# spark/hadoop單機版cluster架設

This note is written by Rudy Lee

## 環境

Vmware Workstaion - Centos7(VMware Workstation 15 Pro)

## 架設配置(以M、s1、s2代稱)

- tibame@master(Memory 16g、Processors 2、HDD 100g)
- tibame@slave1(Memory 5g、Processors 1、HDD 20g)
- tibame@slave2(Memory 5g、Processors 1、HDD 20g)

## 修改hostname的名稱(username@hostname)

```
[tibame@master bin]$ sudo vim /etc/hostname
[sudo] password for tibame:
[tibame@master bin]$ ▮
```
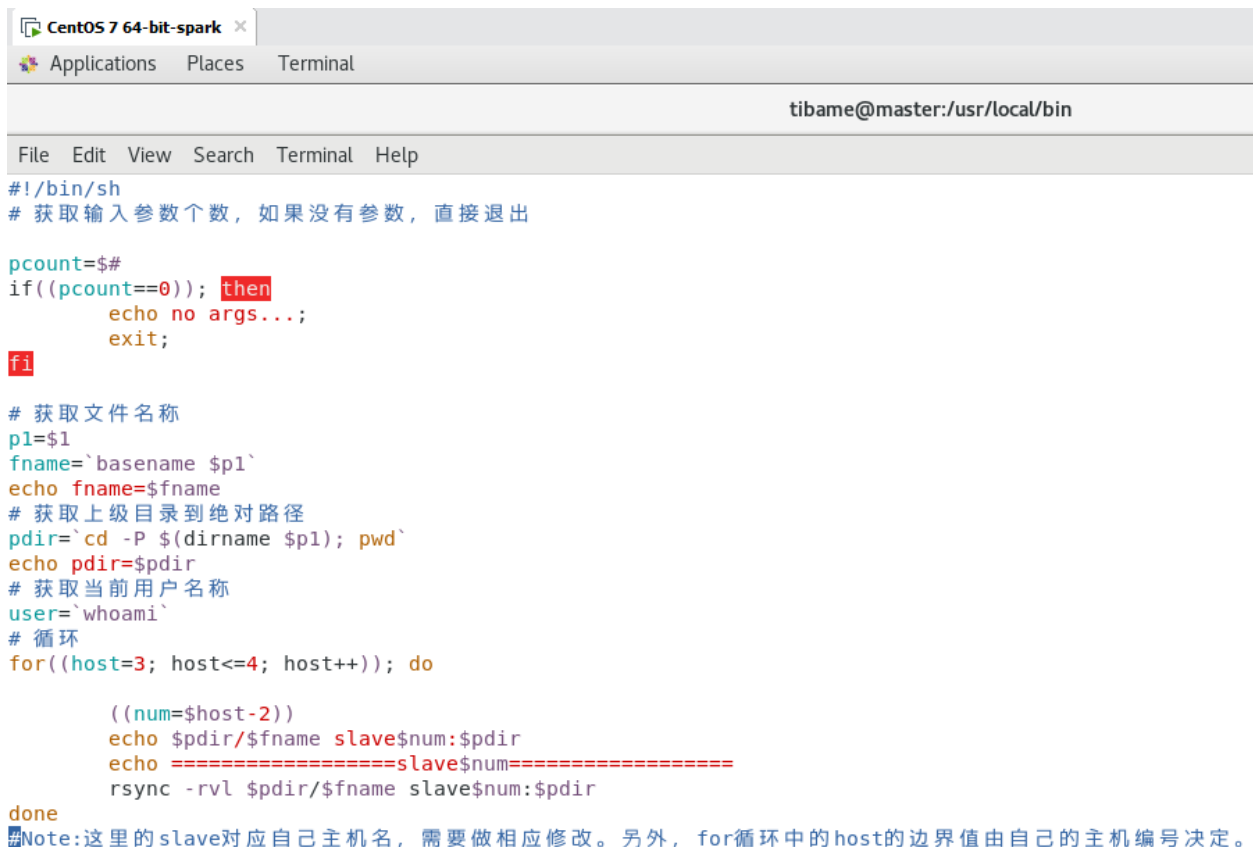
# 修改username的名稱

```
# usermod -l new-name old-name
```

# 執行前先準備xsync的function把主從資料夾同步

cd usr/local/bin

sudo vim xsync

```
#!/bin/sh
# 获取输入参数个数，如果没有参数，直接退出

pcount=$#
if((pcount==0)); then
        echo no args...;
        exit;
fi

# 获取文件名称
p1=$1
fname=`basename $p1`
echo fname=$fname
# 获取上级目录到绝对路径
pdir=`cd -P $(dirname $p1); pwd`
echo pdir=$pdir
# 获取当前用户名称
user=`whoami`
# 循环
for((host=3; host<=4; host++)); do

        ((num=$host-2))
        echo $pdir/$fname slave$num:$pdir
        echo ==================slave$num==================
        rsync -rvl $pdir/$fname slave$num:$pdir
done
#Note:这里的slave对应自己主机名，需要做相应修改。另外，for循环中的host的边界值由自己的主机编号决定。
```

# 建立keygen(方便登入slave,目標把keygen的鑰匙放到salve)

1. ssh-keygen

2. cat id_rsa.pub >> authorized_keys(複製key)

3. chmod 600 authorized_keys(打開權限)

4. ssh slave1

5. scp id_rsa.pub tibame@slave1:~/(把key送到salve1的桌面)

6. ssh-keygen(在slave1下建立.ssh目錄以放置key)

7. mv ~/authorized_keys /.ssh(把桌面的key放到slave1的的ssh資料夾裡面)

8. exit

9. ssh slave2(後續同上)

# 下載spark、hadoop、scala、jdk

💡 cd ~/Downloads

- spark-2.4.6-bin-hadoop2.7.tgz
- hadoop-2.10.0.tar.gz
- scala-2.11.12.tgz
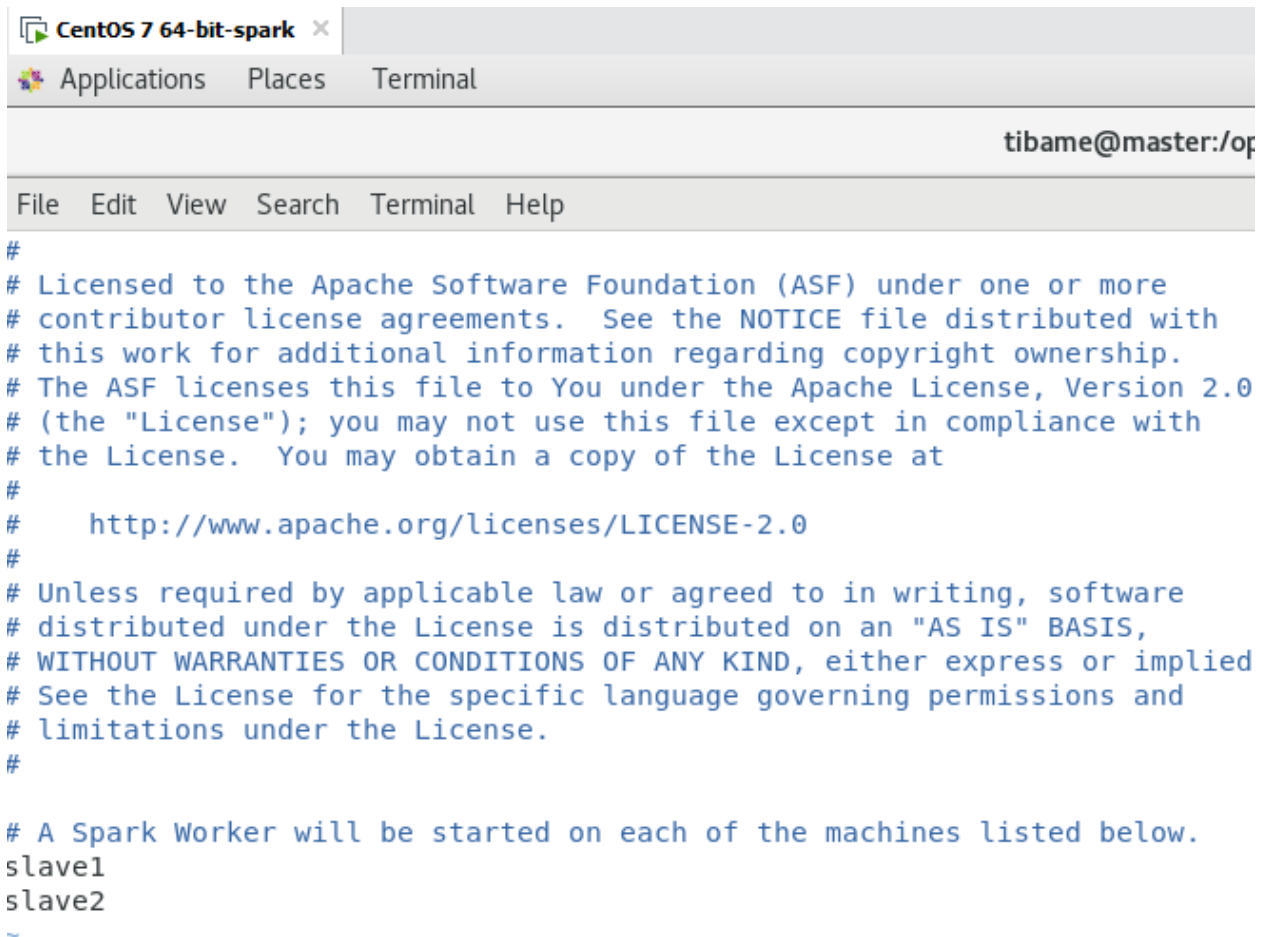- jdk-8u261-linux-x64.tar.gz

💡 tar -zxvf (壓縮檔), mv (解壓完目錄送到家目錄) ~/

# 設定spark環境

💡 cd spark/conf

1. mv slaves.template slaves

2. mv spark-env.sh.template spark-env.sh

3. vim slaves (加入slave1與slave2的worker)

```
CentOS 7 64-bit-spark  ×

Applications   Places   Terminal

                                                    tibame@master:/op

File  Edit  View  Search  Terminal  Help
#
# Licensed to the Apache Software Foundation (ASF) under one or more
# contributor license agreements.  See the NOTICE file distributed with
# this work for additional information regarding copyright ownership.
# The ASF licenses this file to You under the Apache License, Version 2.0
# (the "License"); you may not use this file except in compliance with
# the License.  You may obtain a copy of the License at
#
#    http://www.apache.org/licenses/LICENSE-2.0
#
# Unless required by applicable law or agreed to in writing, software
# distributed under the License is distributed on an "AS IS" BASIS,
# WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied
# See the License for the specific language governing permissions and
# limitations under the License.
#

# A Spark Worker will be started on each of the machines listed below.
slave1
slave2
~
```

4. vim spark-env.sh (在最下面加入)

SPARK_MASTER_WEBUI_PORT(UI介面)

SPARK_MASTER _PORT(執行port)

Applications  Places  Terminal

tibame@master:/opt/module/spark/conf

File  Edit  View  Search  Terminal  Help

```
# Options for the daemons used in the standalone deploy mode
# to bind the master to a different IP address or hostname
#/ SPARK_MASTER_WEBUI_PORT, to use non-default ports for the master
# - SPARK_MASTER_OPTS, to set config properties only for the master (e.g. "-Dx=y")
# - SPARK_WORKER_CORES, to set the number of cores to use on this machine
# - SPARK_WORKER_MEMORY, to set how much total memory workers have to give executors (e.g. 1000m, 2g)
# - SPARK_WORKER_PORT / SPARK_WORKER_WEBUI_PORT, to use non-default ports for the worker
# - SPARK_WORKER_DIR, to set the working directory of worker processes
# - SPARK_WORKER_OPTS, to set config properties only for the worker (e.g. "-Dx=y")
# - SPARK_DAEMON_MEMORY, to allocate to the master, worker and history server themselves (default: 1g).
# - SPARK_HISTORY_OPTS, to set config properties only for the history server (e.g. "-Dx=y")
# - SPARK_SHUFFLE_OPTS, to set config properties only for the external shuffle service (e.g. "-Dx=y")
# - SPARK_DAEMON_JAVA_OPTS, to set config properties for all daemons (e.g. "-Dx=y")
# - SPARK_DAEMON_CLASSPATH, to set the classpath for all daemons
# - SPARK_PUBLIC_DNS, to set the public dns name of the master or workers

# Generic options for the daemons used in the standalone deploy mode
# - SPARK_CONF_DIR      Alternate conf dir. (Default: ${SPARK_HOME}/conf)
# - SPARK_LOG_DIR       Where log files are stored.  (Default: ${SPARK_HOME}/logs)
# - SPARK_PID_DIR       Where the pid file is stored. (Default: /tmp)
# - SPARK_IDENT_STRING  A string representing this instance of spark. (Default: $USER)
# - SPARK_NICENESS      The scheduling priority for daemons. (Default: 0)
# - SPARK_NO_DAEMONIZE  Run the proposed command in the foreground. It will not output a PID file.
# Options for native BLAS, like Intel MKL, OpenBLAS, and so on.
# You might get better performance to enable these options if using native BLAS (see SPARK-21305).
# - MKL_NUM_THREADS=1        Disable multi-threading of Intel MKL
# - OPENBLAS_NUM_THREADS=1   Disable multi-threading of OpenBLAS


export SPARK_MASTER_IP="master"
export SPARK_MASTER_PORT="7077"
export SPARK_MASTER_WEBUI_PORT="8080"
"spark-env.sh" 74L, 4285C
```

5. xsync spark/(同步spark資料夾)

6. sbin/start-all.sh(啟動spark)

7. 查看spark安裝狀態(master:8080)

The screenshot shows the Spark Master web UI:

**Spark Master at spark://master:7077**

URL: spark://master:7077
Alive Workers: 2
Cores in use: 2 Total, 0 Used
Memory in use: 6.1 GB Total, 0.0 B Used
Applications: 0 Running, 9 Completed
Drivers: 0 Running, 0 Completed
Status: ALIVE

**Workers (2)**

| Worker Id | Address | State | Cores | Memory |
|---|---|---|---|---|
| worker-20200907054449-192.168.111.133-37620 | 192.168.111.133:37620 | ALIVE | 1 (0 Used) | 3.5 GB (0.0 B Used) |
| worker-20200907054449-192.168.111.134-45808 | 192.168.111.134:45808 | ALIVE | 1 (0 Used) | 2.6 GB (0.0 B Used) |

**Running Applications (0)**

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

**Completed Applications (9)**

| Application ID | Name | Cores | Memory per Executor | Submitted Time | User | State | Duration |
|---|---|---|---|---|---|---|---|
| app-20200909000324-0007 | Spark shell | 2 | 1024.0 MB | 2020/09/09 00:03:24 | tibame | KILLED | 25.9 h |
| app-20200909235039-0008 | Spark shell | 0 | 1024.0 MB | 2020/09/09 23:50:39 | tibame | KILLED | 2.1 h |
| app-20200907235119-0006 | Spark shell | 2 | 1024.0 MB | 2020/09/07 23:51:19 | tibame | KILLED | 23.5 h |
| app-20200907214031-0005 | Spark shell | 2 | 1024.0 MB | 2020/09/07 21:40:31 | tibame | KILLED | 1.6 h |
| app-20200907061938-0004 | Spark shell | 2 | 1024.0 MB | 2020/09/07 06:19:38 | tibame | KILLED | 15.3 h |
| app-20200907061454-0003 | Spark shell | 2 | 1024.0 MB | 2020/09/07 06:14:54 | tibame | FINISHED | 3.2 min |
| app-20200907060712-0002 | Spark shell | 2 | 1024.0 MB | 2020/09/07 06:07:12 | tibame | KILLED | 8.3 min |
| app-20200907060505-0001 | Spark shell | 2 | 1024.0 MB | 2020/09/07 06:05:05 | tibame | KILLED | 5.9 min |

## 8.啟動spark scala編寫環境

## cd bin

## ./spark-shell spark://master:7077

```
[tibame@master bin]$ ./spark-shell spark://master:7077
20/09/10 02:01:54 WARN NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
Using Spark's default log4j profile: org/apache/spark/log4j-defaults.properties
Setting default log level to "WARN".
To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use setLogLevel(newLevel).
20/09/10 02:02:12 WARN Utils: Service 'SparkUI' could not bind on port 4040. Attempting port 4041.
20/09/10 02:02:12 WARN Utils: Service 'SparkUI' could not bind on port 4041. Attempting port 4042.
Spark context Web UI available at http://master:4042
Spark context available as 'sc' (master = local[*], app id = local-1599728532374).
Spark session available as 'spark'.
Welcome to
      ____              __
     / __/__  ___ _____/ /__
    _\ \/ _ \/ _ `/ __/  '_/
   /___/ .__/\_,_/_/ /_/\_\   version 2.4.6
      /_/

Using Scala version 2.11.12 (Java HotSpot(TM) 64-Bit Server VM, Java 1.8.0_261)
Type in expressions to have them evaluated.
Type :help for more information.

scala> ▮
```

## 9. 確認scala可以執行(輸入sc)

```
beeline.cmd              pyspark2.cmd        sparkR            spark-sql2.cmd
docker-image-tool.sh     pyspark.cmd         sparkR2.cmd       spark-sql.cmd
find-spark-home          run-example         sparkR.cmd        spark-submit
find-spark-home.cmd      run-example.cmd     spark-shell       spark-submit2.cmd
load-spark-env.cmd       spark-class         spark-shell2.cmd  spark-submit.cmd
load-spark-env.sh        spark-class2.cmd    spark-shell.cmd
[tibame@master bin]$ ./spark-shell.sh spark://master:7077
bash: ./spark-shell.sh: No such file or directory
[tibame@master bin]$ ./spark-shell spark://master:7077
20/09/10 02:01:54 WARN NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
Using Spark's default log4j profile: org/apache/spark/log4j-defaults.properties
Setting default log level to "WARN".
To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use setLogLevel(newLevel).
20/09/10 02:02:12 WARN Utils: Service 'SparkUI' could not bind on port 4040. Attempting port 4041.
20/09/10 02:02:12 WARN Utils: Service 'SparkUI' could not bind on port 4041. Attempting port 4042.
Spark context Web UI available at http://master:4042
Spark context available as 'sc' (master = local[*], app id = local-1599728532374).
Spark session available as 'spark'.
Welcome to
      ____              __
     / __/__  ___ _____/ /__
    _\ \/ _ \/ _ `/ __/  '_/
   /___/ .__/\_,_/_/ /_/\_\   version 2.4.6
      /_/

Using Scala version 2.11.12 (Java HotSpot(TM) 64-Bit Server VM, Java 1.8.0_261)
Type in expressions to have them evaluated.
Type :help for more information.

scala> sc.textFile("hdfs://master/input").flatMap(_.split(" ")).map((_,1)).reduceByKey(_+_).collect
res0: Array[(String, Int)] = Array((taiwan,1), (hello,4), (meme,1), (spark,1), (handsome,1))
```

# 設定hadoop環境

1. cd hadoop-2.10.0/etc/hadoop/

2. vim core-site.xml

```xml
<?xml version="1.0" encoding="UTF-8"?>
<?xml-stylesheet type="text/xsl" href="configuration.xsl"?>
<!--
  Licensed under the Apache License, Version 2.0 (the "License");
  you may not use this file except in compliance with the License.
  You may obtain a copy of the License at

    http://www.apache.org/licenses/LICENSE-2.0

  Unless required by applicable law or agreed to in writing, software
  distributed under the License is distributed on an "AS IS" BASIS,
  WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied.
  See the License for the specific language governing permissions and
  limitations under the License. See accompanying LICENSE file.
-->

<!-- Put site-specific property overrides in this file. -->

<configuration>
 <property>
  <name>fs.defaultFS</name>
  <value>hdfs://master</value>
 </property>
</configuration>
```

3. vim hdfs-site.xml

File   Edit   View   Search   Terminal   Help

```xml
<?xml version="1.0" encoding="UTF-8"?>
<?xml-stylesheet type="text/xsl" href="configuration.xsl"?>
<!--
  Licensed under the Apache License, Version 2.0 (the "License");
  you may not use this file except in compliance with the License.
  You may obtain a copy of the License at

    http://www.apache.org/licenses/LICENSE-2.0

  Unless required by applicable law or agreed to in writing, software
  distributed under the License is distributed on an "AS IS" BASIS,
  WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied.
  See the License for the specific language governing permissions and
  limitations under the License. See accompanying LICENSE file.
-->

<!-- Put site-specific property overrides in this file. -->

<configuration>
  <property>
    <name>dfs.namenode.name.dir</name>
    <value>/home/tibame/hdfs/namenode/</value>
  </property>
  <property>
    <name>dfs.datanode.data.dir</name>
    <value>/home/tibame/hdfs/datanode/</value>
  </property>
</configuration>
```

4. vim hadoop-env.sh (同時加入scala與jdk環境變數)

```
# export HADOOP_DFSROUTER_OPTS=""
###

###
# Advanced Users Only!
###

# The directory where pid files are stored. /tmp by default.
# NOTE: this should be set to a directory that can only be written to by
#       the user that will run the hadoop daemons.  Otherwise there is the
#       potential for a symlink attack.
export HADOOP_PID_DIR=${HADOOP_PID_DIR}
export HADOOP_SECURE_DN_PID_DIR=${HADOOP_PID_DIR}

# A string representing this instance of hadoop. $USER by default.
export HADOOP_IDENT_STRING=$USER



export JAVA_HOME=/home/tibame/jdk1.8.0_261
export PATH=$JAVA_HOME/bin:$PATH

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

export SPARK_HOME=/opt/module/spark
export PATH=$SPARK_HOME/bin:$PATH

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

export HADOOP_CONF_DIR=$HADOOP_HOME/etc/hadoop
```

## 5.vim slaves

```
master
slave1
slave2
```

## 6.設定hadoop及PATH環境變數

cd ~

vim ~/.bashrc

```
# .bashrc

# Source global definitions
if [ -f /etc/bashrc ]; then
        . /etc/bashrc
fi

# Uncomment the following line if you don't like systemctl's auto-paging feature:
# export SYSTEMD_PAGER=

# User specific aliases and functions
export JAVA_HOME=/home/tibame/jdk1.8.0_261
export PATH=$JAVA_HOME/bin:$PATH

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

export SPARK_HOME=/opt/module/spark
export PATH=$SPARK_HOME/bin:$PATH

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

export HADOOP_CONF_DIR=$HADOOP_HOME/etc/hadoop
```

source ~/bashrc

> 💡 修改過./bashrc都必須要source才會啟動

7. xsync hadoop-2.10.0(同步hadoop資料夾)

8. 創建hdfs資料夾

cd ~

mkdir hdfs

cd hdfs

mkdir namenode

mkdir datanode

## 9. 首次創建需要初始化

hadoop namenode -format

## 10. 啟動HDFS叢集

cd ~
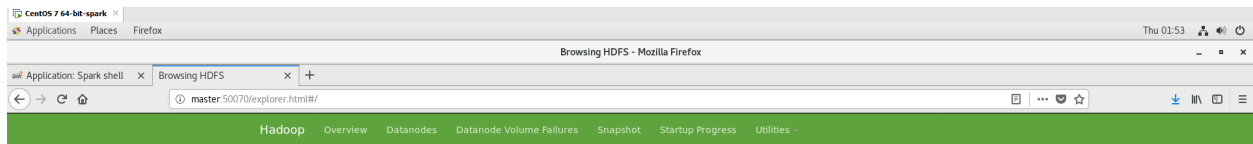cd hadoop-2.10.0/sbin

./start-dfs.sh

## 11. 檢查hdfs的狀態

http://master:50070



## 12. 創造hdfs的資料夾(透過utilities查看)

hadoop fs -mkdir /tmp

hadoop fs -mkdir -p /user/saprk

hadoop fs -ls -R /