# Calculating Pi - PySpark implementation

Emily Weng

# Overview:

1. Use previous project:  Project: Creating MapReduce program to calculating Pi
2. Add the implementation of Pi Calculation using PySpark

# Follow the step for Pi on MapReduce

1. Make your instance and open SSH Browser

   a.
   ```
   Welcome to Ubuntu 20.04.6 LTS (GNU/Linux 5.15.0-1062-gcp x86_64)

    * Documentation:  https://help.ubuntu.com
    * Management:     https://landscape.canonical.com
    * Support:        https://ubuntu.com/pro

    System information as of Mon Jun 24 13:39:13 UTC 2024

     System load:  0.0               Processes:            104
     Usage of /:   19.4% of 9.51GB   Users logged in:      0
     Memory usage: 22%               IPv4 address for ens4: 10.140.0.11
     Swap usage:   0%

   Expanded Security Maintenance for Applications is not enabled.

   0 updates can be applied immediately.

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


   The list of available updates is more than a week old.
   To check for new updates run: sudo apt update


   The programs included with the Ubuntu system are free software;
   the exact distribution terms for each program are described in the
   individual files in /usr/share/doc/*/copyright.

   Ubuntu comes with ABSOLUTELY NO WARRANTY, to the extent permitted by
   applicable law.

   eweng909@instance-20240624-123745:~$
   ```

# Install Java JDK

1. sudo apt-get install openjdk-8-jdk

# Check version

1. $ java -version

```
eweng909@instance-20240624-123745:~$ java -version
openjdk version "1.8.0_412"
OpenJDK Runtime Environment (build 1.8.0_412-8u412-ga-1~20.04.1-b08)
OpenJDK 64-Bit Server VM (build 25.412-b08, mixed mode)
eweng909@instance-20240624-123745:~$
```

# Install ssh, sshd, phsd

1. Check if ssh/sshd/pdsh exists already, if not, install them
   a. which ssh
   b. which sshd
   c. which pshd

```
eweng909@instance-20240624-123745:~$ which ssh
/usr/bin/ssh
eweng909@instance-20240624-123745:~$ which sshd
/usr/sbin/sshd
eweng909@instance-20240624-123745:~$ which pshd
eweng909@instance-20240624-123745:~$ sudo apt-get install pdsh
Reading package lists... Done
Building dependency tree
Reading state information... Done
The following additional packages will be installed:
  genders libgenders0
Suggested packages:
  rdist
The following NEW packages will be installed:
  genders libgenders0 pdsh
0 upgraded, 3 newly installed, 0 to remove and 2 not upgraded.
Need to get 167 kB of archives.
After this operation, 519 kB of additional disk space will be used.
Do you want to continue? [Y/n] y
Get:1 http://asia-east1.gce.archive.ubuntu.com/ubuntu focal/universe a
 kB]
```

# Download Hadoop 3.3.5

1. wget https://dlcdn.apache.org/hadoop/common/hadoop-3.3.5/hadoop-3.3.5.tar.gz

# Unzip the tar file

1. Unzip the tar file
   a.    $  tar xzf hadoop-3.3.5.tar.gz

# Set up the rest of Hadoop Environment:

1. Modify bashrc file and set java and Hadoop environment
2. Configure HDFS

# Prepare input data

1.  $ mkdir PiCalculation
2.  $ cd PiCalculation
3.  $ vi GenerateRandomNumbers.java
4.  $ javac GenerateRandomNumbers.java
5.  $ java -cp . GenerateRandomNumbers

```
eweng909@instance-20240624-123745:~$ mkdir PiCalculation
eweng909@instance-20240624-123745:~$ cd PiCalculation
eweng909@instance-20240624-123745:~/PiCalculation$ vi GenerateRandomNumbers.java
eweng909@instance-20240624-123745:~/PiCalculation$ javac GenerateRandomNumbers.java
eweng909@instance-20240624-123745:~/PiCalculation$ java -cp . GenerateRandomNumbers
How many random numbers to generate:
10
What's the radius?
5
eweng909@instance-20240624-123745:~/PiCalculation$ ls
GenerateRandomNumbers.class   GenerateRandomNumbers.java   PiCalculationInput
```

# Set up paraphrase less SSH

1. ssh-keygen -t rsa -P '' -f ~/.ssh/id_rsa
2. cat ~/.ssh/id_rsa.pub >> ~/.ssh/authorized_keys
3. chmod 0600 ~/.ssh/authorized_keys
4. ssh localhost

```
eweng909@instance-20240624-123745:~/hadoop-3.3.5$ cat ~/.ssh/id_rsa.pub >> ~/.ssh/authorized_keys
eweng909@instance-20240624-123745:~/hadoop-3.3.5$ chmod 0600 ~/.ssh/authorized_keys
eweng909@instance-20240624-123745:~/hadoop-3.3.5$ ssh localhost
Welcome to Ubuntu 20.04.6 LTS (GNU/Linux 5.15.0-1062-gcp x86_64)

 * Documentation:  https://help.ubuntu.com
 * Management:     https://landscape.canonical.com
 * Support:        https://ubuntu.com/pro

 System information as of Mon Jun 24 15:03:04 UTC 2024

  System load:  0.0               Processes:             107
  Usage of /:   47.5% of 9.51GB   Users logged in:       1
  Memory usage: 30%               IPv4 address for ens4: 10.140.0.11
  Swap usage:   0%

 * Strictly confined Kubernetes makes edge and IoT secure. Learn how MicroK8s
   just raised the bar for easy, resilient and secure K8s cluster deployment.

   https://ubuntu.com/engage/secure-kubernetes-at-the-edge

Expanded Security Maintenance for Applications is not enabled.

2 updates can be applied immediately.
2 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

New release '22.04.3 LTS' available.
Run 'do-release-upgrade' to upgrade to it.


Last login: Mon Jun 24 14:53:43 2024 from 35.235.244.80
```

# Make the HDFS directories required to execute MapReduce jobs

1. cd hadoop-3.3.5
2. bin/hdfs namenode -format
3. sbin/start-dfs.sh

# Continue:

1. Update hdfs-site.xml and core-site.xml file
2. wget http://localhost:9870/

```
eweng909@instance-20240624-123745:~/hadoop-3.3.5$ wget http://localhost:9870/
--2024-06-24 15:19:38--  http://localhost:9870/
Resolving localhost (localhost)... 127.0.0.1
Connecting to localhost (localhost)|127.0.0.1|:9870... connected.
HTTP request sent, awaiting response... 302 Found
Location: http://localhost:9870/index.html [following]
--2024-06-24 15:19:38--  http://localhost:9870/index.html
Reusing existing connection to localhost:9870.
HTTP request sent, awaiting response... 200 OK
Length: 1079 (1.1K) [text/html]
Saving to: 'index.html'

index.html                100%[===============================================>]   1.05K  --.-KB/s    in 0s

2024-06-24 15:19:38 (101 MB/s) - 'index.html' saved [1079/1079]
```

# Continue:

1. bin/hdfs dfs -mkdir /user
2. bin/hdfs dfs -mkdir /user/eweng909
3. bin/hdfs dfs -mkdir /user/eweng909/picalculate
4. bin/hdfs dfs -mkdir /user/eweng909/picalculate/input
5. bin/hdfs dfs -put ../PiCalculation/PiCalculationInput /user/eweng909/picalculate/input

```
eweng909@instance-20240624-123745:~/hadoop-3.3.5$ bin/hdfs dfs -mkdir /user
eweng909@instance-20240624-123745:~/hadoop-3.3.5$ bin/hdfs dfs -mkdir /user/eweng909
eweng909@instance-20240624-123745:~/hadoop-3.3.5$ bin/hdfs dfs -mkdir /user/eweng909/picalculate
eweng909@instance-20240624-123745:~/hadoop-3.3.5$ bin/hdfs dfs -mkdir /user/eweng909/picalculate/input
eweng909@instance-20240624-123745:~/hadoop-3.3.5$ bin/hdfs dfs -put ../PiCalculation/PiCalculationInput /user/
weng909/picalculate/input
eweng909@instance-20240624-123745:~/hadoop-3.3.5$
```

# Build PiCalculation java file

1. cd hadoop-3.3.5
2. vi PiCalculation.java

```java
import java.io.*;
import java.util.*;
import java.lang.Object;
import java.net.URI;

import org.apache.hadoop.fs.Path;
import org.apache.hadoop.conf.*;
import org.apache.hadoop.io.*;
import org.apache.hadoop.mapreduce.*;
import org.apache.hadoop.mapreduce.Mapper.Context;
import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
import org.apache.hadoop.mapreduce.lib.input.TextInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
import org.apache.hadoop.mapreduce.lib.output.TextOutputFormat;
import org.apache.hadoop.fs.*;

public class PiCalculation {

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

        private final static IntWritable one = new IntWritable(1);
        private Text word = new Text();
        private int totalLines = 0;

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

            totalLines += 1;
            String line = value.toString();
            line = line.replace("(", "");
            line = line.replace(")", "");
            line = line.replace(",", " ");

            StringTokenizer itr = new StringTokenizer(line);
            int radius = 200;// Same as the one you give in PiDataGenerator stage
            while (itr.hasMoreTokens()) {
                String x, y;
                x = itr.nextToken();
```

# Compile PiCalculation.java and create a jar

1.  javac PiCalculation.java
2.  jar cf wc.jar PiCalculation*class

```
eweng909@instance-20240624-123745:~/hadoop-3.3.5$ javac PiCalculation.java
eweng909@instance-20240624-123745:~/hadoop-3.3.5$ jar cf pi.jar PiCalculation*.class
eweng909@instance-20240624-123745:~/hadoop-3.3.5$ ls
LICENSE-binary    'PiCalculation$IntSumReducer.class'     README.txt    index.html      logs
LICENSE.txt       'PiCalculation$TokenizerMapper.class'    bin          lib            pi.jar
NOTICE-binary      PiCalculation.class                     etc          libexec        sbin
NOTICE.txt         PiCalculation.java                      include      licenses-binary share
eweng909@instance-20240624-123745:~/hadoop-3.3.5$
```

# Results

hadoop jar pi.jar PiCalculation /user/eweng909/picalculate/input /user/eweng909/picalculate/output

# Output

1. bin/hdfs dfs -ls /user/eweng909/picalculate/new_output
2. bin/hdfs dfs -cat /user/eweng909/picalculate/new_output/part-r-00000

```
eweng909@instance-20240624-123745:~/hadoop-3.3.5$ bin/hdfs dfs -ls /user/eweng909/picalculate/new_output
Found 2 items
-rw-r--r--   1 eweng909 supergroup          0 2024-06-24 16:13 /user/eweng909/picalculate/new_output/_SUCCESS
-rw-r--r--   1 eweng909 supergroup         11 2024-06-24 16:13 /user/eweng909/picalculate/new_output/part-r-000
00
eweng909@instance-20240624-123745:~/hadoop-3.3.5$ bin/hdfs dfs -cat /user/eweng909/picalculate/new_output/part-
r-00000
outside 10
eweng909@instance-20240624-123745:~/hadoop-3.3.5$
```

With PySpark

# Using pyspark

1.  wget https://downloads.apache.org/spark/spark-3.5.1/spark-3.5.1-bin-hadoop3.tgz

# Unzip file

1. tar -xvf spark-3.5.1-bin-hadoop3.tgz

# Add directory path into your bash file

```
export SPARK_HOME=/home/eweng909/spark/spark-3.5.1-bin-hadoop3

export PATH=$PATH:$SPARK_HOME/bin
```

# Create python file

```python
from pyspark.sql import SparkSession
import random

def inside(p):
    x, y = random.random(), random.random()
    return x*x + y*y < 1

if __name__ == "__main__":
    spark = SparkSession.builder.appName("PiCalculation").getOrCreate()
    sc = spark.sparkContext

    num_samples = 1000000
    count = sc.parallelize(range(0, num_samples)).filter(inside).count()
    pi = 4 * count / num_samples
    print(f"Pi is roughly {pi}")

    spark.stop()
```

# Run the command for it to work

1. spark-submit picalculation.py
2. It should run the results.

```
24/06/25 13:28:16 INFO TaskSchedulerImpl: Killing all running tasks in stage 0: Stage finished
24/06/25 13:28:16 INFO DAGScheduler: Job 0 finished: count at /home/eweng909/pialculation.py:13, took 3.129067
s
Pi is roughly 3.141996
```