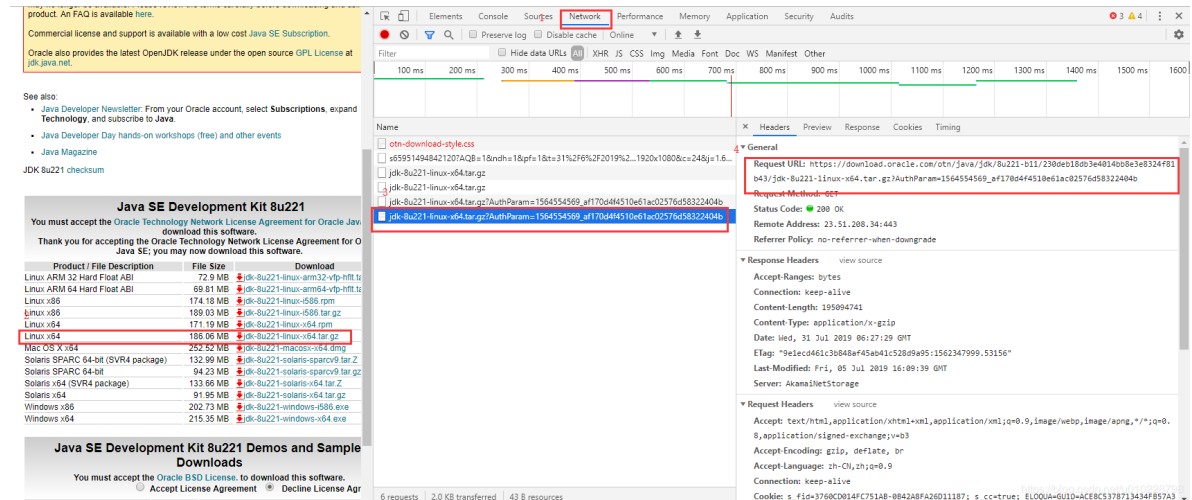


# 安装jdk

使用wget时会验证失败，因为需要登陆Oracle

所以可以用google浏览器桌面式安装时使用 network的方式获取下载链接



```
1 wget url
```

- 创建文件夹(根据自己喜好)

```
1 mkdir /usr/local/software/java/文件名
```

- 解压至自己的目录

```
1 mv 安装包路径 /usr/local/software/java
2 tar xvfz 安装包路径
```

- 配置环境变量

```
1 vi /etc/profile
2
3 export JAVA_HOME=/usr/local/software/java/jdk1.8.0_231
4 export CLASSPATH=$JAVA_HOME/lib/
5 export PATH=$PATH:$JAVA_HOME/bin
6
7 source /etc/profile
8
9 java -version
```

# 安装Scala

从[这里](#)下载对应版本的安装包

```
1 wget https://downloads.lightbend.com/scala/2.13.1/scala-2.13.1.tgz
```

## Other resources

You can find the installer download links for other operating systems, as well as documentation and source code archives for Scala 2.13.1 below.

Archive	System	Size
<a href="#">scala-2.13.1.tgz</a>	Mac OS X, Unix, Cygwin	18.77M
<a href="#">scala-2.13.1.msi</a>	Windows (msi installer)	115.13M
<a href="#">scala-2.13.1.zip</a>	Windows	18.81M
<a href="#">scala-2.13.1.deb</a>	Debian	582.81M
<a href="#">scala-2.13.1.rpm</a>	RPM package	115.52M
<a href="#">scala-docs-2.13.1.txz</a>	API docs	48.58M
<a href="#">scala-docs-2.13.1.zip</a>	API docs	99.67M
<a href="#">scala-sources-2.13.1.tar.gz</a>	Sources	

```
1 vi /etc/profile
2 export SCALA_HOME=/usr/local/software/scala-2.13.1
3 export PATH=$PATH:$SCALA_HOME/bin
4
5 source /etc/profile
6 scala -version
```

## 配置静态ip及DNS

```
1 vi /etc/sysconfig/network-scripts/ifcfg-ens33
2
3 TYPE=Ethernet
4 PROXY_METHOD=none
5 BROWSER_ONLY=no
6 BOOTPROTO=static # 修改
7 DEFROUTE=yes
8 IPV4_FAILURE_FATAL=no
9 IPV6INIT=yes
10 IPV6_AUTOCONF=yes
11 IPV6_DEFROUTE=yes
12 IPV6_FAILURE_FATAL=no
13 IPV6_ADDR_GEN_MODE=stable-privacy
14 NAME=ens33
15 UUID=5de61f71-a87c-4961-a678-7d9a69c4704f
16 DEVICE=ens33
17 ONBOOT=yes # 修改
18 IPV6_PRIVACY=no
19 #ip
20 IPADDR=192.168.197.135
21 #网关
```

```
22 | GATEWAY=192.168.197.2
23 | #子网掩码
24 | PREFIX0=24
25 | #使用主的DNS
26 | DNS1=8.8.8.8
27 | #备用的DNS
28 | DNS2=8.8.4.4
```

```
1 | systemctl restart network.service # 生效
2 | ping baidu.com # 测试
```

## 关闭防火墙

```
1 | systemctl stop firewalld # 临时关闭防火墙
2 | systemctl disable firewalld # 禁止开机启动
```

## 配置hosts文件

```
1 | vim /etc/hosts
2 |
3 | 192.168.197.129 master
4 | 192.168.197.130 slaver1
```



```
zzh@master: ~
文件(F) 编辑(E) 查看(V) 搜索(S) 终端(T) 帮助(H)
127.0.0.1      localhost
127.0.1.1      master

# The following lines are desirable for IPv6 capable hosts
::1          ip6-localhost ip6-loopback
fe00::0      ip6-localnet
ff00::0      ip6-mcastprefix
ff02::1      ip6-allnodes
ff02::2      ip6-allrouters

192.168.197.129 master
192.168.197.130 slaver1
~
~
~
```

## 添加用户账号

在所有的主机下均建立一个账号admin用来运行hadoop，并将其添加至sudoers中

```
1 useradd admin # 添加用户通过手动输入修改密码
2 passwd admin # 更改用户 admin 的密码
```

设置admin用户具有root权限 修改 /etc/sudoers 文件，找到下面一行，在root下面添加一行

```
1 visudo
2
3 ## Allow root to run any commands anywhere
4 root    ALL=(ALL)    ALL
5 admin   ALL=(ALL)    ALL
```

用admin帐号登录，然后用命令 su - ，切换用户即可获得root权限进行操作。

```
1 su - admin # 切换admin用户
2 su - # 切换root用户
```

# ssh免密

## • 配置master和slaver

- 切换root用户
- 配置每台主机的SSH免密码登录环境
- 在每台主机上生成公钥和私钥对

```
1 # 对每台主机
2 ssh-keygen -t rsa
```

- 将slaver上的id\_rsa.pub发送给master

```
1 scp ~/.ssh/id_rsa.pub root@master:~/.ssh/id_rsa.pub.slaver1
```

- 在master上，将所有公钥加载到用于认证的公钥文件authorized\_key中，并查看生成的文件

```
1 cat ~/.ssh/id_rsa.pub* >> ~/.ssh/authorized_keys
2 cd .ssh
3 ls
```

- 将master上的公钥文件authorized\_key分发给slaver

```
1 scp ~/.ssh/authorized_keys root@slaver1:~/.ssh/
```

- 最后使用SSH命令，检验是否能免密码登录。

```
1 ssh slaver1
```

- 进入admin用户，免密

```
1 su - admin
2 ssh-keygen -t rsa
3 ssh-copy-id master #本机也需要
4 ssh-copy-id slaver1
```

```
1
2
```













## 创建软件安装目录

root用户下

```
1 mkdir /usr/local/software
2 chown -R admin:admin /usr/local/software #赋予权限
```

## 安装Hadoop

- 从[这里](#)下载对应版本的安装包

<u>Name</u>	<u>Last modified</u>	<u>Size</u>	<u>Description</u>
 <a href="#">Parent Directory</a>		-	
 <a href="#">current/</a>	2019-10-20 20:44	-	
 <a href="#">current2/</a>	2019-09-24 06:32	-	
 <a href="#">hadoop-2.10.0/</a>	2019-10-29 18:22	-	
 <a href="#">hadoop-2.7.7/</a>	2018-07-19 18:12	-	
 <a href="#">hadoop-2.8.5/</a>	2018-09-18 09:13	-	
 <a href="#">hadoop-2.9.2/</a>	2018-11-19 21:45	-	
 <a href="#">hadoop-3.1.3/</a>	2019-10-20 20:44	-	
 <a href="#">hadoop-3.2.1/</a>	2019-09-24 06:32	-	
 <a href="#">stable/</a>	2019-09-24 06:32	-	
 <a href="#">stable2/</a>	2019-09-24 06:32	-	
 <a href="#">readme.txt</a>	2015-04-20 18:32	184	

```
1 wget http://apache.claz.org/hadoop/common/hadoop-3.2.1/hadoop-3.2.1.tar.gz
```

- 配置环境变量

```

1 vi /etc/profile
2
3 export HADOOP_HOME=/usr/local/software/hadoop-3.2.1
4 export HADOOP_CONF_DIR=$HADOOP_HOME/etc/hadoop
5 export PATH=$PATH:$HADOOP_HOME/bin
6
7 source /etc/profile

```

- 修改hadoop-env.sh、yarn-env.sh，添加JAVA\_HOME

```

1 export JAVA_HOME=/usr/local/software/java/jdk1.8.0_231

```

```

52 # The java implementation to use. By default, this environment
53 # variable is REQUIRED on ALL platforms except OS X!
54 export JAVA_HOME=/home/zzh/software/java/jdk1.8.0_231
55

```

```

169 # export YARN_CONTAINER_RUNTIME_DOCKER_RUN_OVERRIDE_DISABLE=true
170
171 export JAVA_HOME=/home/zzh/software/java/jdk1.8.0_231

```

- 配置 core-site.xml

```

1 <configuration>
2   <property>
3     <name>fs.defaultFS</name>
4     <value>hdfs://master:9000</value>
5   </property>
6   <property>
7     <name>hadoop.tmp.dir</name>
8     <value>/usr/local/software/hadoop-3.2.1/data</value>
9   </property>
10 </configuration>

```

- 配置 hdfs-site.xml (创建好datanode namenode目录)

```

1 <configuration>
2   <!--指定hdfs数据的冗余份数 默认是3-->
3   <property>
4     <name>dfs.replication</name>
5     <value>2</value>
6   </property>
7   <property>
8     <name>dfs.name.dir</name>
9     <value>/usr/local/software/hadoop-3.2.1/hdfs/namenode</value>
10  </property>
11  <property>
12    <name>dfs.data.dir</name>
13    <value>/usr/local/software/hadoop-3.2.1/hdfs/datanode</value>
14  </property>
15 </configuration>

```

- 配置 mapred-site.xml

```

1 <configuration>
2   <property>
3     <name>mapreduce.framework.name</name>
4     <value>yarn</value>
5   </property>
6   <property>
7     <name>mapred.job.tracker</name>
8     <value>http://master:9001</value>
9   </property>
10 </configuration>

```

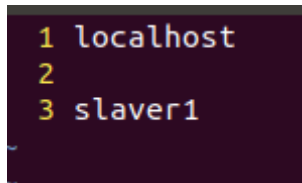
- 配置 yarn-site.xml

```

1 <configuration>
2
3 <!-- Site specific YARN configuration properties -->
4   <property>
5     <name>yarn.nodemanager.aux-services</name>
6     <value>mapreduce_shuffle</value>
7   </property>
8   <property>
9     <name>yarn.resourcemanager.hostname</name>
10    <value>master</value>
11  </property>
12 </configuration>

```

- 修改workers(hadoop2.x的是slavers)，配置slaver节点的ip或hostname



```

1 localhost
2
3 slaver1

```

- 将配置好的hadoop-3.2.1文件分发给所有slaver

```

1 scp -r /usr/local/software/hadoop-3.2.1/
  admin@slaver1:/usr/local/software/

```

scp目标服务器注意写法（用户名@主机ip）

- 进入bin目录并初始化

```

1 cd /usr/local/software/hadoop-3.2.1/bin
2 ./hadoop namenode -format

```

- 启动hadoop集群

```

1 cd /usr/local/software/hadoop-3.2.1/sbin
2 start-dfs.sh
3 start-yarn.sh

```

## 安装conda

- 由于ubuntu自带了python3.6，就不用安装了，可以安装pip

```
1 | sudo apt install python3-pip
```

- 选择适合自己的版本，用wget命令下载

```
1 | wget -c https://repo.continuum.io/miniconda/Miniconda3-latest-Linux-x86_64.sh
```

这里选择的是 latest-Linux 版本，所以下载的程序会随着python的版本更新而更新

- ```
1 | bash Miniconda3-latest-Linux-x86_64.sh #运行
```

- 换清华源

```
1 | vi ~/.condarc
2
3 | channels:
4 |   - https://mirrors.ustc.edu.cn/anaconda/pkg/main/
5 |   - https://mirrors.ustc.edu.cn/anaconda/cloud/conda-forge/
6 |   - https://mirrors.tuna.tsinghua.edu.cn/anaconda/pkg/free/
7 |   - defaults
8 | show_channel_urls: true
```

- 创建环境

```
1 | cd /usr/local/software/miniconda3/bin
2 | ./conda create -n env_name list of packages
3
4 | # ps: conda create -n my_env numpy
```

在这里，`-n env_name` 设置环境的名称（`-n` 是指名称），而 `list of packages` 是要安装在环境中的包的列表。

创建环境时，可以指定要安装在环境中的 Python 版本。这在你同时使用 Python 2.x 和 Python 3.x 中的代码时很有用。

```
1 | conda create -n env_name python=3.6
```

- 进入环境

```
1 | cd /usr/local/software/miniconda3/bin/
2
3
4 | #在 Linux 上使用下面进入环境:
5 | source activate my_env
6 | #在 windows 上使用下面进入环境:
7 | activate my_env
```

可以使用 `conda list` 检查安装的包



```
1 # 安装包的命令
2 conda install package_name
```

- 离开环境

```
1 #在 Linux 上请键入:
2 source deactivate。
3 #在 windows 上键入:
4 deactivate
```

参考链接：<https://blog.csdn.net/abc13526222160/article/details/84624334>

# Jupyter运行conda环境

## 安装jupyter

```
1 pip3 install jupyter # 安装命令
2
3 jupyter notebook --generate-config # 在当前用户路径 ~/创建 .jupyter目录
```

- 进入python环境获取密匙

```
1 from notebook.auth import passwd
2 passwd() # 'sha1:0a07349659f1:85120e097cbc4e3ccb91c7936640fb8be95d1537'
```

- 进入jupyter目录下的 jupyter\_notebook\_config.py 文件

```
1 vi ~/.jupyter/ jupyter_notebook_config.py
```

```
1 c.NotebookApp.ip='*' # 就是设置所有ip皆可访问
2 c.NotebookApp.password = u'sha:ce...' # 刚才复制的那个密文'
3 c.NotebookApp.open_browser = False # 禁止自动打开浏览器
4 c.NotebookApp.port =8888 #随便指定一个端口
5 c.NotebookApp.notebook_dir = u'' # 指定工作目录
```

- 进入待激活环境

```
1 service activate test
```

- 安装nb\_conda

```
1 conda install nb_conda
```

- 添加虚拟环境至notebook

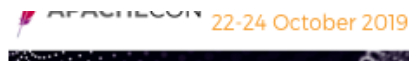
```
1 python -m ipykernel install --user --name 虚拟环境名 --display-name "显示环境名"
2 #python -m ipykernel install --user --name test --display-name "python test"
```

- 在待激活环境里启动notebook

```
1 jupyter notebook
```

## 安装Spark

- 从[这里](#)下载对应版本的安装包



### Download Apache Spark™

1. Choose a Spark release:
2. Choose a package type:
3. Download Spark: [spark-3.0.0-preview-bin-hadoop3.2.tgz](#)
4. Verify this release using the 3.0.0-preview [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 https://www-eu.apache.org/dist/spark/spark-3.0.0-preview/spark-3.0.0-preview-bin-hadoop3.2.tgz
```

- 解压后进入conf目录

```
1 cd /usr/local/software/spark/conf
```

在该目录下，看到很多文件都是以template结尾的，这是因为spark给我们提供的是模板配置文件，我们可以先拷贝一份，然后将.template给去掉，变成真正的配置文件后再编辑。

```
zzh@master:~/software/spark/conf$ ls
fairscheduler.xml.template  slaves.template
log4j.properties.template  spark-defaults.conf.template
metrics.properties.template spark-env.sh.template
```

- 配置spark-env.sh，该文件包含spark的各种运行环境

```

1 | cp spark-env.sh.template spark-env.sh
2 | vi spark-env.sh
3
4 | export SCALA_HOME=/usr/local/software/scala-2.13.1
5 | export JAVA_HOME=/usr/local/software/java/jdk1.8.0_231
6 | export HADOOP_INSTALL=/usr/local/software/hadoop-3.2.1
7 | export HADOOP_CONF_DIR=$HADOOP_INSTALL/etc/hadoop
8 | SPARK_MASTER_IP=master # 得用ip格式
9 | SPARK_MASTER_HOST=master # 得用ip格式
10 | SPARK_LOCAL_IP=master
11 | SPARK_LOCAL_DIRS=/usr/local/software/spark-3.0.0-preview-bin-hadoop3.2
12 | SPARK_DRIVER_MEMORY=2G

```

- 配置slaves文件

```

1 | cp slaves.template slaves
2 | vi slaves
3
4 | slaver1

```

```

18 # Hadoop
19 localhost
20
21 slaver1

```

- 将配置好的spark-2.3.1文件分发给所有slaver

```

1 | scp -r /usr/local/software/spark/ admin@slaver1:/usr/local/software/

```

- 启动spark集群

- 切换用户admin
- 进入hadoop目录，在该目录下，启动 hadoop 文件管理系统 HDFS以及启动 hadoop 任务管理器 YARN。

```

1 | cd /usr/local/software/hadoop-3.2.1/sbin
2 | start-dfs.sh
3 | start-yarn.sh

```

- 启动spark

```

1 | cd /usr/local/software/spark/sbin
2 | ./start-all.sh

```

- 查看Spark集群信息

- 使用jps命令
- 查看spark管理界面，在浏览器中输入：<http://master:8080>

- 运行 spark-shell，可以进入 Spark 的 shell 控制台

- 停止运行集群

停止集群时，运行sbin/stop-all.sh停止Spark集群，运行sbin/stop-dfs.sh来关闭hadoop 文件管理系统 HDFS，最后运行sbin/stop-yarn.sh来关闭hadoop 任务管理器 YARN。

# Jupyter里运行PySpark

---

用findSpark包

```
1 # 需要配置好SPARK_HOME环境变量
2 import findspark
3 findspark.init()
4
5 # 下面再使用pyspark函数
```

## 参考链接

---

[https://blog.csdn.net/qg\\_15349687/article/details/82748074](https://blog.csdn.net/qg_15349687/article/details/82748074)

## 踩坑

---

### root@zzh-0: Permission denied (publickey,password)

---

解决办法，复制导入公钥就可以了，SSH链接需要使用公钥认证：

切换到ssh目录：cd ~/.ssh/

```
1 ssh-keygen -t rsa -P "" (回车)
2 cat $HOME/.ssh/id_rsa.pub >> $HOME/.ssh/authorized_keys
```

再次启动HDFS就可以了

```
1 sbin/start-dfs.sh
```

参考链接：<https://yq.aliyun.com/articles/695939>

### but there is no HDFS\_NAMENODE\_USER defined. Aborting operation

---

```

root@zzh-1:/usr/local/software/hadoop/hadoop-3.1.3# ./sbin/start-all.sh
Starting namenodes on [zzh-1]
ERROR: Attempting to operate on hdfs namenode as root
ERROR: but there is no HDFS_NAMENODE_USER defined. Aborting operation.
Starting datanodes
ERROR: Attempting to operate on hdfs datanode as root
ERROR: but there is no HDFS_DATANODE_USER defined. Aborting operation.
Starting secondary namenodes [zzh-1]
ERROR: Attempting to operate on hdfs secondarynamenode as root
ERROR: but there is no HDFS_SECONDARYNAMENODE_USER defined. Aborting operation.
WARNING: YARN_CONF_DIR has been replaced by HADOOP_CONF_DIR. Using value of YARN_CONF_DIR.
Starting resourcemanager
ERROR: Attempting to operate on yarn resourcemanager as root
ERROR: but there is no YARN_RESOURCEMANAGER_USER defined. Aborting operation.
Starting nodemanagers
ERROR: Attempting to operate on yarn nodemanager as root
ERROR: but there is no YARN_NODEMANAGER_USER defined. Aborting operation.
root@zzh-1:/usr/local/software/hadoop/hadoop-3.1.3#

```

在start-dfs.sh和stop-dfs.sh中：

```

1 HDFS_DATANODE_USER=root
2 HADOOP_SECURE_DN_USER=hdfs
3 HDFS_NAMENODE_USER=root
4 HDFS_SECONDARYNAMENODE_USER=root

```

```

#!/usr/bin/env bash
HDFS_DATANODE_USER=root
HADOOP_SECURE_DN_USER=hdfs
HDFS_NAMENODE_USER=root
HDFS_SECONDARYNAMENODE_USER=root
# Licensed to the Apache Software Foundation (ASF) under one or more

```

在start-yarn.sh和stop-yarn.sh中：

```

1 YARN_RESOURCEMANAGER_USER=root
2 HADOOP_SECURE_DN_USER=yarn
3 YARN_NODEMANAGER_USER=root

```

```

#!/usr/bin/env bash
YARN_RESOURCEMANAGER_USER=root
HADOOP_SECURE_DN_USER=yarn
YARN_NODEMANAGER_USER=root
# Licensed to the Apache Software Foundation (ASF) under one or more

```

## The authenticity of host '0.0.0.0 (0.0.0.0)' can't be established.

解决方案关闭SELINUX

```
1  -- 关闭SELINUX
2  vim /etc/selinux/config
3
4  -- 注释掉
5  #SELINUX=enforcing
6  #SELINUXTYPE=targeted
7  - 添加
8  SELINUX=disabled
```

[<https://www.oracle.com/technetwork/java/javase/downloads/jdk8-downloads-2133151.html> ]: