

## Spark 安裝介紹

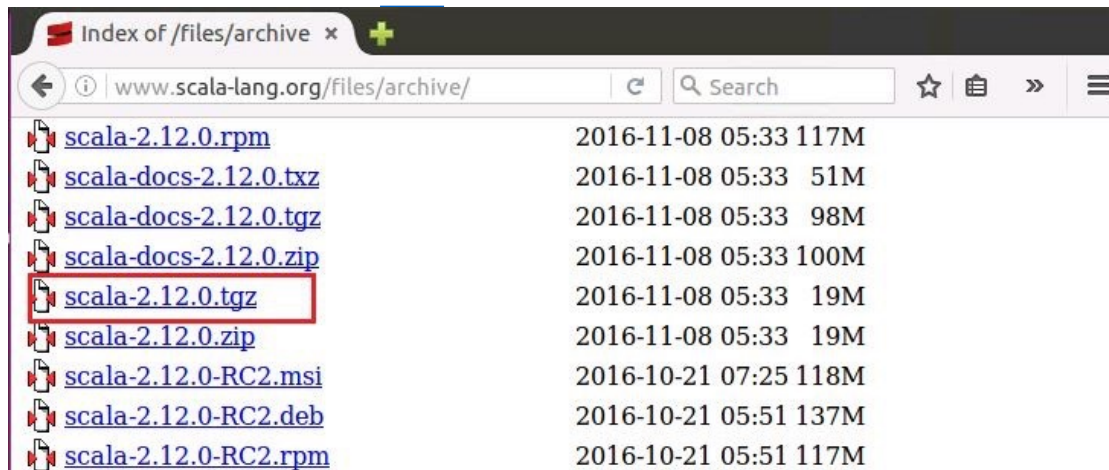
Spark Cluster Manager可以執行在下列模式：

1. 本機執行(Local Machine)：於本機執行，適合入門學習，測試用。
2. Spark Standalone cluster：由Spark提供的cluster管理模式，若沒有架設Hadoop Multi Node cluster，可以用本模式操作HDFS。
3. Hadoop YARN：於YARN上執行，由YARN進行多台機器的資源管理。
4. 雲端執行：針對更大型規模的計算工作，可以將Spark程式在雲端執行，如AWS的EC2平台。

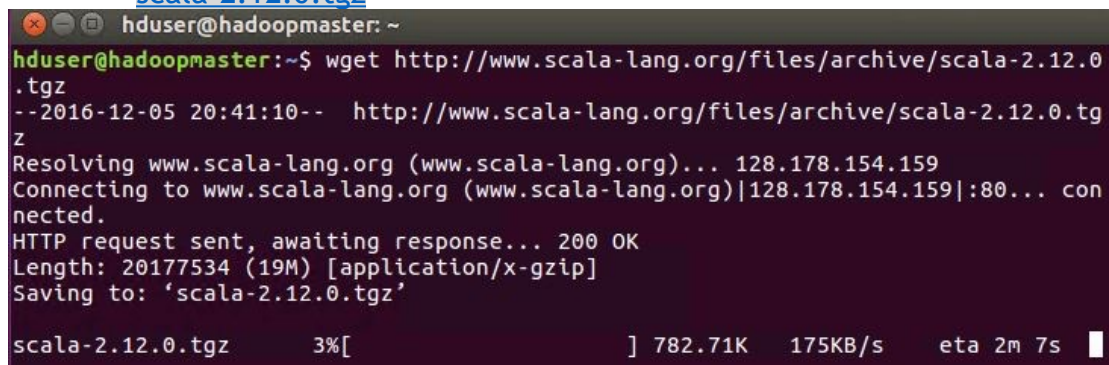
本文將教導本機執行和Spark Standalone cluster和YARN的安裝和執行方式。

### 1. Scala安裝

- Spark可以用python、Java等多種語言執行，本文選擇scala為主
- 下載 Scala，可於網址看到不同版本的Scala



- 執行 `wget http://www.scala-lang.org/files/archive/scala-2.12.0.tgz`



- 解壓縮Scala，輸入 `tar xvf scala-2.12.0.tgz`

```
hduser@hadoopmaster: ~
hduser@hadoopmaster:~$ tar xvf scala-2.12.0.tgz
scala-2.12.0/
scala-2.12.0/man/
scala-2.12.0/man/man1/
scala-2.12.0/man/man1/scala.1
scala-2.12.0/man/man1/scalap.1
scala-2.12.0/man/man1/fsc.1
scala-2.12.0/man/man1/scaladoc.1
scala-2.12.0/man/man1/scalac.1
scala-2.12.0/bin/
scala-2.12.0/bin/scalac
```

- 搬移至/usr/local下，輸入 `sudo mv scala-2.12.0 /usr/local/scala`

```
hduser@hadoopmaster: ~
hduser@hadoopmaster:~$ sudo mv scala-2.12.0 /usr/local/scala
[sudo] password for hduser:
hduser@hadoopmaster:~$
```

- 編輯 ~/.bashrc，輸入 `sudo gedit ~/.bashrc`

```
Export SCALA_HOME=/usr/local/scala
Export PATH=$PATH:$SCALA_HOME/bin
```

```
hduser@hadoopmaster:~$ sudo gedit ~/.bashrc

*.bashrc
~/
Save

export YARN_HOME=$HADOOP_HOME
export HADOOP_COMMON_LIB_NATIVE_DIR=$HADOOP_HOME/lib/native
export HADOOP_OPTS="-Djava.library.path=$HADOOP_HOME/lib"
export JAVA_LIBRARY_PATH=$HADOOP_HOME/lib/native:$JAVA_LIBRARY_PATH
#Hadoop Variables
#SCALA Variables
export SCALA_HOME=/usr/local/scala
export PATH=$PATH:$SCALA_HOME/bin
#SCALA Variables
```

- 讓 ~/.bashrc 生效，輸入 `source ~/.bashrc`

```
hduser@hadoopmaster: ~
hduser@hadoopmaster:~$ source ~/.bashrc
hduser@hadoopmaster:~$
```

- 到此即可執行Scala，輸入 `scala`，測試輸入程式執行

```
hduser@hadoopmaster: ~
hduser@hadoopmaster:~$ scala
Welcome to Scala 2.12.0 (OpenJDK 64-Bit Server VM, Java 1.8.0_111).
Type in expressions for evaluation. Or try :help.

scala> 1+1
res0: Int = 2

scala> :q
hduser@hadoopmaster:~$
```

## 2. 安裝Spark

- 到Spark網址下載Spark，注意需配合Hadoop版本來選擇Spark版本

Downloads | Apach... x +

spark.apache.org/downloads.html

## Download Apache Spark™

1. Choose a Spark release: **2.0.2 (Nov 14 2016) v**
2. Choose a package type: **Pre-built for Hadoop 2.7 and later v**
3. Choose a download type: **Select Apache Mirror v**
4. Download Spark: **spark-2.0.2-bin-hadoop2.7.tgz**
5. Verify this release using the [2.0.2 signatures and checksums](#) and [project release KEYS](#).

Apache Download ... x +

www.apache.org/dyn/closer.lua/spark/sj

The Apache Way

Contribute

ASF Sponsors

We suggest the following mirror site for your download:

**<http://apache.stu.edu.tw/spark/spark-2.0.2/spark-2.0.2-bin-hadoop2.7.tgz>**

Other mirror sites are suggested below. Please use the backup mirrors only to download PGP and MD5 signatures to [verify your downloads](#) or if no other mirrors are working.

- 下載Spark，輸入 `wget http://apache.stu.edu.tw/spark/spark-2.0.2/spark-2.0.2-bin-hadoop2.7.tgz`
- 解壓縮，輸入 `tar xzf spark-2.0.2-bin-hadoop2.7.tgz`
- 搬移至 `/usr/local/spark` 下，輸入 `sudo mv spark-2.0.2-bin-hadoop2.7 /usr/local/spark`

```
hduser@hadoopmaster: ~
hduser@hadoopmaster:~$ wget http://apache.stu.edu.tw/spark/spark-2.0.2/spark-2.0.2-bin-hadoop2.7.tgz
--2016-12-05 21:01:24-- http://apache.stu.edu.tw/spark/spark-2.0.2/spark-2.0.2-bin-hadoop2.7.tgz
Resolving apache.stu.edu.tw (apache.stu.edu.tw)... 120.119.118.1, 2001:e10:c41:e
eee::1
Connecting to apache.stu.edu.tw (apache.stu.edu.tw)|120.119.118.1|:80... connect
ed.
HTTP request sent, awaiting response... 200 OK
Length: 187426587 (179M) [application/x-gzip]
Saving to: 'spark-2.0.2-bin-hadoop2.7.tgz'

spark-2.0.2-bin-had 100%[=====] 178.74M 257KB/s in 10m 15s
2016-12-05 21:11:40 (297 KB/s) - 'spark-2.0.2-bin-hadoop2.7.tgz' saved [18742658
7/187426587]

hduser@hadoopmaster:~$ tar xzf spark-2.0.2-bin-hadoop2.7.tgz
hduser@hadoopmaster:~$ sudo mv spark-2.0.2-bin-hadoop2.7 /usr/local/spark
[sudo] password for hduser:
hduser@hadoopmaster:~$
```

- 編輯 `~/.bashrc`，輸入 `sudo gedit ~/.bashrc`

Export `SPARK_HOME=/usr/local/spark`  
Export `PATH=$PATH:$SPARK_HOME/bin`

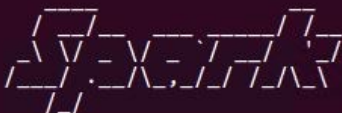


```
hduser@hadoopmaster: ~  
hduser@hadoopmaster:~$ sudo gedit ~/.bashrc  
#Hadoop Variables  
#SCALA Variables  
export SCALA_HOME=/usr/local/scala  
export PATH=$PATH:$SCALA_HOME/bin  
#SCALA Variables  
#SPARK Variables  
export SPARK_HOME=/usr/local/spark  
export PATH=$PATH:$SPARK_HOME/bin  
#SPARK Variables
```

- 讓設定生效，輸入 `source ~/.bashrc`

```
hduser@hadoopmaster: ~  
hduser@hadoopmaster:~$ source ~/.bashrc  
hduser@hadoopmaster:~$
```

- 啟動spark-shell，輸入 `spark-shell`

```
hduser@hadoopmaster: ~  
hduser@hadoopmaster:~$ spark-shell  
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).  
16/12/05 21:19:21 WARN NativeCodeLoader: Unable to load native-hadoop library fo  
r your platform... using builtin-java classes where applicable  
16/12/05 21:19:26 WARN SparkContext: Use an existing SparkContext, some configur  
ation may not take effect.  
Spark context Web UI available at http://192.168.59.137:4040  
Spark context available as 'sc' (master = local[*], app id = local-1480943965187  
).  
Spark session available as 'spark'.  
Welcome to  
 version 2.0.2  
Using Scala version 2.11.8 (OpenJDK 64-Bit Server VM, Java 1.8.0_111)  
type in expressions to have them evaluated.  
Type :help for more information.  
scala>
```

- 設定spark-shell互動介面的顯示訊息，因為預設會顯示過多訊息，影响閱讀。
  - i. `cd /usr/local/spark/conf`
  - ii. `cp log4j.properties.template log4j.properties`
  - iii. 編輯 `log4j.properties`，輸入 `sudo gedit log4j.properties`







```
scala> val textFile = sc.textFile("hdfs://hadoopmaster:9000/user/hduser/test/README.txt")
textFile: org.apache.spark.rdd.RDD[String] = hdfs://hadoopmaster:9000/user/hduser/test/README.txt MapPartitionsRDD[5] at textFile at <console>:24

scala> textFile.count()
res2: Long = 31

scala> █
```

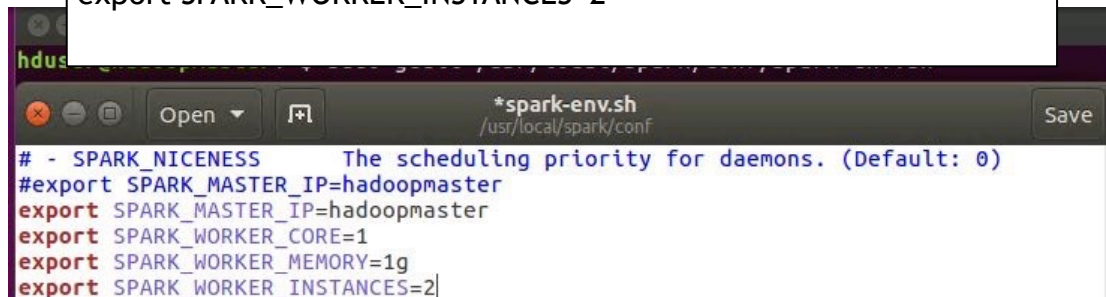
#### 4. 在Spark standalone cluster環境執行

- 自樣本建立spark-env.sh檔案，輸入`cp /usr/local/spark/conf/spark-env.sh.template /usr/local/spark/conf/spark-env.sh`

```
hduser@hadoopmaster:~$ cp /usr/local/spark/conf/spark-env.sh.template /usr/local/spark/conf/spark-env.sh
```

- 設定spark-env.sh，設定每個worker的資源分配，輸入`sudo gedit /usr/local/spark/conf/spark-env.sh`(注意：**每個worker的記憶體不得低於1G**，否則無法運行)

```
export SPARK_MASTER_IP=hadoopmaster
export SPARK_WORKER_CORE=1
export SPARK_WORKER_MEMORY=1g
export SPARK_WORKER_INSTANCES=2
```



```
*spark-env.sh
/usr/local/spark/conf

# - SPARK_NICENESS      The scheduling priority for daemons. (Default: 0)
#export SPARK_MASTER_IP=hadoopmaster
export SPARK_MASTER_IP=hadoopmaster
export SPARK_WORKER_CORE=1
export SPARK_WORKER_MEMORY=1g
export SPARK_WORKER_INSTANCES=2
```

- 將hadoopmaster的Spark程式複製到HadoopSlave1，輸入以下指令
  - ssh hadoopslave1
  - sudo mkdir /usr/local/spark
  - sudo chown hduser:hduser /usr/local/spark
  - exit

```

hduser@hadoopmaster:~$ ssh hadoopslave1
Welcome to Ubuntu 16.04.1 LTS (GNU/Linux 4.4.0-31-generic x86_64)

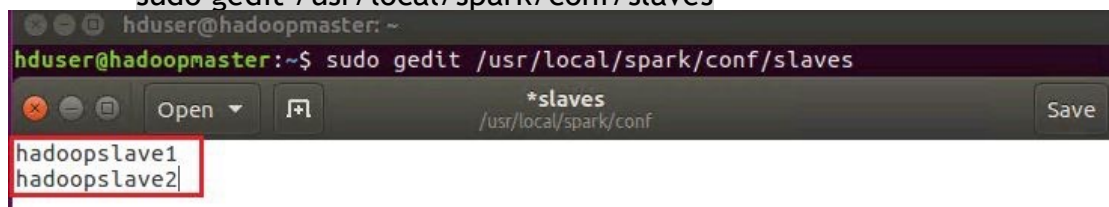
 * Documentation:  https://help.ubuntu.com
 * Management:    https://landscape.canonical.com
 * Support:       https://ubuntu.com/advantage

196 packages can be updated.
4 updates are security updates.

*** System restart required ***
Last login: Thu Dec  8 20:24:01 2016 from 192.168.59.137
hduser@hadoopslave1:~$ sudo mkdir /usr/local/spark
[sudo] password for hduser:
hduser@hadoopslave1:~$ sudo chown hduser:hduser /usr/local/spark
hduser@hadoopslave1:~$ exit
logout
Connection to hadoopslave1 closed.
hduser@hadoopmaster:~$ sudo scp -r /usr/local/spark hduser@hadoopslave1:/usr/local
The authenticity of host 'hadoopslave1 (192.168.59.134)' can't be established.
ECDSA key fingerprint is SHA256:l5HfVz2GKon2xpmavQSLRqfvPdxuogiqaf/Xjx5XV3E.
Are you sure you want to continue connecting (yes/no)? yes
Warning: Permanently added 'hadoopslave1,192.168.59.134' (ECDSA) to the list of
known hosts.
hduser@hadoopslave1's password:
NOTICE                               100% 24KB 24.2KB/s 00:00
hello.txt                           100% 13   0.0KB/s 00:00
userlib-0.1.zip                      100% 668  0.7KB/s 00:00

```

- 仿上面步驟將Spark複製到HadoopSlave2
- 編輯slaves檔案，設定Spark Standalone cluster有那些伺服器，輸入  
sudo gedit /usr/local/spark/conf/slaves



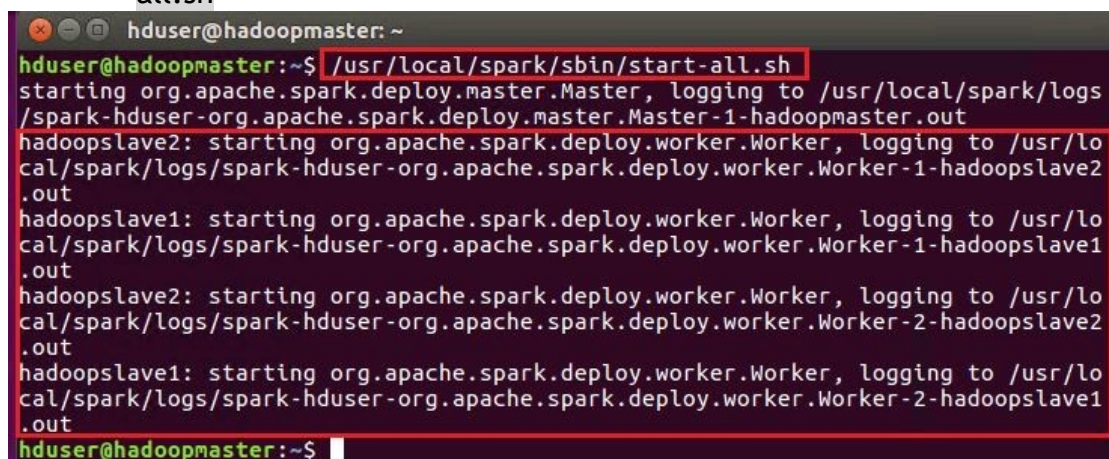
```

hduser@hadoopmaster:~$ sudo gedit /usr/local/spark/conf/slaves

```

The screenshot shows a terminal window where the command `sudo gedit /usr/local/spark/conf/slaves` is executed. Below it, a window titled `*slaves` displays the content of the `slaves` file, which lists `hadoopslave1` and `hadoopslave2` as cluster members.

- 啟動Spark Standalone cluster，輸入 /usr/local/spark/sbin/start-all.sh



```

hduser@hadoopmaster:~$ /usr/local/spark/sbin/start-all.sh
starting org.apache.spark.deploy.master.Master, logging to /usr/local/spark/logs
/spark-hduser-org.apache.spark.deploy.master.Master-1-hadoopmaster.out
hadoopslave2: starting org.apache.spark.deploy.worker.Worker, logging to /usr/local/spark/logs/spark-hduser-org.apache.spark.deploy.worker.Worker-1-hadoopslave2.out
hadoopslave1: starting org.apache.spark.deploy.worker.Worker, logging to /usr/local/spark/logs/spark-hduser-org.apache.spark.deploy.worker.Worker-1-hadoopslave1.out
hadoopslave2: starting org.apache.spark.deploy.worker.Worker, logging to /usr/local/spark/logs/spark-hduser-org.apache.spark.deploy.worker.Worker-2-hadoopslave2.out
hadoopslave1: starting org.apache.spark.deploy.worker.Worker, logging to /usr/local/spark/logs/spark-hduser-org.apache.spark.deploy.worker.Worker-2-hadoopslave1.out
hduser@hadoopmaster:~$

```

The screenshot shows a terminal window where the command `/usr/local/spark/sbin/start-all.sh` is executed. The output shows the Spark Master and Worker nodes starting up, with log files being created for each node.

- 可由上圖得知，共啟動了4個worker
- 執行Spark-shell，輸入 spark-shell --master spark://hadoopmaster:7077



```
hduser@hadoopmaster: ~
hduser@hadoopmaster:~$ spark-shell --master spark://hadoopmaster:7077
16/12/09 20:13:14 WARN NativeCodeLoader: Unable to load native-hadoop library fo
r your platform... using builtin-java classes where applicable
16/12/09 20:13:16 WARN SparkConf:
SPARK_WORKER_INSTANCES was detected (set to '2').
This is deprecated in Spark 1.0+.

Please instead use:
- ./spark-submit with --num-executors to specify the number of executors
- Or set SPARK_EXECUTOR_INSTANCES
- spark.executor.instances to configure the number of instances in the spark co
nfig.

16/12/09 20:13:34 WARN SparkContext: Use an existing SparkContext, some configur
ation may not take effect.
Spark context Web UI available at http://192.168.59.137:4040
Spark context available as 'sc' (master = spark://hadoopmaster:7077, app id = ap
p-20161209201330-0001).
Spark session available as 'spark'.
Welcome to

  ____  __
 / ___/ /_
/ /   / __/
/ /___/ __/
 \___/_/

 version 2.0.2

Using Scala version 2.11.8 (OpenJDK 64-Bit Server VM, Java 1.8.0_111)
Type in expressions to have them evaluated.
Type :help for more information.

scala>
```

- 查看Spark Standalone WebUI，於瀏覽器輸入 <http://hadoopmaster:8080>

Spark Master at spark:/... x +

hadoopmaster:8080

Applications: 1 Running, 1 Completed  
Drivers: 0 Running, 0 Completed  
Status: ALIVE

Workers

Worker Id	Address	State	Cores	Memory
worker-20161215193743-192.168.59.135-34289	192.168.59.135:34289	ALIVE	4 (4 Used)	1024.0 MB (1024.0 MB Used)
worker-20161215193743-192.168.59.135-36856	192.168.59.135:36856	ALIVE	4 (4 Used)	1024.0 MB (1024.0 MB Used)
worker-20161215193745-192.168.59.134-36926	192.168.59.134:36926	ALIVE	4 (4 Used)	1024.0 MB (1024.0 MB Used)
worker-20161215193745-192.168.59.134-40080	192.168.59.134:40080	ALIVE	4 (4 Used)	1024.0 MB (1024.0 MB Used)

Running Applications

Application ID	Name	Cores	Memory per Node	Submitted Time	User	State	Duration
app-20161215194328-0001	(kill) Spark shell	16	1024.0 MB	2016/12/15 19:43:28	hduser	RUNNING	1.3 min

- 讀取本地檔案
  - v. `val textFile=sc.textFile("file:/usr/local/spark/README.md")`
  - vi. `textFile.count`

```
scala> val textFile=sc.textFile("file:/usr/local/spark/README.md")
textFile: org.apache.spark.rdd.RDD[String] = file:/usr/local/spark/README.md MapPartitionsRDD
[1] at textFile at <console>:24

scala> textFile.count
res0: Long = 99
```

- 讀取HDFS的檔案(假設在HDFS上有一個/user/hduser/test/README.txt的檔案)

- vii. `val textFile=sc.textFile("hdfs://hadoopmaster:9000/user/hduser/test/README.txt")`
- viii. `textFile.count`

```
scala> val textFile=sc.textFile("hdfs://hadoopmaster:9000/user/hduser/test/README.txt")
textFile: org.apache.spark.rdd.RDD[String] = hdfs://hadoopmaster:9000/user/hduser/test/README.txt MapPartitionsRDD[1] at textFile at <console>:24

scala> textFile.count
res0: Long = 31
```

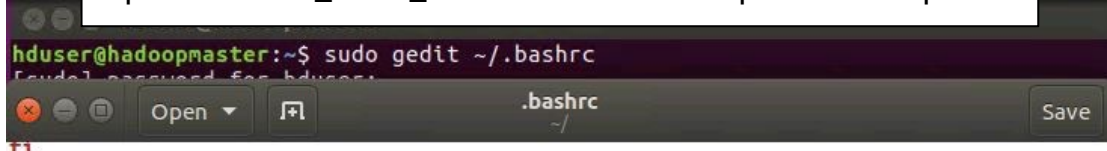
- 停止Spark Standalone cluster，輸入 `/usr/local/spark/sbin/stop-all.sh`

```
hduser@hadoopmaster:~$ /usr/local/spark/sbin/stop-all.sh
hadoopslave1: stopping org.apache.spark.deploy.worker.Worker
hadoopslave2: stopping org.apache.spark.deploy.worker.Worker
hadoopslave1: stopping org.apache.spark.deploy.worker.Worker
hadoopslave2: stopping org.apache.spark.deploy.worker.Worker
stopping org.apache.spark.deploy.master.Master
```

## 5. 在Hadoop YARN執行spark

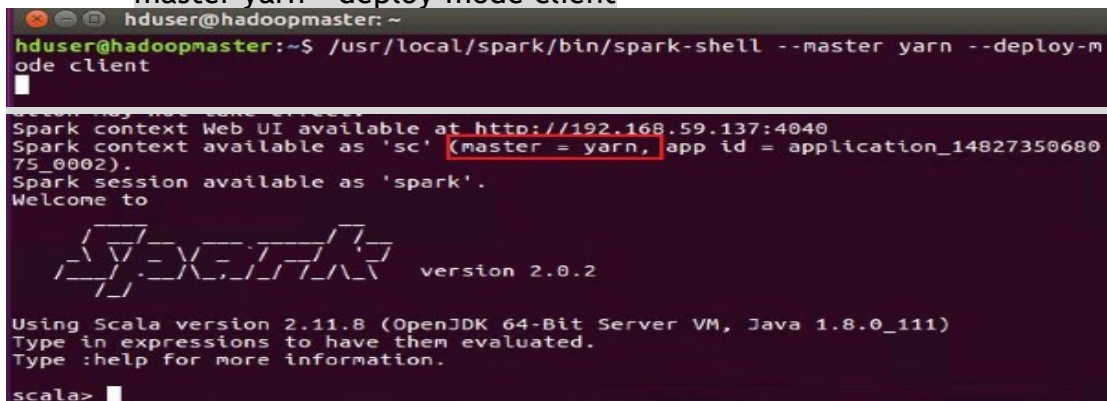
- 執行YARN需指定HADOOP\_CONF\_DIR參數，輸入 `sudo gedit ~/.bashrc`，加入下列參數

```
export HADOOP_CONF_DIR=/usr/local/hadoop/etc/hadoop
```



```
#Hadoop Variables
export JAVA_HOME=/usr/lib/jvm/java-8-openjdk-amd64
export HADOOP_HOME=/usr/local/hadoop
export HADOOP_CONF_DIR=/usr/local/hadoop/etc/hadoop
export PATH=$PATH:$HADOOP_HOME/bin
export PATH=$PATH:$HADOOP_HOME/sbin
export HADOOP_MAPRED_HOME=$HADOOP_HOME
```

- 讓~/.bashrc生效，輸入 `source ~/.bashrc`
- 在YARN上執行spark-shell，輸入 `/usr/local/spark/bin/spark-shell --master yarn --deploy-mode client`



- 測試讀取本機資料，輸入
  - i. `val textFile=sc.textFile("file:/usr/local/spark/README.md")`
  - ii. `textFile.count`

```
scala> val textFile=sc.textFile("file:/usr/local/spark/README.md")
textFile: org.apache.spark.rdd.RDD[String] = file:/usr/local/spark/README.md Map
PartitionsRDD[1] at textFile at <console>:24

scala> textFile.count
res0: Long = 99
```

- 測試讀取HDFS資料，輸入

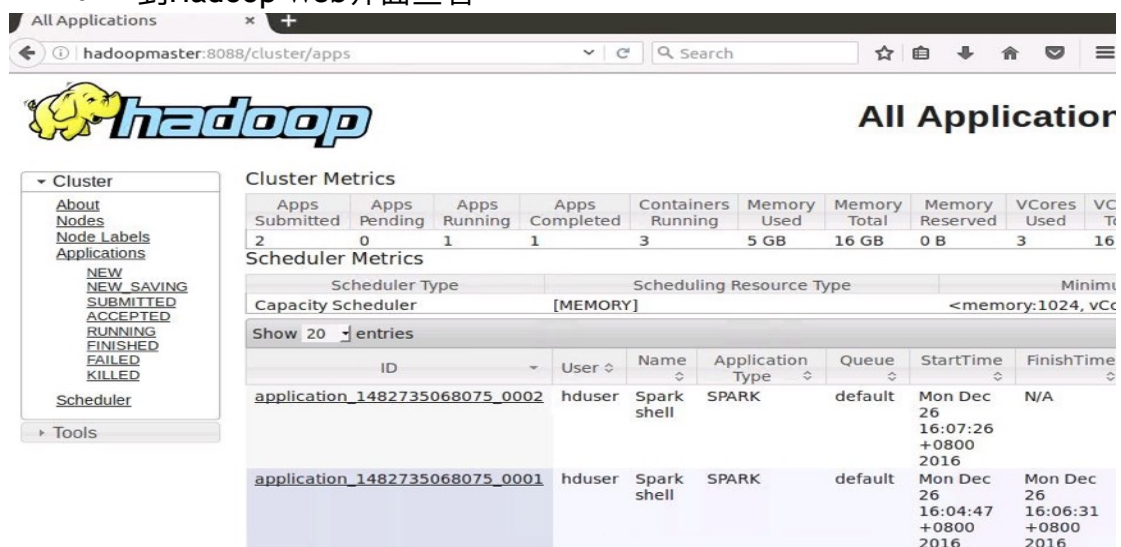
iii. `val textFile=sc.textFile("hdfs://hadoopmaster:9000/user/hduser/test/README.txt")`

iv. `textFile.count`

```
scala> val textFile=sc.textFile("hdfs://hadoopmaster:9000/user/hduser/test/README.txt")
textFile: org.apache.spark.rdd.RDD[String] = hdfs://hadoopmaster:9000/user/hduser/test/README.txt MapPartitionsRDD[3] at textFile at <console>:24

scala> textFile.count
res1: Long = 31
```

- 到Hadoop Web介面查看



The screenshot shows the Hadoop Web UI interface. On the left is a navigation menu with options like Cluster, About, Nodes, Node Labels, Applications, and Scheduler. The main content area displays 'Cluster Metrics' and 'Scheduler Metrics'. Below these, there is a table listing applications.

Cluster Metrics										
Apps Submitted	Apps Pending	Apps Running	Apps Completed	Containers Running	Memory Used	Memory Total	Memory Reserved	VCores Used	VCores Total	
2	0	1	1	3	5 GB	16 GB	0 B	3	16	

Scheduler Metrics							
Scheduler Type		Scheduling Resource Type				Minimum	
Capacity Scheduler		[MEMORY]				<memory:1024, vCores:16	
Show 20 entries							
ID	User	Name	Application Type	Queue	StartTime	FinishTime	
application_1482735068075_0002	hduser	Spark shell	SPARK	default	Mon Dec 26 16:07:26 +0800 2016	N/A	
application_1482735068075_0001	hduser	Spark shell	SPARK	default	Mon Dec 26 16:04:47 +0800 2016	Mon Dec 26 16:06:31 +0800 2016	