



Assignment NO - 11

| | |
|----------|--|
| Page No. | |
| Date | |



Title - Big data Analytics 2.

Design a distributed application using map-reduce which processes a long file of a system.

objective - Students should be able to processes a log file of a system using the Hadoop map-reduce framework.

~~Prerequisite - 1. Basic of java programming~~

~~Theory - steps to install Hadoop for distributed environment.~~

Java code for processes a log file a system

Steps to install hadoop for distributed environment

Initially create one folder store input file

(access - log - short - csv), sales country Reducer Java,
sales Country Driver .java files)

Step 1) Go to hadoop home directory and Format
the NameNode.

cd hadoop - 2.7.3.

bin/hadoop namenode - format.

Step 2) once the Namenode is the centerpiece of an
HDFS file system it keeps the directory tree
of all files stored in the HDFS and tracks
all the files stored across cluster.

l.hadoop - daemon .sh start namenode

a) start Data Node . Resource manager

Resource manager is the master that attributes all the available cluster resource and thus helps in managing the distributed apps running on the YARN system. Its work is to manage each Node Managers and the each applications running each Node Managers and the each application's AppMaster.

3) Start Data Node -

on startup a DataNode connects to the namenode and it responds to the requests from the Namenode for different operations.

4) Start Node Manager.

The Node manager in each machine framework is the agent which is responsible for managing container, monitoring the resource usage and reporting the same to the resource manager.

5) Start JobHistory Server -

Job History Server is responsible for servicing all job history related request from client / mr yarn - daemon.

Sh start nodemanager.

Step 3 - To check that all the Hadoop services are up and running run the below command.

jps

Step 4 - cd

Step 5 - sudo mkdir mapreduce - Vijay

- Step 6 - sudo chmod 777 - R mapreduce - vijay /
- Step 7 - sudo chown - R vijay mapreduce - vijay /
- Step 8 - sudo cp / home / vijay / Desktop / log files / * .w*
- mapreduce - vijay /
- Step 9 - cd mapreduce - vijay /
- Step 10 - ls
- Step 11 - sudo chmod +r *
- Step 12 - export .
- Step 13 - Javae - d - SalesMapper. jar a sales Country
Reducer. java Sales Country Driver. java
- Step 14 - ls
- Step 15 - cd sales Country /
- Step 16 - ls (check is class files are created)
- Step 17 - cd - -
- Step 18 - gedit Manifest .txt
- Step 19 - Jar - Cfm mapreduce - vijay .jar Manifest .txt
Sales Country / *. class
- Step 20 - ls
- Step 21 - cd
- Step 22 - cd mapreduce - vijay /
- Step 23 - sudo mkdir / input 200 .
- Step 24 - sudo cp access - log - short .csv / input 200
- Step 25 - SHADOOP - HOME / bin / hdfs dfs - put / input
200 /
- Step 26 - SHADOOP - HOME / bin / hadoop jar mapreduce -
vijay .jar / input 200 / output 200
- Step 27 - hadoop fs - ls / output 200 .

Step 28 - hadoop fs - cat /out /part 0000 .

Step 29 - Now open the Mozilla browser and go to localhost : 50070 /dfshealth.html to check the Namenode interface .

~~conclusion - Hence we study the design a distributed app in using Map-reduce which processes a log file of a system.~~

| TS | PR | VC | VA | RH | Total marks | Sign |
|----|----|----|----|----|-------------|------|
| 52 | 52 | 01 | 52 | 02 | 99 | L |

Title - A simple program in SCALA using Apache spark framework.

objective - write a simple program in SCALA using Apache spark framework.

Theory

- Steps to install Scala
- Apache Spark Framework Installation
- Source code

1. Install scala

Step 1) Java - Version

Step 2) Install Scala from the apt repository by running the following commands to search for Scala and install it.

sudo apt search scala \Rightarrow Search for the package
sudo apt install scala \Rightarrow install the package.

Step 3) - To verify the installation of Scala run the following command , scala -version.

2. Apache Spark Framework Installation.

Apache Spark is an open-source, distributed processing system used for big data workloads. It utilizes in memory caching, and optimized query execution for fast analytics queries.

against data of any size.

Step 1) Now go to the official Apache spark download page and grab the latest version (i.e. 3.2.1) at the time of writing this article. Alternatively, you can use the wget command to download the file directly in the terminal.

wget

~~https://apache-mirror.wuchna.com/spark/spark-3.2.1/spark-3.2.1-bin-hadoop2.7.tgz~~

Step 2) Extract the Apache Spark tar file. tar - xvzf spark-3.2.1-bin-hadoop2.7.tgz

Step 3) Move the extracted spark directory to /opt directory. sudo mv.

Spark-3.2.1-bin-hadoop2.7 /opt/spark

Configure Environmental variables for Spark

Step 4) Now you have to set a few environmental variables in .profile file before starting up the Spark.

echo "export SPARK_HOME = /opt/spark" >> ~/.profile

echo "export
PATH = \$PATH:/opt/spark/bin:/opt/spark/sbin"
in >> ~/.profile. echo "export"



against data of any size.

Step 1 - Now go to the official Apache spark download page and grab the latest version (i.e. 3.2.1) at the time of writing this article. Alternatively, you can use the wget command to download the file directly in the terminal.

wget

~~https://apache.mirror.wuchna.com/spark/spark-3.2.1/spark-3.2.1-bin-hadoop2.7.tgz~~

Step 2) Extract the Apache Spark tar file. tar -xvzf spark-3.2.1-bin-hadoop2.7.tgz

Step 3) Move the extracted spark directory to opt directory. sudo mv

spark-3.2.1-bin-hadoop2.7 /opt/spark

Configure Environmental variables for Spark

Step 4 - Now you have to set a few environmental variables in .profile file before starting up the Spark.

echo "export SPARK_HOME=/opt/spark" >> ~/.profile

echo "export
PATH=\$PATH:/opt/spark/bin:/opt/spark/sbin
in" >> ~/.profile echo "export

~ profile.

Step 5 :- To make sure that these new environment variables are reachable within the shell and available to apache spark , it is also mandatory to run the following command to take recent changes into effect.

Source ~ profile.

Step 6 - ls -l /opt /spark .

Start Apache browser spark in Ubuntu

Step 7) Run the following command to start the Spark master service and slave service

Start - master.sh.

Start - workers.sh spark :11 Localhost : 7077.

(if workers not starting then remove and install openssh:

Sudo apt - get remove openssh-client openssh-server

Sudo apt - get install openssh-client openssh-server.

Step 8) once the service is started go to the browser and type the following URL Access spark page from the page , you can see my master and slave service is started.

http :11 localhost : 8080

... run analytics queries



| | |
|----------|--|
| Page No. | |
| Date | |

step 9 - you can also check if spark shell works fine by launching the spark-shell command .

sudo apt install snapd

snap find "intellij"

sudo snap install intellij-idea-community
--classic

start intellij IDE Community Edition.

Conclusion - Hence we study have simple program using SCALA Apache spark framework.

| PR | PR | UC | VA | RN | Total | Sign |
|-----|-----|-----|-----|-----|-------|-------|
| TS | | | | | | marks |
| (2) | (2) | (2) | (2) | (2) | (10) | |
| 02 | 02 | 07 | 02 | 09 | 09 | ✓ |