

大数据课程软件统一说明

一、文档说明

为了统一我们的操作系统与软件环境，我们统一课前基本软件环境，实现全程学习当中的软件版本都是一致的

二、VmWare与linux版本

VmWare版本：

VmWare版本不做要求，使用VmWare10版本以上即可，关于VmWare的安装，直接使用安装包一直下一步安装即可，且安装包当中附带破解密钥，进行破解即可使用

linux版本

linux统一使用centos

centos统一使用centos7.6 64位版本

种子文件下载地址：http://mirrors.aliyun.com/centos/7.6.1810/isos/x86_64/CentOS-7-x86_64-DVD-1810.torrent

三、使用VmWare来安装linux软件

参见视频操作说明

四、三台linux服务器环境准备

使用三台linux服务器，来做统一的环境准备，通过统一环境，实现所有授课环节环境统一

三台机器IP设置

参见视频设置

三台机器修改ip地址：

```
1 vi /etc/sysconfig/network-scripts/ifcfg-ens33
2
3 BOOTPROTO="static"
4 IPADDR=192.168.52.100
5 NETMASK=255.255.255.0
6 GATEWAY=192.168.52.1
7 DNS1=8.8.8.8
```

准备三台linux机器，IP地址分别设置成为

第一台机器IP地址：192.168.52.100

第二台机器IP地址：192.168.52.110

第三台机器IP地址：192.168.52.120

三台机器关闭防火墙

三台机器在root用户下执行以下命令关闭防火墙

```
1 systemctl stop firewalld
2 systemctl disable firewalld
```

三台机器关闭selinux

三台机器在root用户下执行以下命令关闭selinux

三台机器执行以下命令，关闭selinux

```
1 vim /etc/selinux/config
2
3 SELINUX=disabled
```

三台机器更改主机名

三台机器分别更改主机名

第一台主机名更改为：node01.kaikeba.com

第二台主机名更改为：node02.kaikeba.com

第三台主机名更改为：node03.kaikeba.com

第一台机器执行以下命令修改主机名

```
1 vim /etc/hostname
2 node01.kaikeba.com
```

第二台机器执行以下命令修改主机名

```
1 vim /etc/hostname
2 node02.kaikeba.com
```

第三台机器执行以下命令修改主机名

```
1 vim /etc/hostname
2 node03.kaikeba.com
```

三台机器更改主机名与IP地址映射

三台机器执行以下命令更改主机名与IP地址映射关系

```
1 vim /etc/hosts
2
3 192.168.52.100 node01.kaikeba.com node01
4 192.168.52.110 node02.kaikeba.com node02
5 192.168.52.120 node03.kaikeba.com node03
```

三台机器同步时间

三台机器执行以下命令定时同步阿里云服务器时间

```
1 yum -y install ntpdate
2 crontab -e
3 */1 * * * * /usr/sbin/ntpdate time1.aliyun.com
4
5
```

三台机器添加普通用户

三台linux服务器统一添加普通用户hadoop，并给以sudo权限，用于以后所有的大数据软件的安装并统一设置普通用户的密码为 123456

```
1 useradd hadoop
2 passwd hadoop
```

三台机器为普通用户添加sudo权限

```
1 visudo
2
3 hadoop ALL=(ALL) ALL
```

三台定义统一目录

定义三台linux服务器软件压缩包存放目录，以及解压后安装目录，三台机器执行以下命令，创建两个文件夹，一个用于存放软件压缩包目录，一个用于存放解压后目录

```
1 mkdir -p /kkb/soft # 软件压缩包存放目录
2 mkdir -p /kkb/install # 软件解压后存放目录
3 chown -R hadoop:hadoop /kkb # 将文件夹权限更改为hadoop用户
```

三台机器安装jdk

使用hadoop用户来重新连接三台机器，然后使用hadoop用户来安装jdk软件

上传压缩包到第一台服务器的/kkb/soft下面，然后进行解压，配置环境变量即可，三台机器都依次安装即可

```
1 cd /kkb/soft/
2
3 tar -zxf jdk-8u141-linux-x64.tar.gz -C /kkb/install/
4 sudo vim /etc/profile
5
6
7 #添加以下配置内容，配置jdk环境变量
8 export JAVA_HOME=/kkb/install/jdk1.8.0_141
9 export PATH=$JAVA_HOME/bin:$PATH
```

hadoop用户免密码登录

三台机器在hadoop用户下执行以下命令生成公钥与私钥

```
1 ssh-keygen -t rsa
2 三台机器在hadoop用户下，执行以下命令将公钥拷贝到node01服务器上面去
3 ssh-copy-id node01
4 node01在hadoop用户下，执行以下命令，将authorized_keys拷贝到node02与node03服务器
5 cd /home/hadoop/.ssh/
6 scp authorized_keys node02:$PWD
7 scp authorized_keys node03:$PWD
```

三台机器关机重启

三台机器在root用户下执行以下命令，实现关机重启

```
1 reoot -h now
```

五、三台机器安装zookeeper集群

注意事项：**三台机器一定要保证时钟同步**

第一步：下载zookeeper的压缩包，下载网址如下

<http://archive.cloudera.com/cdh5/cdh/5/>

我们在这个网址下载我们使用的zk版本为[zookeeper-3.4.5-cdh5.14.2.tar.gz](http://archive.cloudera.com/cdh5/cdh/5/zookeeper-3.4.5-cdh5.14.2.tar.gz)

下载完成之后，上传到我们的node01的/kkb/soft路径下准备进行安装

第二步：解压

node01执行以下命令解压zookeeper的压缩包到node01服务器的/kkb/install路径下去，然后准备进行安装

```
1 | cd /kbb/soft
2 |
3 | tar -zxvf zookeeper-3.4.5-cdh5.14.2.tar.gz -C /kbb/install/
```

第三步：修改配置文件

第一台机器修改配置文件

```
1 | cd /kbb/install/zookeeper-3.4.5-cdh5.14.2/conf
2 |
3 | cp zoo_sample.cfg zoo.cfg
4 |
5 | mkdir -p /kbb/install/zookeeper-3.4.5-cdh5.14.2/zkdatas
6 |
7 | vim zoo.cfg
8 |
9 | dataDir=/kbb/install/zookeeper-3.4.5-cdh5.14.2/zkdatas
10 |
11 | autopurge.snapRetainCount=3
12 |
13 | autopurge.purgeInterval=1
14 |
15 | server.1=node01:2888:3888
16 | server.2=node02:2888:3888
17 | server.3=node03:2888:3888
```

第四步：添加myid配置

在第一台机器的/kbb/install/zookeeper-3.4.5-cdh5.14.2/zkdatas/

这个路径下创建一个文件，文件名为myid,文件内容为1

```
1 | echo 1 > /kbb/install/zookeeper-3.4.5-cdh5.14.2/zkdatas/myid
```

第五步：安装包分发并修改myid的值

安装包分发到其他机器

```
1 | 第一台机器上面执行以下两个命令
2 |
3 | scp -r /kbb/install/zookeeper-3.4.5-cdh5.14.2/ node02:/kbb/install/
4 |
5 | scp -r /kbb/install/zookeeper-3.4.5-cdh5.14.2/ node03:/kbb/install/
6 |
7 | 第二台机器上修改myid的值为2
8 |
9 | 直接在第二台机器任意路径执行以下命令
10 |
11 | echo 2 > /kbb/install/zookeeper-3.4.5-cdh5.14.2/myid
12 |
13 |
14 |
15 | 第三台机器上修改myid的值为3
```

```
16
17 直接在第三台机器任意路径执行以下命令
18
19 echo 3 > /kkb/install/zookeeper-3.4.5-cdh5.14.2/myid
```

第六步：三台机器启动zookeeper服务

三台机器启动zookeeper服务

这个命令三台机器都要执行

```
1 /kkb/install/zookeeper-3.4.5-cdh5.14.2/bin/zkServer.sh start
2
3 查看启动状态
4
5 /kkb/install/zookeeper-3.4.5-cdh5.14.2/bin/zkServer.sh status
```

六，hadoop环境安装

1、CDH软甲版本重新进行编译

1、为什么要编译hadoop

由于CDH的所有安装包版本都给出了对应的软件版本，一般情况下是不需要自己进行编译的，但是由于cdh给出的hadoop的安装包没有提供带C程序访问的接口，所以我们在使用本地库（本地库可以用来做压缩，以及支持C程序等等）的时候就会出问题，好了废话不多说，接下来看如何编译

2、编译环境的准备

2.1：准备linux环境

准备一台linux环境，内存4G或以上，硬盘40G或以上，我这里使用的是Centos6.9 64位的操作系统（注意：一定要使用64位的操作系统）

2.2：虚拟机联网，关闭防火墙，关闭selinux

```
1 关闭防火墙命令：
2
3 service iptables stop
4 chkconfig iptables off
5
6 关闭selinux
7 vim /etc/selinux/config
8
```

```
# This file controls the state of SELinux on the system.
# SELINUX= can take one of these three values:
#   enforcing - SELinux security policy is enforced.
#   permissive - SELinux prints warnings instead of enforcing.
#   disabled - No SELinux policy is loaded.
#SELINUX=enforcing
SELINUX=disabled
# SELINUXTYPE= can take one of these two values:
#   targeted - Targeted processes are protected,
#   mls - Multi Level Security protection.
SELINUXTYPE=targeted
```

2.3 : 安装jdk1.7

注意：亲测证明[hadoop-2.6.0-cdh5.14.2](#) 这个版本的编译，只能使用jdk1.7，如果使用jdk1.8那么就会报错

注意：这里一定不要使用jdk1.8，亲测jdk1.8会出现错误

将我们jdk的安装包上传到/kkb/soft（我这里使用的是jdk1.7.0_71这个版本）

解压我们的jdk压缩包

统一两个路径

```
1 mkdir -p /kkb/soft
2 mkdir -p /kkb/install
3 cd /kkb/soft
4 tar -zxvf jdk-7u71-linux-x64.tar.gz -C ../servers/
```

配置环境变量

```
1 vim /etc/profile
2
3 export JAVA_HOME=/kkb/install/jdk1.7.0_71
4
5 export PATH=$JAVA_HOME/bin:$PATH
```

```
1 让修改立即生效
2
3 source /etc/profile
```

2.4 : 安装maven

这里使用maven3.x以上的版本应该都可以，不建议使用太高的版本，强烈建议使用3.0.5的版本即可

将maven的安装包上传到/kkb/soft

然后解压maven的安装包到/kkb/install

```
1 cd /kkb/soft/
2
3 tar -zxvf apache-maven-3.0.5-bin.tar.gz -C ../servers/
```

配置maven的环境变量

```
1 vim /etc/profile
2
3 export MAVEN_HOME=/kkb/install/apache-maven-3.0.5
4
5 export MAVEN_OPTS="-Xms4096m -Xmx4096m"
6
7 export PATH=$MAVEN_HOME/bin:$PATH
```

```
1 | 让修改立即生效
2 |
3 | source /etc/profile
```

2.5 : 安装findbugs

下载findbugs

```
1 | cd /kkb/soft
2 |
3 | wget --no-check-certificate
  | https://sourceforge.net/projects/findbugs/files/findbugs/1.3.9/findbugs-
  | 1.3.9.tar.gz/download -O findbugs-1.3.9.tar.gz
4 |
5 |
6 |
7 | 解压findbugs
8 |
9 | tar -zxvf findbugs-1.3.9.tar.gz -C ../install/
10 |
11 |
12 |
13 | 配置findbugs的环境变量
14 |
15 | vim /etc/profile
16 |
17 | export JAVA_HOME=/kkb/install/jdk1.7.0_75
18 |
19 | export PATH=$JAVA_HOME/bin:$PATH
20 |
21 |
22 |
23 | export MAVEN_HOME=/kkb/install/apache-maven-3.0.5
24 |
25 | export PATH=$MAVEN_HOME/bin:$PATH
26 |
27 |
28 | export FINDBUGS_HOME=/kkb/install/findbugs-1.3.9
29 | export PATH=$FINDBUGS_HOME/bin:$PATH
```

```
1 | 让修改立即生效
2 |
3 | source /etc/profile
```

2.6 : 在线安装一些依赖包


```

1 yum install autoconf automake libtool cmake
2
3 yum install ncurses-devel
4
5 yum install openssl-devel
6
7 yum install lzo-devel zlib-devel gcc gcc-c++
8
9
10
11 bzip2压缩需要的依赖包
12
13 yum install -y bzip2-devel

```

2.7 : 安装protobuf

protobuf下载百度网盘地址

<https://pan.baidu.com/s/1pJlZubT>

下载之后上传到 /kkb/soft

解压protobuf并进行编译

```

1 cd /kkb/soft
2
3 tar -zxvf protobuf-2.5.0.tar.gz -C ../servers/
4
5 cd /kkb/install/protobuf-2.5.0
6
7 ./configure
8
9 make && make install

```

2.8、安装snappy

snappy下载地址：

<http://code.google.com/p/snappy/>

```

1 cd /kkb/soft/
2
3 tar -zxf snappy-1.1.1.tar.gz -C ../servers/
4
5 cd ../servers/snappy-1.1.1/
6
7 ./configure
8
9 make && make install

```

2.9 : 下载cdh源码准备编译

源码下载地址为：

<http://archive.cloudera.com/cdh5/cdh/5/hadoop-2.6.0-cdh5.14.2-src.tar.gz>

下载源码进行编译

```
1 cd /kkb/soft
2
3 wget http://archive.cloudera.com/cdh5/cdh/5/hadoop-2.6.0-cdh5.14.2-
  src.tar.gz
4
5 tar -zxvf hadoop-2.6.0-cdh5.14.2-src.tar.gz -C ../servers/
6
7 cd /kkb/install/hadoop-2.6.0-cdh5.14.2
8
9 编译不支持snappy压缩:
10
11 mvn package -Pdist,native -DskipTests -Dtar
12
13
14
15 编译支持snappy压缩:
16
17 mvn package -DskipTests -Pdist,native -Dtar -Drequire.snappy -e -X
18
19 编译完成之后我们需要的压缩包就在下面这个路径里面
```

2.10 : 常见编译错误

如果编译时候出现这个错误：

An Ant BuildException has occurred: exec returned: 2

这是因为tomcat的压缩包没有下载完成，需要自己下载一个对应版本的apache-tomcat-6.0.53.tar.gz的压缩包放到指定路径下面去即可

这两个路径下面需要放上这个tomcat的 压缩包

```
1 /kkb/install/hadoop-2.6.0-cdh5.14.2/hadoop-hdfs-project/hadoop-hdfs-
  httpfs/downloads
2
3 /kkb/install/hadoop-2.6.0-cdh5.14.2/hadoop-common-project/hadoop-
  kms/downloads
```

2、hadoop集群的安装

安装环境服务部署规划

服务器IP	192.168.52.100	192.168.52.110	192.168.52.120
HDFS	NameNode		
HDFS	SecondaryNameNode		
HDFS	DataNode	DataNode	DataNode
YARN	ResourceManager		
YARN	NodeManager	NodeManager	NodeManager

服务器IP	NodeManager	NodeManager	NodeManager
历史日志服务器	192.168.52.100 JobHistoryServer	192.168.52.110	192.168.52.120

第一步：上传压缩包并解压

将我们重新编译之后支持snappy压缩的hadoop包上传到第一台服务器并解压

第一台机器执行以下命令

```
1 cd /kbb/soft/
2
3
4
5 tar -zxvf hadoop-2.6.0-cdh5.14.2_after_compile.tar.gz -C ../install/
```

第二步：查看hadoop支持的压缩方式以及本地库

第一台机器执行以下命令

```
1 cd /kbb/install/hadoop-2.6.0-cdh5.14.2
2
3 bin/hadoop checknative
```

如果出现openssl为false，那么所有机器在线安装openssl即可，执行以下命令，虚拟机联网之后就可以在线进行安装了

```
1 yum -y install openssl-devel
```

第三步：修改配置文件

修改core-site.xml

第一台机器执行以下命令

```
1 cd /kbb/install/hadoop-2.6.0-cdh5.14.2/etc/hadoop
2
3 vim core-site.xml
4
5 <configuration>
6     <property>
7         <name>fs.defaultFS</name>
8         <value>hdfs://node01:8020</value>
9     </property>
10    <property>
11        <name>hadoop.tmp.dir</name>
12        <value>/kbb/install/hadoop-2.6.0-
13        cdh5.14.2/hadoopDatas/tempDatas</value>
14    </property>
15    <!-- 缓冲区大小，实际工作中根据服务器性能动态调整 -->
16    <property>
17        <name>io.file.buffer.size</name>
18        <value>4096</value>
19    </property>
20    <!-- 开启hdfs的垃圾桶机制，删除掉的数据可以从垃圾桶中回收，单位分钟 -->
21    <property>
```

```

21         <name>fs.trash.interval</name>
22         <value>10080</value>
23     </property>
24 </configuration>

```

修改hdfs-site.xml

第一台机器执行以下命令

```

1  cd /kbb/install/hadoop-2.6.0-cdh5.14.2/etc/hadoop
2
3  vim hdfs-site.xml
4
5  <configuration>
6      <!-- NameNode存储元数据信息的路径，实际工作中，一般先确定磁盘的挂载目录，然后多个目
        录用，进行分割 -->
7      <!-- 集群动态上下线
8          <property>
9              <name>dfs.hosts</name>
10             <value>/kbb/install/hadoop-2.6.0-
11             cdh5.14.2/etc/hadoop/accept_host</value>
12         </property>
13         <property>
14             <name>dfs.hosts.exclude</name>
15             <value>/kbb/install/hadoop-2.6.0-
16             cdh5.14.2/etc/hadoop/deny_host</value>
17         </property>
18         -->
19         <property>
20             <name>dfs.namenode.secondary.http-address</name>
21             <value>node01:50090</value>
22         </property>
23         <property>
24             <name>dfs.namenode.http-address</name>
25             <value>node01:50070</value>
26         </property>
27         <property>
28             <name>dfs.namenode.name.dir</name>
29             <value>file:///kbb/install/hadoop-2.6.0-
30             cdh5.14.2/hadoopDatas/namenodeDatas</value>
31         </property>
32         <!-- 定义dataNode数据存储的节点位置，实际工作中，一般先确定磁盘的挂载目录，然后多
33             个目录用，进行分割 -->
34         <property>
35             <name>dfs.datanode.data.dir</name>
36             <value>file:///kbb/install/hadoop-2.6.0-
37             cdh5.14.2/hadoopDatas/datanodeDatas</value>
38         </property>
39         <property>
40             <name>dfs.namenode.edits.dir</name>
41             <value>file:///kbb/install/hadoop-2.6.0-
42             cdh5.14.2/hadoopDatas/dfs/nn/edits</value>
43         </property>
44         <property>
45             <name>dfs.namenode.checkpoint.dir</name>
46             <value>file:///kbb/install/hadoop-2.6.0-
47             cdh5.14.2/hadoopDatas/dfs/snn/name</value>

```

```

41     </property>
42     <property>
43         <name>dfs.namenode.checkpoint.edits.dir</name>
44         <value>file:///kbb/install/hadoop-2.6.0-
cdh5.14.2/hadoopData/dfs/nn/snn/edits</value>
45     </property>
46     <property>
47         <name>dfs.replication</name>
48         <value>2</value>
49     </property>
50     <property>
51         <name>dfs.permissions</name>
52         <value>false</value>
53     </property>
54 <property>
55     <name>dfs.blocksize</name>
56     <value>134217728</value>
57 </property>
58 </configuration>
59

```

修改hadoop-env.sh

第一台机器执行以下命令

```

1  cd /kbb/install/hadoop-2.6.0-cdh5.14.2/etc/hadoop
2
3  vim hadoop-env.sh
4  export JAVA_HOME=/kbb/install/jdk1.8.0_141
5

```

修改mapred-site.xml

第一台机器执行以下命令

```

1  cd /kbb/install/hadoop-2.6.0-cdh5.14.2/etc/hadoop
2
3  vim mapred-site.xml
4
5  <configuration>
6      <property>
7          <name>mapreduce.framework.name</name>
8          <value>yarn</value>
9      </property>
10     <property>
11         <name>mapreduce.job.ubertask.enable</name>
12         <value>true</value>
13     </property>
14     <property>
15         <name>mapreduce.jobhistory.address</name>
16         <value>node01:10020</value>
17     </property>
18     <property>
19         <name>mapreduce.jobhistory.webapp.address</name>
20         <value>node01:19888</value>

```

```
21     </property>
22 </configuration>
23
```

修改yarn-site.xml

第一台机器执行以下命令

```
1  cd /kbb/install/hadoop-2.6.0-cdh5.14.2/etc/hadoop
2
3  vim yarn-site.xml
4
5  <configuration>
6      <property>
7          <name>yarn.resourcemanager.hostname</name>
8          <value>node01</value>
9      </property>
10     <property>
11         <name>yarn.nodemanager.aux-services</name>
12         <value>mapreduce_shuffle</value>
13     </property>
14 </configuration>
15
```

修改slaves文件

第一台机器执行以下命令

```
1  cd /kbb/install/hadoop-2.6.0-cdh5.14.2/etc/hadoop
2
3  vim slaves
4  node02
5  node03
6
```

第四步：创建文件存放目录

第一台机器执行以下命令

node01机器上面创建以下目录

```
1  mkdir -p /kbb/install/hadoop-2.6.0-cdh5.14.2/hadoopDatas/tempDatas
2  mkdir -p /kbb/install/hadoop-2.6.0-cdh5.14.2/hadoopDatas/namenodeDatas
3  mkdir -p /kbb/install/hadoop-2.6.0-cdh5.14.2/hadoopDatas/datanodeDatas
4  mkdir -p /kbb/install/hadoop-2.6.0-cdh5.14.2/hadoopDatas/dfs/nn/edits
5  mkdir -p /kbb/install/hadoop-2.6.0-cdh5.14.2/hadoopDatas/dfs/snn/name
6  mkdir -p /kbb/install/hadoop-2.6.0-cdh5.14.2/hadoopDatas/dfs/nn/snn/edits
```

第五步：安装包的分发

第一台机器执行以下命令

```
1 | cd /kbb/install/
2 |
3 | scp -r hadoop-2.6.0-cdh5.14.2/ node02:$PWD
4 | scp -r hadoop-2.6.0-cdh5.14.2/ node03:$PWD
5 |
```

第六步：配置hadoop的环境变量

三台机器都要进行配置hadoop的环境变量

三台机器执行以下命令

```
1 | vim /etc/profile
2 |
3 | export HADOOP_HOME=/kbb/install/hadoop-2.6.0-cdh5.14.2
4 | export PATH=$HADOOP_HOME/bin:$HADOOP_HOME/sbin:$PATH
5 |
6 | 配置完成之后生效
7 |
8 | source /etc/profile
9 |
```

第七步：集群启动

要启动 Hadoop 集群，需要启动 HDFS 和 YARN 两个集群。

注意：首次启动HDFS时，必须对其进行格式化操作。本质上是一些清理和准备工作，因为此时的HDFS在物理上还是不存在的。

```
1 | bin/hdfs namenode -format或者bin/hadoop namenode -format
```

单个节点逐一启动

```
1 | 在主节点上使用以下命令启动 HDFS NameNode:
2 | hadoop-daemon.sh start namenode
3 |
4 | 在每个从节点上使用以下命令启动 HDFS DataNode:
5 | hadoop-daemon.sh start datanode
6 |
7 | 在主节点上使用以下命令启动 YARN ResourceManager:
8 | yarn-daemon.sh start resourcemanager
9 |
10 | 在每个从节点上使用以下命令启动 YARN nodemanager:
11 | yarn-daemon.sh start nodemanager
12 |
13 | 以上脚本位于$HADOOP_PREFIX/sbin/目录下。如果想要停止某个节点上某个角色，只需要把命令中的
    | start 改为stop 即可。
14 |
```

脚本一键启动

如果配置了 `etc/hadoop/slaves` 和 `ssh` 免密登录，则可以使用程序脚本启动所有Hadoop 两个集群的相关进程，在主节点所设定的机器上执行。

启动集群

node01节点上执行以下命令

```
1  第一台机器执行以下命令
2
3  cd /kbb/install/hadoop-2.6.0-cdh5.14.2/
4  sbin/start-dfs.sh
5  sbin/start-yarn.sh
6  sbin/mr-jobhistory-daemon.sh start historyserver
7
8  停止集群:
9
10 sbin/stop-dfs.sh
11
12 sbin/stop-yarn.sh
13
```

第八步：浏览器查看启动页面

hdfs集群访问地址

<http://192.168.52.100:50070/dfshealth.html#tab-overview>

yarn集群访问地址

<http://192.168.52.100:8088/cluster>

jobhistory访问地址：

<http://192.168.52.100:19888/jobhistory>.

