# Course: Big Data
## *Lab 04*
# PySpark - RDD

## Question 1:

Based on the tutorial of PySpark, students install PySpark in Ubuntu.
- Define the environment variable: JAVA_HOME
- Define the environment variable: SPARK_HOME
- Start the pyspark-shell and write an instruction to print down the PySpark version
- Take the screenshot and insert it into the table below.



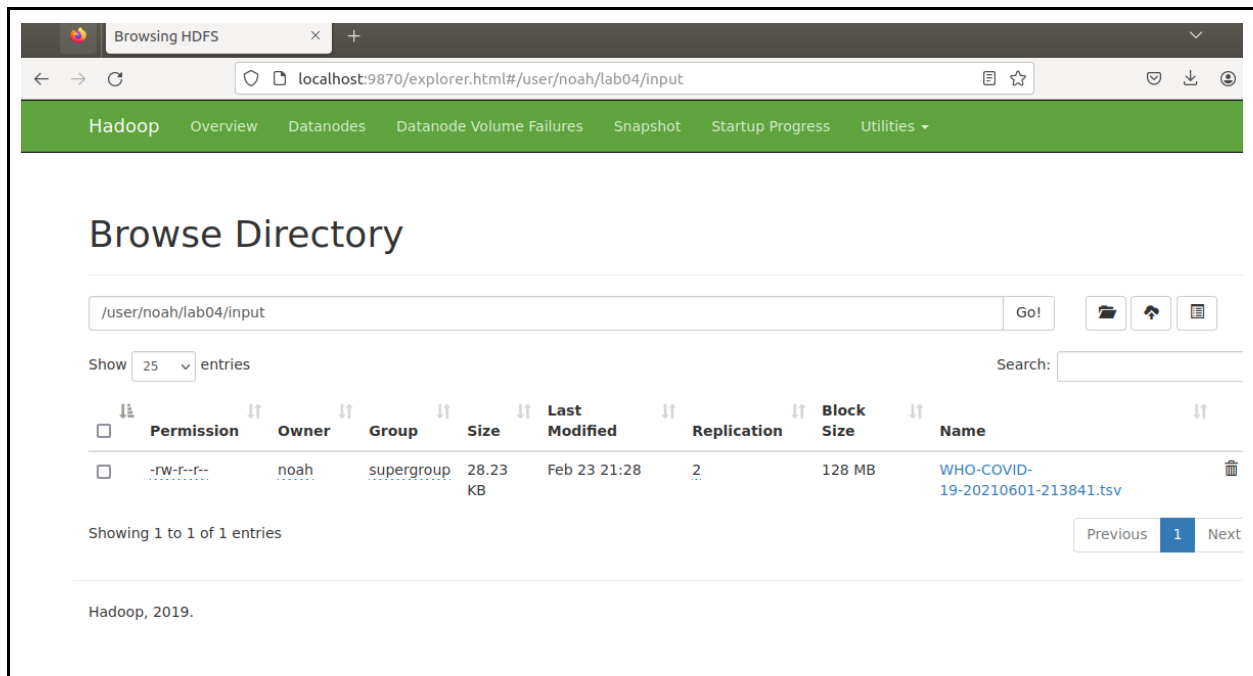## Question 2:

Given a tsv file WHO-COVID-19-20210601-213841.tsv which is corresponding to the WHO Coronavirus (COVID-19) Dashboard.

Students are required to create a folder, named **lab04**, in HDFS and then copy the tsv to **lab04/input/**

Take a screenshot to show the content of **lab04/input/** in HDFS

## Question 3:

Write a PySpark program, located in **ASEANCaseCount.py**, to count the number of cumulative total cases among ASEAN countries (*South-East Asia Region in the given data table*) using RDDs.

- Insert your source code into the table below.

```
from pyspark.sql import SparkSession

def filter_south_east_asia(line):
parts = line.split('\t')
region = parts[1].strip()
return region == "South-East Asia"

def parse_and_filter(line):
parts = line.split('\t')
country = parts[0]
cumulative_total_str = parts[2].replace(',','')
cumulative_total = float(cumulative_total_str)
return country, cumulative_total

if __name__ == "__main__":
# Create instance of SparkSession
spark = SparkSession.builder.appName("ASEANCaseCount").getOrCreate()

input_path = 'hdfs://localhost:9000/user/hoang/lab04/input/WHO-COVID-19-20210601-
```

```
213841.tsv'

input_rdd = spark.sparkContext.textFile(input_path)

filtered_rdd = input_rdd.filter(filter_south_east_asia)

asean_rdd = filtered_rdd.map(parse_and_filter)
print(asean_rdd.count())

cumulative_total = asean_rdd.map(lambda x: x[1]).reduce(lambda x, y: x + y)

print("Cumulative total cases among ASEAN countries in South-East Asia Region:",
cumulative_total)

spark.stop()
```

- Take a screenshot of the terminal to visualize the program result.
  <mark>Result: Cumulative total cases among ASEAN countries in South-East Asia Region: 31923614 cases.</mark>

```
t.py:33, took 0.198013 s
Cumulative total cases among ASEAN countries in South-East Asia Region: 31923614.0
24/02/27 21:22:21 INFO SparkUI: Stopped Spark web UI at http://192.168.146.121:4040
24/02/27 21:22:22 INFO MapOutputTrackerMasterEndpoint: MapOutputTrackerMasterEndpoint stopped!
24/02/27 21:22:22 INFO MemoryStore: MemoryStore cleared
24/02/27 21:22:22 INFO BlockManager: BlockManager stopped
24/02/27 21:22:22 INFO BlockManagerMaster: BlockManagerMaster stopped
24/02/27 21:22:22 INFO OutputCommitCoordinator$OutputCommitCoordinatorEndpoint: OutputCommitCoordinat
or stopped!
24/02/27 21:22:22 INFO SparkContext: Successfully stopped SparkContext
24/02/27 21:22:22 INFO ShutdownHookManager: Shutdown hook called
24/02/27 21:22:22 INFO ShutdownHookManager: Deleting directory /tmp/spark-558ffb4d-3258-4e7e-8c3e-5f5
c957ffa8b
24/02/27 21:22:22 INFO ShutdownHookManager: Deleting directory /tmp/spark-04787d68-bbc7-42eb-b8ad-870
593448da7/pyspark-0bc6e0e1-191d-47a7-9f55-198c6bd2926c
```