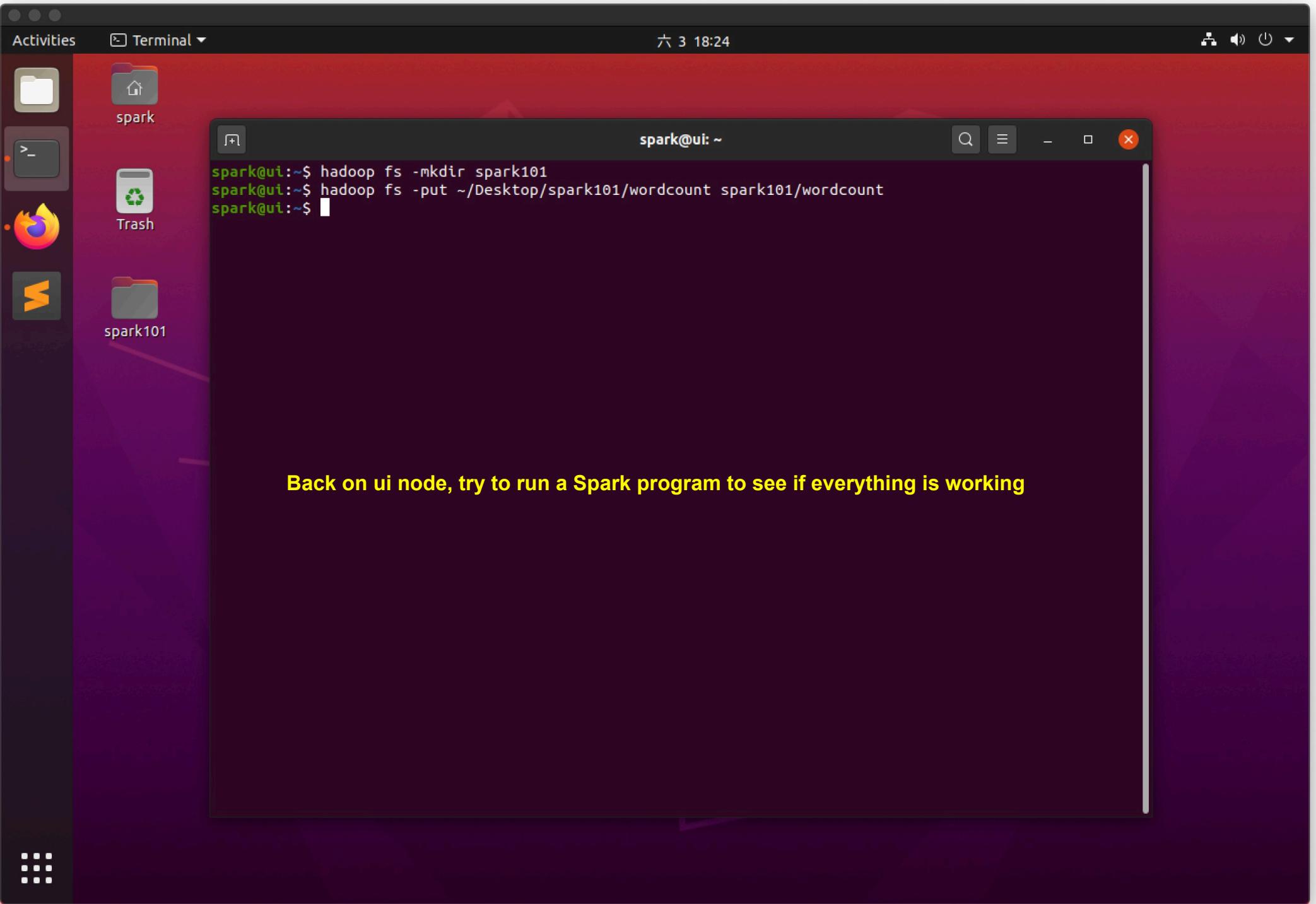
▼ Running Applications (0)

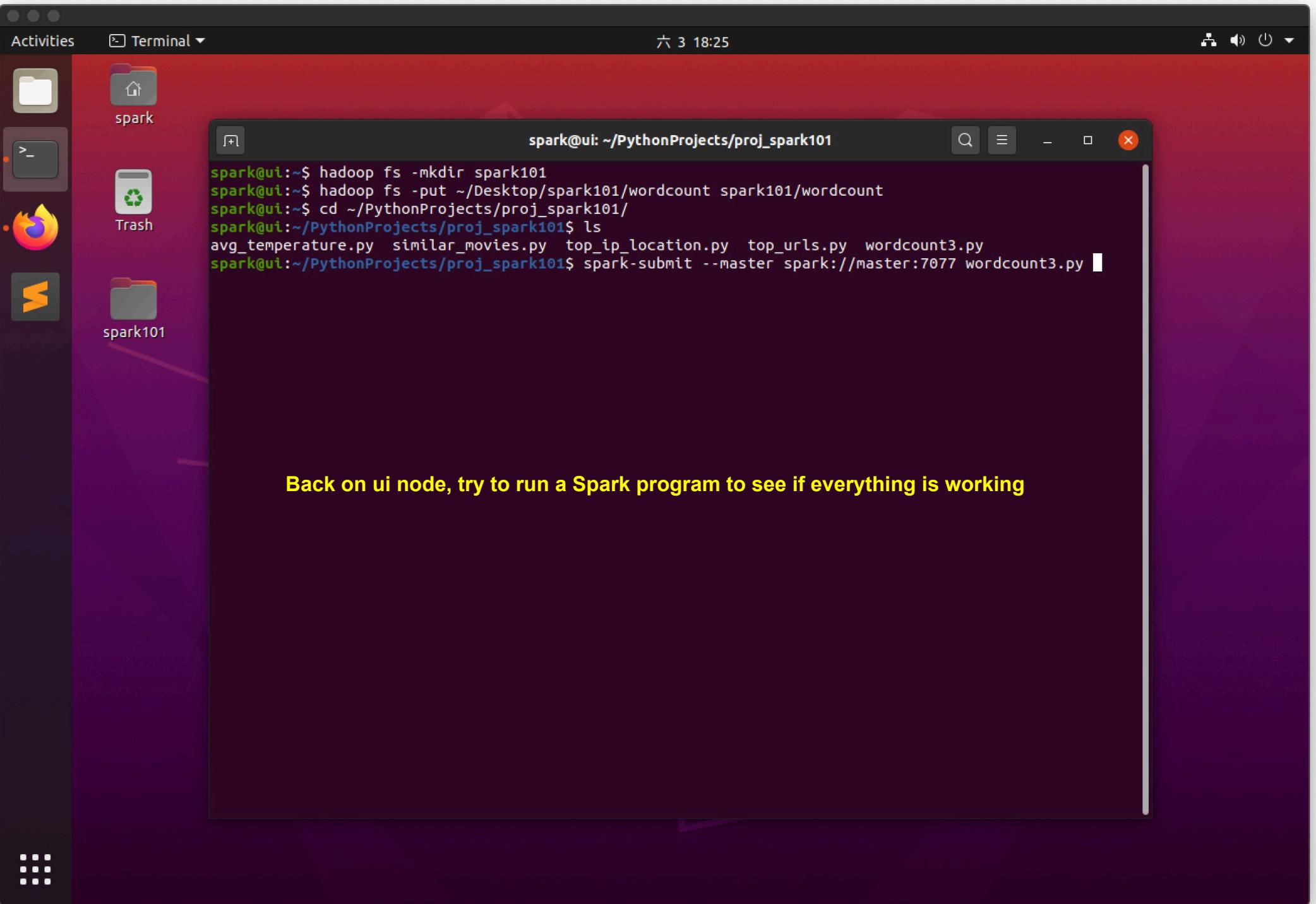
Application ID	Name	Cores	Memory per Executor	Submitted Time	User	State	Duration
----------------	------	-------	---------------------	----------------	------	-------	----------

▼ Completed Applications (0)

Application ID	Name	Cores	Memory per Executor	Submitted Time	User	State	Duration
----------------	------	-------	---------------------	----------------	------	-------	----------







Activities Firefox Web Browser ▾

六 3 18:26

Browsing HDFS - Mozilla Firefox

Browsing HDFS | Spark Master at spark:// | +

master:50070/explorer.html#/user/spark/spark101/wordcount/output2

Hadoop Overview Datanodes Datanode Volume Failures Snapshot Startup Progress Utilities

File information - part-00000

Download Head the file (first 32K) Tail the file (last 32K)

Block information -- Block 0

Block ID: 1073741836

Block Pool ID: BP-2102378775-127.0.1.1-1590919722642

Generation Stamp: 1012

Size: 139888

Availability:

- worker1
- worker2

File contents

```
('project', 100)
('gutenberg', 93)
('ebook', 11)
('of', 15008)
('peace', 115)
('leo', 4)
('tolstoy', 13)
('..', 35206)
```

Go! Search: Block Size Name

128 MB _SUCCESS

128 MB part-00000

128 MB part-00001

Previous 1 Next

Browse Directory

/user/spark/spark101/wordcount/output2

Show 25 entries

	Permission	Owner
<input type="checkbox"/>	-rw-r--r--	spark
<input type="checkbox"/>	-rw-r--r--	spark
<input type="checkbox"/>	-rw-r--r--	spark

Showing 1 to 3 of 3 entries

Hadoop, 2019.