

搭建SPARK开发环境

任务概述

本章任务，主要是要建立对Spark引擎的一个整体认识，宏观上了解Spark的体系结构，知道Spark的应用场景，解决了什么问题。

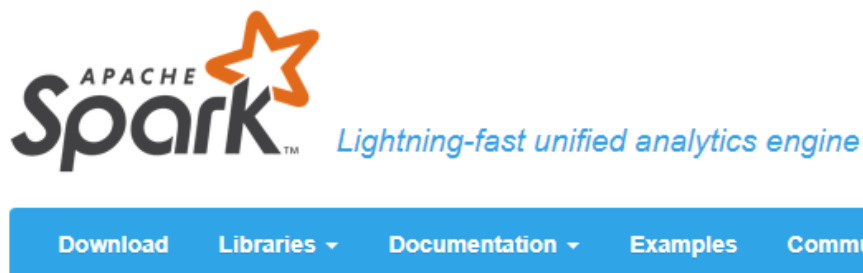
同时，还需掌握如何搭建Linux集群和Spark集群。

在基础设施建立完毕后，还需要知道如何在单机模式、伪分布式模式、集群模式上运行Spark应用。

1. 下载spark

spark有提供与hadoop相对应的预版本，下载时请根据自己需要进行选择，这里我们选择不依赖于具体hadoop的spark版本，在使用时我们再根据自己hadoop版本设置相应的环境参数。

spark目前最新的稳定版是2.4.5,版本3.0还是预览发行状态。



Download Apache Spark™

1. Choose a Spark release
2. Choose a package type
3. Download Spark [spark-2.4.5-bin-without-hadoop.tgz](http://mirror.bit.edu.cn/apache/spark/spark-2.4.5/spark-2.4.5-bin-without-hadoop.tgz)
4. Verify this release using the 2.4.5 [signatures](#), [checksums](#) and [project release KEYS](#).

Note that, Spark is pre-built with Scala 2.11 except version 2.4.2, which is pre-built with Scala 2.12.

```
1 | wget http://mirror.bit.edu.cn/apache/spark/spark-2.4.5/spark-2.4.5-bin-  
  | without-hadoop.tgz
```

注意：原连接使用的是https协议，下载速度可能很慢，这里改为使用http协议，速度会快很多，如果http协议被禁止不能使用，请改回https协议进行下载。

2. 解压spark文件到/opt,并建立相应的soft links文件

```
1 | tar zxvf spark-2.4.5-bin-without-hadoop.tgz  
2 | ln -sf $spark-2.4.5-bin-without-hadoop spark  
3 | ln -sf /opt/spark /usr/local/spark  
4 |
```

3.检查java及相关环境变量设置，这里我们仍以hadoop身份来运行spark，当然也可以单独指定一个帐号来运行，如果是单独一个运行帐号，后续在使用hdfs时，需要注意对hdfs分布式文件系统设置相应的权限。

```
1 java -version
2 rpm -qa | grep java    ##----> 1.8.0  jdk
3 env | grep JAVA
4 env | grep HOME
5
```

环境变量设置

vi .bash_profile 或 vi .bashrc

```
1 # User specific environment and startup programs
2 JAVA_HOME=/usr/lib/jvm/java
3 JRE_HOME=/usr/lib/jvm/jre
4 export JAVA_HOME
5 export JRE_HOME
6 export HADOOP_HOME=/opt/hadoop
7
8 #export SPARK_HOME=/usr/local/spark
9 export SPARK_HOME=/opt/spark
10 export PYTHONPATH=$SPARK_HOME/python:$SPARK_HOME/python/lib/py4j-0.10.8.1-
    src.zip
11 export PYSARK_PYTHON=python3
12
13 PATH=$PATH:$HOME/bin:$JAVA_HOME/bin:$JRE_HOME/bin:$HADOOP_HOME/bin:$HADOOP_
    HOME/sbin:$SPARK_HOME/bin
14 export PATH
15
16 alias vi=vim
```

4.编辑配置文件spark-env.sh

```
1 cp spark-env.sh.template  spark-env.sh
2 vi  spark-env.sh
```

spark-env.sh内容:

```
1 export JAVA_HOME=/usr/lib/jvm/java
2 export HADOOP_CONF_DIR=/opt/hadoop/etc/hadoop
3 export HADOOP_HOME=/opt/hadoop/
4
5 export SPARK_MASTER_HOST=master.lab.hwadee.com
6 #export SPARK_MASTER_PORT=7077
7 #export SPARK_MASTER_WEBUI_PORT=8080
8 #export SPARK_WORKER_CORES=1
9 #export SPARK_WORKER_MEMORY=1024M
10
11 export SPARK_PID_DIR=/var/lib/spark
12 export SPARK_LOG_DIR=/var/log/spark
13 export LD_LIBRARY_PATH=$HADOOP_HOME/lib/native/:$LD_LIBRARY_PATH
14 export SPARK_DIST_CLASSPATH=$( ${HADOOP_HOME}/bin/hadoop classpath)
```

```
15  
16
```

创建PID及日志目录并修改属主

```
1  mkdir /var/log/spark  
2  chown hadoop:hadoop /var/log/spark  
3  mkdir /var/lib/spark  
4  chown hadoop:hadoop /var/lib/spark
```

5.加入节点

```
1  vim  slaves  
2  datanode01.lab.hwadee.com  
3  datanode02.lab.hwadee.com  
4  datanode03.lab.hwadee.com  
5
```

6.同步资料到其它节点

修改文件属主为运行帐号:

```
1  chown hadoop:hadoop /opt/spark* -R
```

分发spark文件到其它节点:

a.)分发修改后登录脚本

```
1  scp ~hadoop/.bash_profile hadoop@datanode01.lab.hwadee.com:~/  
2  scp ~hadoop/.bash_profile hadoop@datanode02.lab.hwadee.com:~/  
3  scp ~hadoop/.bash_profile hadoop@datanode03.lab.hwadee.com:~/
```

b.) 分发spark程序到其它节点

```
1  scp -r /opt/spark-2.4.5-bin-without-hadoop  
hadoop@datanode01.lab.hwadee.com:/opt  
2  scp -r /opt/spark-2.4.5-bin-without-hadoop  
hadoop@datanode01.lab.hwadee.com:/opt  
3  scp -r /opt/spark-2.4.5-bin-without-hadoop  
hadoop@datanode01.lab.hwadee.com:/opt
```

c.)创建相关的link文件和目录

```
1  ssh datanode01 "cd /opt;ln -sf /spark-2.4.5-bin-without-hadoop hadoop;ln -  
sf /opt/spark /usr/local/spark"  
2  ssh datanode02 "cd /opt;ln -sf /spark-2.4.5-bin-without-hadoop hadoop;ln -  
sf /opt/spark /usr/local/spark"  
3  ssh datanode03 "cd /opt;ln -sf /spark-2.4.5-bin-without-hadoop hadoop;ln -  
sf /opt/spark /usr/local/spark"
```

```
1 | ssh datanode01 "mkdir /var/lib/spark;chown hadoop:hadoop /var/lib/spark"
2 | ssh datanode01 "mkdir /var/log/spark;chown hadoop:hadoop /var/log/spark"
3 |
4 | ssh datanode02 "mkdir /var/lib/spark;chown hadoop:hadoop /var/lib/spark"
5 | ssh datanode02 "mkdir /var/log/spark;chown hadoop:hadoop /var/log/spark"
6 |
7 | ssh datanode03 "mkdir /var/lib/spark;chown hadoop:hadoop /var/lib/spark"
8 | ssh datanode03 "mkdir /var/log/spark;chown hadoop:hadoop /var/log/spark"
```

7.启动spark

```
1 | /opt/spark/sbin/start-all.sh
```

8.测试验证:

put 文件到hdfs

vi 123.txt

```
1 | aaaa
2 | bbb
3 | dd
4 | ddd
5 | cc
6 | sss    sss
7 |
```

vi 456.txt

```
1 | 1111
2 | 2222
3 | 3333
4 | 4444
5 | 5555
```

```
1 | hadoop fs -put 123.txt /
2 | hadoop fs -put 456.txt /
3 | hadoop fs -ls /
```

9.运行 spark-shell

```
1  
2 val f1 = sc.textFile("hdfs://master.lab.hwadee.com/123.txt")  
3 val f2 = spark.read.textFile("/456.txt")  
4 f1.count()  
5 f2.count()  
6 f1.first()  
7 f2.first()
```

10.提交任务到spark

```
1 /opt/spark/bin/spark-submit --master spark://master.lab.hwadee.com:7077 --  
  class org.apache.spark.examples.SparkPi /opt/spark/examples/jars/spark-  
  examples_2.11-2.4.5.jar 2>/dev/null
```