

1 环境准备

在虚拟机安装Cent OS 6.x的操作系统（我这里是CentOS 6.6）。

1.1 修改主机名

修改主机名为spark1234，并修改配置文件确保重启后主机名依然生效。

在配置文件/etc/sysconfig/network修改主机名为spark1234。

```
1 # hostname spark1234
2 # vim /etc/sysconfig/network
3 NETWORKING=yes
4 HOSTNAME=spark1234
```

配置主机名解析：

在/etc/hosts配置主机ip和主机名的解析

```
[root@spark1234 ~]# vim /etc/hosts
```

```
127.0.0.1    localhost localhost.localdomain localhost4 localhost4.localdomain4
::1         localhost localhost.localdomain localhost6 localhost6.localdomain6
192.168.9.14 spark1234
```

1.2 关闭防火墙

关闭iptables：

```
1 # service iptables stop
2 # chkconfig iptables off
3 # chkconfig --list | grep iptables
4 iptables          0:关闭    1:关闭    2:关闭    3:关闭    4:关闭    5:关闭    6:关闭
5
```

关闭selinux：

```
1 # setenforce 0
2 # vim /etc/sysconfig/selinux
```

```
root@spark1234 ~]# vim /etc/sysconfig/selinux

# 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=disabled
# SELINUXTYPE= can take one of these two values:
#   targeted - Targeted processes are protected,
#   mls - Multi Level Security protection.
SELINUXTYPE=targeted
```

1.3 安装jdk

1) 首先卸载openjdk

```
1 --查看java版本
2 [root@spark1234 ~]# java -version
3 java version "1.7.0_45"
4 OpenJDK Runtime Environment (rhel-2.4.3.3.el6-x86_64 u45-b15)
5 OpenJDK 64-Bit Server VM (build 24.45-b08, mixed mode)
6
7 --查看安装源
8 [root@spark1234 ~]# rpm -qa | grep java
9 java-1.7.0-openjdk-1.7.0.45-2.4.3.3.el6.x86_64
10 tzdata-java-2013g-1.el6.noarch
11 java-1.6.0-openjdk-1.6.0.0-1.66.1.13.0.el6.x86_64
12
13 -- 卸载
14 [root@spark1234 ~]# rpm -e --nodeps java-1.7.0-openjdk-1.7.0.45-2.4.3.3.el6.x86_64
15 [root@spark1234 ~]# rpm -e --nodeps tzdata-java-2013g-1.el6.noarch
16 [root@spark1234 ~]# rpm -e --nodeps java-1.6.0-openjdk-1.6.0.0-
17 1.66.1.13.0.el6.x86_64
18
```

```
19 --验证是否卸载成功
20 [root@spark1234 ~]# rpm -qa | grep java
21 [root@spark1234 ~]# java -version
22 -bash: /usr/bin/java: 没有那个文件或目录
```

2) 安装java

```
1 -- 下载并解压java源码包
2 [root@spark1234 java]# mkdir /usr/local/java
3 [root@spark1234 java]# mv jdk-7u79-linux-x64.tar.gz /usr/local/java
4 [root@spark1234 java]# cd /usr/local/java
5 [root@spark1234 java]# tar xvf jdk-7u79-linux-x64.tar.gz
6 [root@spark1234 java]# ls
7 jdk1.7.0_79  jdk-7u79-linux-x64.tar.gz
8 [root@spark1234 java]#
```

3) 设置环境变量

```
1 [root@spark1234 java]# vim /etc/profile
2 [root@spark1234 java]# tail /etc/profile
3 export JAVA_HOME=/usr/local/java/jdk1.7.0_79
4 export JRE_HOME=/usr/local/java/jdk1.7.0_79/jre
5 export
6 CLASSPATH=.:$JAVA_HOME/lib/dt.jar:$JAVA_HOME/lib/tools.jar:$JRE_HOME/lib:$CLASSPATH
7
8 export PATH=$JAVA_HOME/bin:$PATH
9
10 -- 生效环境变量
11 [root@spark1234 ~]# source /etc/profile
```

4) 验证

```
1 -- 验证
2 [root@spark1234 ~]# java -version
3 java version "1.7.0_79"
4 Java(TM) SE Runtime Environment (build 1.7.0_79-b15)
5 Java HotSpot(TM) 64-Bit Server VM (build 24.79-b02, mixed mode)
6 [root@spark1234 ~]# javac -version
7 javac 1.7.0_79
```

1.4 创建hadoop用户并配置免密码认证

(1)、创建hadoop用户

```
1 -- 创建hadoop用户
2 [root@spark1234 ~]# useradd hadoop
3 -- 设置hadoop用户的密码为hadoop
4 [root@spark1234 ~]# echo "hadoop" | passwd --stdin hadoop
```

(2) 配置免密码登陆

如果不配置密码也可以，但是每次启动hadoop服务都需要输入密码，建议配置

```
1 -- 切换到hadoop用户
2 [root@spark1234 ~]# su - hadoop
3
4 -- 生成公钥，一路回车
5 [hadoop@spark1234 ~]$ /usr/bin/ssh-keygen -t rsa -N ""
6
7 -- 将公钥内容写入文件authorized_keys
8 [hadoop@spark1234 ~]$ cat ~/.ssh/id_rsa.pub >> ~/.ssh/authorized_keys
9
10 -- 验证authorized_keys文件的权限，必须是644，如果不是，需要修改，否则免密失败
11 [hadoop@spark1234 ~]$ ll ~/.ssh/authorized_keys
12 -rw-rw-r-- 1 hadoop hadoop 398 11月 14 21:32 /home/hadoop/.ssh/authorized_keys
13 -- 修改authorized_keys文件的权限
14 [hadoop@spark1234 ~]$ chmod 644 ~/.ssh/authorized_keys
15
16 -- 验证免密
17 [hadoop@spark1234 ~]$ ssh spark1234
18 Last login: Tue Nov 14 21:33:59 2017 from spark1234
19 [hadoop@spark1234 ~]$
```

2. 大数据相关组件安装

2.1 软件清单

软件	版本
hadoop	hadoop-2.6.0-cdh5.6.0
hive	hive-1.1.0-cdh5.6.0
spark	spark-1.6.3-bin-2.6.0-cdh5.6.0
mysql	mysql-5.6

(1) hadoop、hive、spark可从cdh的官网下载，链接：<http://archive.cloudera.com/cdh5/cdh/5/>

(2) mysql采用rpm包安装

下载链接：https://cdn.mysql.com//Downloads/MySQL-5.6/MySQL-5.6.38-1.el6.x86_64.rpm-bundle.tar

2.2 hadoop安装

(1) 解压安装包

```
1 $ tar -xzf hadoop-2.6.0-cdh5.6.0.tar.gz -C /home/hadoop/app/
```

(2) 设置环境变量

在hadoop用户下，在~/.bashrc文件增加两行：

```
1 [hadoop@spark1234 ~]$ vim ~/.bashrc
2 增加如下两行内容：
3 export HADOOP_HOME=/home/hadoop/app/hadoop-2.6.0-cdh5.6.0
4 export PATH=$HADOOP_HOME/bin:$HADOOP_HOME/sbin:$PATH
```

```
[hadoop@spark1234 ~]$ cat ~/.bashrc
# .bashrc
```

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

```
# User specific aliases and functions
export HADOOP_HOME=/home/hadoop/app/hadoop-2.6.0-cdh5.6.0
export PATH=$HADOOP_HOME/bin:$HADOOP_HOME/sbin:$PATH
[hadoop@spark1234 ~]$
```

生效：

```
1 [hadoop@spark1234 ~]$ source ~/.bashrc
```

(3)、修改core-site.xml文件

切换到目录：\$HADOOP_HOME/etc/hadoop

编辑core-site.xml，增加配置：

```
1 <configuration>
2 <property>
3     <name>fs.defaultFS</name>
4     <value>hdfs://spark1234:8020</value>
5 </property>
6 <property>
7     <name>hadoop.tmp.dir</name>
8     <value>/home/hadoop/app/hadoop-2.6.0-cdh5.6.0/tmp</value>
9 </property>
10 </configuration>
```

新建配置文件里面配置的目录：

一定要在hadoop用户下创建，或者hadoop用户有权限操作这个目录

```
1 [hadoop@spark1234 ~]# mkdir /home/hadoop/app/hadoop-2.6.0-cdh5.6.0/tmp
```

(4) 修改配置文件hdfs-site.xml

切换到目录：\$HADOOP_HOME/etc/hadoop

编辑hdfs-site.xml，增加配置：

```
1 <configuration>
2 <property>
3     <name>dfs.replication</name>
4     <value>1</value>
5 </property>
6 </configuration>
```

(5) 修改配置文件hadoop-env.sh

切换到目录：\$HADOOP_HOME/etc/hadoop

修改JAVA_HOME的配置：

export JAVA_HOME=\${JAVA_HOME} 修改为： export
JAVA_HOME=/usr/local/java/jdk1.7.0_79

```
# The only required environment variable is JAVA_HOME. All  
# optional. When running a distributed configuration it :  
# set JAVA_HOME in this file, so that it is correctly def:  
# remote nodes.
```

```
# The java implementation to use.
```

```
export JAVA_HOME=/usr/local/java/jdk1.7.0_79
```

```
# The jsvc implementation to use. Jsvc is required to run  
# that bind to privileged ports to provide authentication
```

(6) 修改配置文件mapred-site.xml

切换到目录：\$HADOOP_HOME/etc/hadoop

```
1 [hadoop@spark1234 hadoop]$ cp mapred-site.xml.template mapred-site.xml
2 [hadoop@spark1234 hadoop]$
3 [hadoop@spark1234 hadoop]$ vim mapred-site.xml
```

编辑mapred-site.xml，增加配置：

```
1 <configuration>
2   <property>
3     <name>mapreduce.framework.name</name>
4     <value>yarn</value>
5   </property>
6 </configuration>
```

(7) 修改配置文件yarn-site.xml

增加如下配置：

```
1 <configuration>
2   <property>
3     <name>yarn.log-aggregation-enable</name>
4     <value>true</value>
5   </property>
```

```

6      <property>
7          <name>yarn.log.server.url</name>
8          <value>http://spark1234:19888/jobhistory/job/</value>
9      </property>
10     <property>
11         <name>yarn.nodemanager.aux-services</name>
12         <value>mapreduce_shuffle</value>
13     </property>
14 </configuration>

```

(8) 格式化

```
1 [hadoop@spark1234 ~]$ hdfs namenode -format
```

```

7/11/14 23:33:40 INFO namenode.FSNamesystem: dfs.namenode.safemode.threshold-pct = 0.9990000128746033
7/11/14 23:33:40 INFO namenode.FSNamesystem: dfs.namenode.safemode.min.datanodes = 0
7/11/14 23:33:40 INFO namenode.FSNamesystem: dfs.namenode.safemode.extension = 30000
7/11/14 23:33:40 INFO metrics.TopMetrics: NNTop conf: dfs.namenode.top.window.num.buckets = 10
7/11/14 23:33:40 INFO metrics.TopMetrics: NNTop conf: dfs.namenode.top.num.users = 10
7/11/14 23:33:40 INFO metrics.TopMetrics: NNTop conf: dfs.namenode.top.windows.minutes = 1,5,25
7/11/14 23:33:40 INFO namenode.FSNamesystem: Retry cache on namenode is enabled
7/11/14 23:33:40 INFO namenode.FSNamesystem: Retry cache will use 0.03 of total heap and retry cache entry expiry time is 600000 millis
7/11/14 23:33:40 INFO util.GSet: Computing capacity for map NameNodeRetryCache
7/11/14 23:33:40 INFO util.GSet: VM type = 64-bit
7/11/14 23:33:40 INFO util.GSet: 0.029999999329447746% max memory 966.7 MB - 297.0 KB
7/11/14 23:33:40 INFO util.GSet: capacity = 2^15 = 32768 entries
7/11/14 23:33:40 INFO namenode.NNConf: ACLs enabled? false
7/11/14 23:33:40 INFO namenode.NNConf: XAttrs enabled? true
7/11/14 23:33:40 INFO namenode.NNConf: Maximum size of an xattr: 16384
7/11/14 23:33:40 INFO namenode.FSImage: Allocated new BlockPoolId: BP-1078578980-192.168.9.14-1510673620265
7/11/14 23:33:40 INFO common.Storage: Storage directory /home/hadoop/app/hadoop-2.6.0-cdh5.6.0/tmp/dfs/name has been successfully formatted.
7/11/14 23:33:40 INFO namenode.NNStorageRetentionManager: Going to retain 1 images with txid >= 0
7/11/14 23:33:40 INFO util.ExitUtil: Exiting with status 0
7/11/14 23:33:40 INFO namenode.NameNode: SHUTDOWN_MSG:
.....
SHUTDOWN_MSG: Shutting down NameNode at spark1234/192.168.9.14
.....

```

(9) 启动服务

```

1 -- 启动hdfs服务
2 [hadoop@spark1234 ~]$ start-dfs.sh
3 -- 启动yarn服务
4 [hadoop@spark1234 ~]$ start-yarn.sh
5 -- 验证服务
6 [hadoop@spark1234 ~]$ jps

```

使用jps查看进程：


```
[hadoop@spark1234 ~]$ jps
8040 NameNode
8328 SecondaryNameNode
8133 DataNode
8470 ResourceManager
8569 NodeManager
8882 Jps
```

查看hdfs的web ui :

<http://192.168.9.14:50070/dfshealth.html#tab-overview>

The screenshot shows the HDFS web UI Overview page. The top navigation bar includes links for Hadoop, Overview (selected), Datanodes, Datanode Volume Failures, Snapshot, Startup Progress, and Utilities. The main content area is titled 'Overview 'spark1234:8020' (active)'. It contains a table with the following information:

Started:	Tue Nov 14 23:35:19 CST 2017
Version:	2.6.0-cdh5.6.0, rc282dc6c30e7d5d27410cabb328d60fc24266d9
Compiled:	2016-01-29T05:41Z by jenkins from Unknown
Cluster ID:	CID-cc4ef9b1-034b-4431-b4f0-2e0b06211b54
Block Pool ID:	BP-1078578980-192.168.9.14-1510673620265

Below the table is a 'Summary' section. It states: 'Security is off.', 'Safemode is off.', '1 files and directories, 0 blocks = 1 total filesystem object(s).', 'Heap Memory used 35.72 MB of 51.78 MB Heap Memory. Max Heap Memory is 966.69 MB.', and 'Non Heap Memory used 29.88 MB of 31 MB Committed Non Heap Memory. Max Non Heap Memory is 130 MB.' At the bottom, a table shows 'Configured Capacity: 34.92 GB'.

查看yarn的web ui :

<http://192.168.9.14:8088/cluster>

The screenshot shows the YARN web UI Cluster Metrics page. The top navigation bar includes links for About, Nodes, Applications (selected), NEW, NEW SAVING, SUBMITTED, ACCEPTED, RUNNING, FINISHED, FAILED, KILLED, and Scheduler. The main content area is titled 'All Applications'. It contains a table with the following information:

Apps Submitted	Apps Pending	Apps Running	Apps Completed	Containers Running	Memory Used	Memory Total	Memory Reserved	VCores Used	VCores Total
0	0	0	0	0	0 B	8 GB	0 B	0	8

Below the table is a 'User Metrics for dr.who' section. It contains a table with the following information:

Apps Submitted	Apps Pending	Apps Running	Apps Completed	Containers Running	Containers Pending	Containers Reserved	Memory
0	0	0	0	0	0	0	0 B

Below the table is a 'Show 20 entries' section. It contains a table with the following information:

ID	User	Name	Application Type	Queue	StartTime	FinishTime	State	FinalStatus	Running Containers
No data available in table									

At the bottom, it says 'Showing 0 to 0 of 0 entries'.

2.3 MySQL安装

(1) 卸载现有的mysql

```

1 [root@spark1234 ~]# rpm -qa | grep mysql
2 mysql-5.1.73-3.el6_5.x86_64
3 mysql-libs-5.1.73-3.el6_5.x86_64
4 mysql-devel-5.1.73-3.el6_5.x86_64
5 [root@spark1234 ~]#
6 [root@spark1234 ~]# rpm -e --nodeps mysql-5.1.73-3.el6_5.x86_64
7 [root@spark1234 ~]# rpm -e --nodeps mysql-libs-5.1.73-3.el6_5.x86_64
8 [root@spark1234 ~]# rpm -e --nodeps mysql-devel-5.1.73-3.el6_5.x86_64
9

```

(2)、将下载的[MySQL-5.6.38-1.el6.x86_64.rpm-bundle.tar](#)上传到主机

```

1 --解压
2 # tar -xf MySQL-5.6.38-1.el6.x86_64.rpm-bundle.tar
3
4 [root@spark1234 test]# ll
5 总用量 445708
6 -rw-r--r-- 1 root root 228198400 11月 14 13:52 MySQL-5.6.38-1.el6.x86_64.rpm-
bundle.tar
7 -rw-r--r-- 1 7155 31415 19045216 9月 14 19:00 MySQL-client-5.6.38-
1.el6.x86_64.rpm
8 -rw-r--r-- 1 7155 31415 3391312 9月 14 19:00 MySQL-devel-5.6.38-
1.el6.x86_64.rpm
9 -rw-r--r-- 1 7155 31415 90407828 9月 14 19:01 MySQL-embedded-5.6.38-
1.el6.x86_64.rpm
10 -rw-r--r-- 1 7155 31415 57551568 9月 14 19:01 MySQL-server-5.6.38-
1.el6.x86_64.rpm
11 -rw-r--r-- 1 7155 31415 1964616 9月 14 19:01 MySQL-shared-5.6.38-
1.el6.x86_64.rpm
12 -rw-r--r-- 1 7155 31415 3969744 9月 14 19:01 MySQL-shared-compat-5.6.38-
1.el6.x86_64.rpm
13 -rw-r--r-- 1 7155 31415 51861916 9月 14 19:01 MySQL-test-5.6.38-1.el6.x86_64.rpm
14
15 -- 安装
16 [root@spark1234 test]# rpm -ivh MySQL*.rpm
17
18 -- 修改配置文件的位置
19 [root@spark1234 test]# cp /usr/share/mysql/my-default.cnf /etc/my.cnf
20

```

(3) 初始化MySQL并设置密码

```

1  -- 初始化mysql
2  [root@spark1234 test]# /usr/bin/mysql_install_db
3
4  -- 启动mysql
5  [root@spark1234 test]# service mysql start
6
7  -- 查看root账号的密码
8  [root@spark1234 test]# cat /root/.mysql_secret
9  # The random password set for the root user at Tue Nov 14 22:01:34 2017 (local
    time): KAxcgzc3R_qVJTB7
10
11 -- 进入MySQL， 并修改密码
12 [root@spark1234 test]# mysql -uroot -pKAxcgzc3R_qVJTB7
13
14 -- 设置初始密码为123456
15 mysql> SET PASSWORD = PASSWORD('123456');
16 Query OK, 0 rows affected (0.03 sec)
17
18 mysql> exit
19

```

设置权限：

```

1  mysql> grant all privileges on *.* to root@'localhost' identified by '123456' with
    grant option;
2  mysql> grant all privileges on *.* to root@'127.0.0.1' identified by '123456' with
    grant option;
3  mysql> grant all privileges on *.* to root@'%' identified by '123456' with grant
    option;
4  mysql> flush privileges;

```

(4) 设置开机启动

```

1  [root@spark1234 test]# chkconfig mysql on
2  [root@spark1234 test]# chkconfig --list | grep mysql
3  mysql          0:关闭    1:关闭    2:启用    3:启用    4:启用    5:启用    6:关闭

```

2.4 hive的安装

(1) 解压安装包

```
1 $ tar -xzf hive-1.1.0-cdh5.6.0.tar.gz -C /home/hadoop/app/
```

(2) 设置环境变量

在hadoop用户下，在~/.bashrc文件增加两行：

```
1 [hadoop@spark1234 ~]$ vim ~/.bashrc
2 增加如下两行内容：
3 export HIVE_HOME=/home/hadoop/app/hive-1.1.0-cdh5.6.0
4 export PATH=$HIVE_HOME/bin:$PATH
```

```
hadoop@spark1234 ~$ cat ~/.bashrc
# .bashrc

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

# User specific aliases and functions
export HADOOP_HOME=/home/hadoop/app/hadoop-2.6.0-cdh5.6.0
export PATH=$HADOOP_HOME/bin:$HADOOP_HOME/sbin:$PATH
export HIVE_HOME=/home/hadoop/app/hive-1.1.0-cdh5.6.0
export PATH=$HIVE_HOME/bin:$PATH
[hadoop@spark1234 ~]$
[hadoop@spark1234 ~]$
```

生效：

```
1 [hadoop@spark1234 ~]$ source ~/.bashrc
```

(3)、在mysql创建hive使用的表

```
1 -- 创建hive数据库
2 mysql> create database hive;
3 -- 设置编码，一定要设置成latin1，否则hive建表和删表会卡住
4 mysql> alter database hive character set latin1;
5
```

```

[hadoop@spark1234 ~]$
[hadoop@spark1234 ~]$ mysql -uroot -p123456
Warning: Using a password on the command line interface can be insecure.
Welcome to the MySQL monitor.  Commands end with ; or \g.
Your MySQL connection id is 10
Server version: 5.6.38 MySQL Community Server (GPL)

Copyright (c) 2000, 2017, Oracle and/or its affiliates. All rights reserved.

Oracle is a registered trademark of Oracle Corporation and/or its
affiliates. Other names may be trademarks of their respective
owners.

Type 'help;' or '\h' for help. Type '\c' to clear the current input statement.

mysql> create database hive;
Query OK, 1 row affected (0.06 sec)

mysql> alter database hive character set latin1;
Query OK, 1 row affected (0.00 sec)

mysql> |

```

(4) 修改hive的配置文件

在 \$HIVE_HOME/conf 目录创建文件hive-site.xml

配置如下：

```

1  <?xml version="1.0" encoding="UTF-8"?>
2
3  <!--Autogenerated by Cloudera Manager-->
4  <configuration>
5      <property>
6          <name>javax.jdo.option.ConnectionURL</name>
7          <value>jdbc:mysql://127.0.0.1:3306/hive?
createDatabaseIfNotExist=true</value>
8      </property>
9
10     <property>
11         <name>javax.jdo.option.ConnectionDriverName</name>
12         <value>com.mysql.jdbc.Driver</value>
13     </property>
14
15     <property>
16         <name>javax.jdo.option.ConnectionUserName</name>
17         <value>root</value>
18     </property>
19
20     <property>
21         <name>javax.jdo.option.ConnectionPassword</name>
22         <value>123456</value>
23     </property>
24 </configuration>

```

(5) 将mysql的jdk驱动包复制到hive的lib目录下

```
1 $ cp mysql-connector-java-5.1.44-bin.jar $HIVE_HOME/lib/
```

(6) 进入hive

进入hive，并验证可用：

```
1 [hadoop@spark1234 soft]$ hive
2 hive> show tables;
3 hive> create table test(id int, name string);
4 hive> insert into test values (1, "abc");
5 hive> select * from test;
```

```
[hadoop@spark1234 soft]$ hive
```

```
which: no hbase in (/home/hadoop/app/hive-1.1.0-cdh5.6.0/bin:/home/hadoop/app/hadoop-2.6.0-
in:/home/hadoop/app/hadoop-2.6.0-cdh5.6.0/bin:/home/hadoop/app/hadoop-2.6.0-cdh5.6.0/sbin:/
usr/local/bin:/bin:/usr/bin:/usr/local/sbin:/usr/sbin:/sbin:/home/hadoop/bin)
```

```
Logging initialized using configuration in jar:file:/home/hadoop/app/hive-1.1.0-cdh5.6.0/li
s
```

WARNING: Hive CLI is deprecated and migration to Beeline is recommended.

```
hive>
```

```
> show tables;
```

```
OK
```

```
Time taken: 1.127 seconds
```

```
hive> create table test(id int, name string);
```

```
OK
```

```
Time taken: 1.173 seconds
```

```
hive>
```

```
> insert into test values (1, 'abc');
```

```
Query ID = hadoop_20171114235050_aacb55f1-1786-46bf-8880-dbdde924630e
```

```
Total jobs = 3
```

```
Launching Job 1 out of 3
```

```
Number of reduce tasks is set to 0 since there's no reduce operator
```

```
Starting Job = job_1510673764000_0001, Tracking URL = http://spark1234:8088/pro
```

```
Kill Command = /home/hadoop/app/hadoop-2.6.0-cdh5.6.0/bin/hadoop job -kill job
```

```
Hadoop job information for Stage-1: number of mappers: 1; number of reducers: 0
```

```
2017-11-14 23:50:33,599 Stage-1 map = 0%, reduce = 0%
```

```
2017-11-14 23:50:48,872 Stage-1 map = 100%, reduce = 0%, Cumulative CPU 1.91 s
```

```
MapReduce Total cumulative CPU time: 1 seconds 910 msec
```

```
Ended Job = job_1510673764000_0001
```

```
Stage-4 is selected by condition resolver.
```

```
Stage-3 is filtered out by condition resolver.
```

```
Stage-5 is filtered out by condition resolver.
```

```
Moving data to: hdfs://spark1234:8020/user/hive/warehouse/test/.hive-staging_hi
```

```
Loading data to table default.test
```

```
Table default.test stats: [numFiles=1, numRows=1, totalSize=6, rawDataSize=5]
```

```
MapReduce Jobs Launched:
```

```
Stage-Stage-1: Map: 1 Cumulative CPU: 1.91 sec HDFS Read: 3421 HDFS Write:
```

```
Total MapReduce CPU Time Spent: 1 seconds 910 msec
```

```
OK
```

```
Time taken: 45.14 seconds
```

```
hive>
```

```
> select * from test;
```

```
OK
```

```
1 abc
```

```
Time taken: 0.43 seconds, Fetched: 1 row(s)
```

```
hive> |
```

2.5 spark安装

(1) 解压安装包

```
1 $ tar -xzf spark-1.6.3-bin-2.6.0-cdh5.6.0.tgz -C /home/hadoop/app/
```

(2) 设置环境变量

在hadoop用户下，在~/.bashrc文件增加两行：

```
1 export SPARK_HOME=/home/hadoop/app/spark-1.6.3-bin-2.6.0-cdh5.6.0
2 export PATH=$SPARK_HOME/bin:$SPARK_HOME/sbin:$PATH
```

```
[hadoop@spark1234 ~]$ cat .bashrc
# .bashrc

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

# User specific aliases and functions
export HADOOP_HOME=/home/hadoop/app/hadoop-2.6.0-cdh5.6.0
export PATH=$HADOOP_HOME/bin:$HADOOP_HOME/sbin:$PATH
export HIVE_HOME=/home/hadoop/app/hive-1.1.0-cdh5.6.0
export PATH=$HIVE_HOME/bin:$PATH
export SPARK_HOME=/home/hadoop/app/spark-1.6.3-bin-2.6.0-cdh5.6.0
export PATH=$SPARK_HOME/bin:$SPARK_HOME/sbin:$PATH
[hadoop@spark1234 ~]$
```

环境变量生效：

```
1 [hadoop@spark1234 ~]$ source .bashrc
```

(3) 修改配置文件

进入目录：\$SPARK_HOME/conf目录

```
1 [hadoop@spark1234 conf]$ cp spark-env.sh.template spark-env.sh
```

增加指定hadoop配置文件的配置：

```
export HADOOP_CONF_DIR=/home/hadoop/app/hadoop-2.6.0-cdh5.6.0/etc/hadoop
```

```
# 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: ${SPARK_HOME}/
# - SPARK_PID_DIR       Where the pid file is stored. (Default: /tmp)
# - SPARK_IDENT_STRING  A string representing this instance of spark. (Default:
# - SPARK_NICENESS       The scheduling priority for daemons. (Default: 0)
export HADOOP_CONF_DIR=/home/hadoop/app/hadoop-2.6.0-cdh5.6.0/etc/hadoop
[hadoop@spark1234 conf]$
```

(4) 将hive的配置文件复制到spark的配置目录下
这样spark就可以直接读取hive中的表

```
1 [hadoop@spark1234 conf]$ cp $HIVE_HOME/conf/hive-site.xml $SPARK_HOME/conf/
```

(5) 进入spark

```
1 $ spark-shell --master yarn --jars /home/hadoop/app/hive-1.1.0-cdh5.6.0/lib/mysql-connector-java-5.1.44-bin.jar
```

```
have its own datastore table.
17/11/15 00:33:59 INFO DataNucleus.Query: Reading in results for query 'org.datanucleus.store.rdbms.query.SQLQ
s closing
17/11/15 00:33:59 INFO metastore.MetaStoreDirectSql: Using direct SQL, underlying DB is MySQL
17/11/15 00:33:59 INFO metastore.ObjectStore: Initialized ObjectStore
17/11/15 00:33:59 INFO metastore.HiveMetaStore: Added admin role in metastore
17/11/15 00:33:59 INFO metastore.HiveMetaStore: Added public role in metastore
17/11/15 00:33:59 INFO metastore.HiveMetaStore: No user is added in admin role, since config is empty
17/11/15 00:33:59 INFO metastore.HiveMetaStore: 0: get_all_databases
17/11/15 00:33:59 INFO HiveMetaStore.audit: ugi=hadoop ip=unknown-ip-addr cmd=get_all_databases
17/11/15 00:34:00 INFO metastore.HiveMetaStore: 0: get_functions: db=default pat=*
17/11/15 00:34:00 INFO HiveMetaStore.audit: ugi=hadoop ip=unknown-ip-addr cmd=get_functions: db=default
17/11/15 00:34:00 INFO DataNucleus.Datastore: The class 'org.apache.hadoop.hive.metastore.model.MResourceUri'
s not have its own datastore table.
17/11/15 00:34:00 INFO session.SessionState: Created local directory: /tmp/5b1487fc-e731-4b4d-8346-0c3aad8e21c
17/11/15 00:34:00 INFO session.SessionState: Created HDFS directory: /tmp/hive/hadoop/5b1487fc-e731-4b4d-8346-
17/11/15 00:34:00 INFO session.SessionState: Created local directory: /tmp/hadoop/5b1487fc-e731-4b4d-8346-0c3a
17/11/15 00:34:00 INFO session.SessionState: Created HDFS directory: /tmp/hive/hadoop/5b1487fc-e731-4b4d-8346-
17/11/15 00:34:00 INFO repl.SparkILoop: Created sql context (with Hive support)..
SQL context available as sqlContext.
```

```
scala>
```

```
scala>
```

```
scala> |
```

测试读取hive中的表：

```
1 scala> sqlContext.sql("select * from test").show
```


17/11/15 00:34:54 INFO sc

```
+---+-----+
| id|name|
+---+-----+
|  1| abc|
+---+-----+
```

scala> |

3. 开发工具搭建

3.1 下载IDE开发工具idea

从官网下载，如果有教育网邮箱，可以申请免费注册使用。

软件：

软件	版本
IDEA	ideaIU-2017.1
Scala的IDEA插件	scala-intellij-bin-2017.1.20

将两个软件上传到主机

3.2 安装IDEA，配置Scala插件

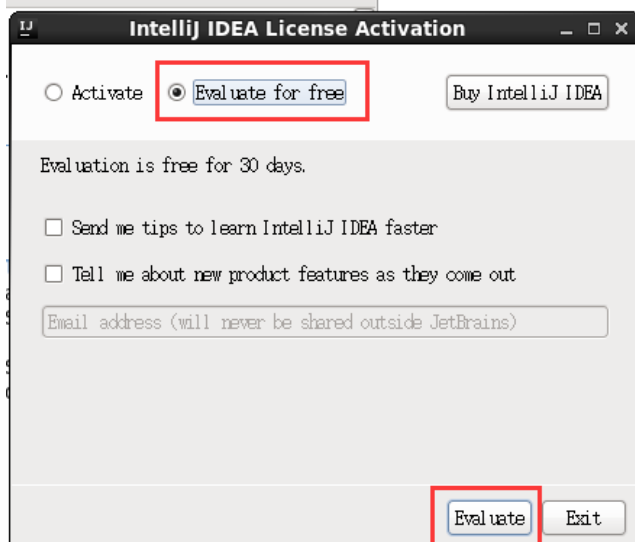
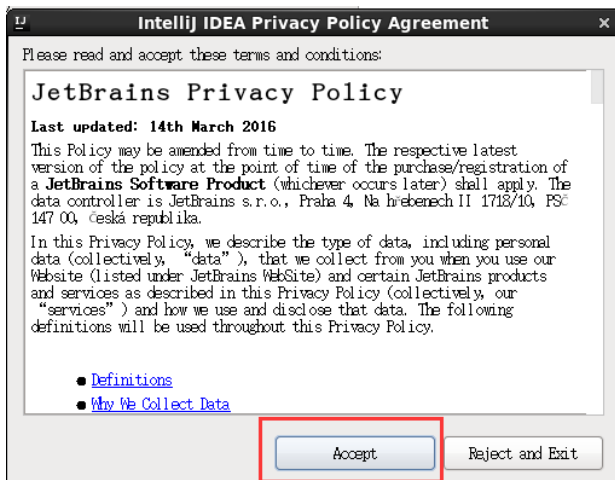
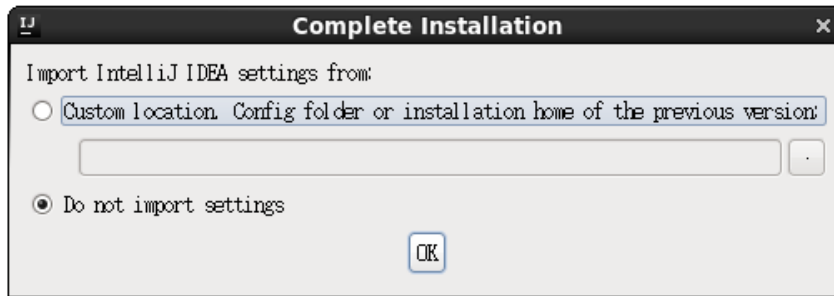
(1) 启动idea

```
[hadoop@spark1234 app]$ 
[hadoop@spark1234 app]$ tar -xzf ideaIU-2017.1.tar.gz
[hadoop@spark1234 app]$ 
[hadoop@spark1234 app]$ ls
hadoop-2.6.0-cdh5.6.0  idea-IU-171.3780.107  spark-1.6.3-bin-2.6.0-cdh5.6.0
hive-1.1.0-cdh5.6.0   ideaIU-2017.1.tar.gz
[hadoop@spark1234 app]$ 
[hadoop@spark1234 app]$ cd idea-IU-171.3780.107/
[hadoop@spark1234 idea-IU-171.3780.107]$ ls
bin  build.txt  help  Install-Linux-tar.txt  ire  lib  license  plugins  redistrib
[hadoop@spark1234 idea-IU-171.3780.107]$ ./bin/idea.sh
```

(2) 配置idea向导

6.0

redist



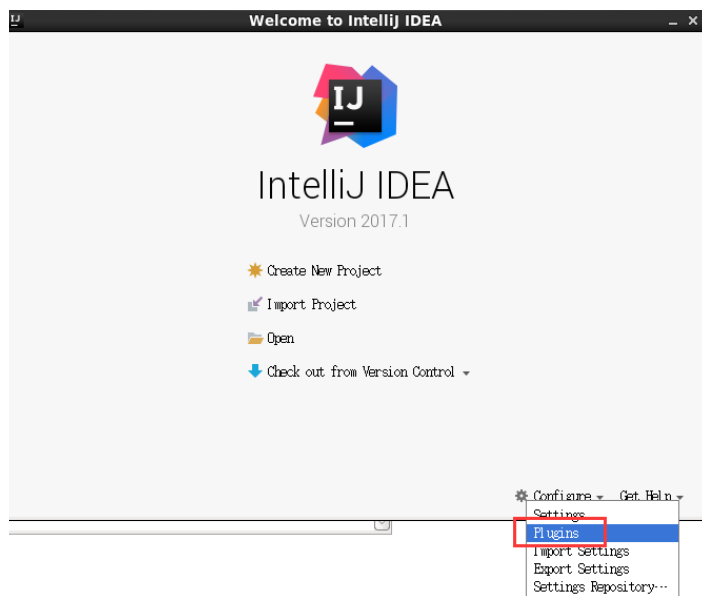
spark-1.6.3-bin-2.6.0-cdh5.6.0



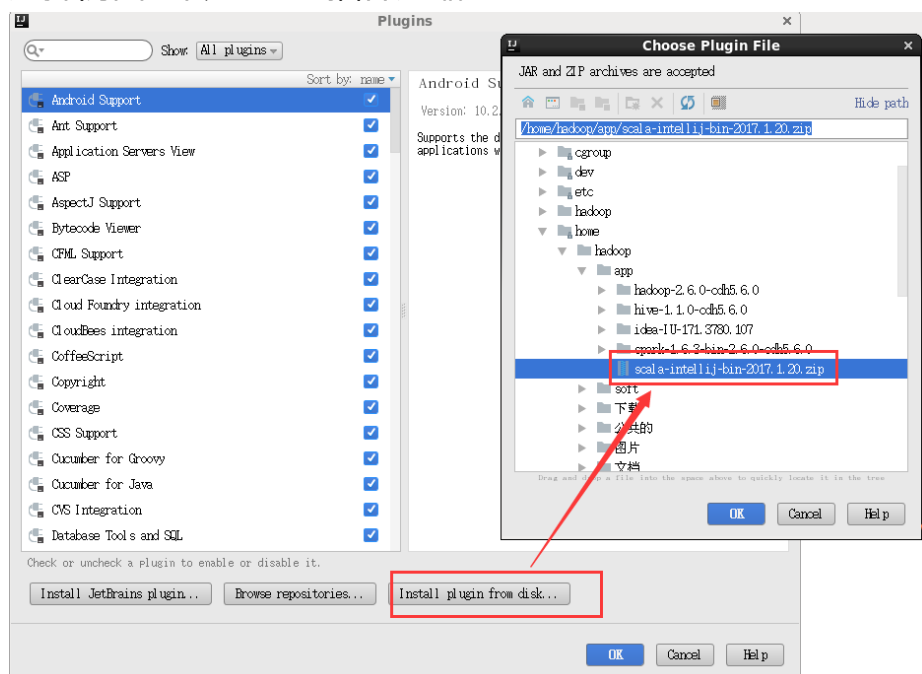


(3) 配置scala插件

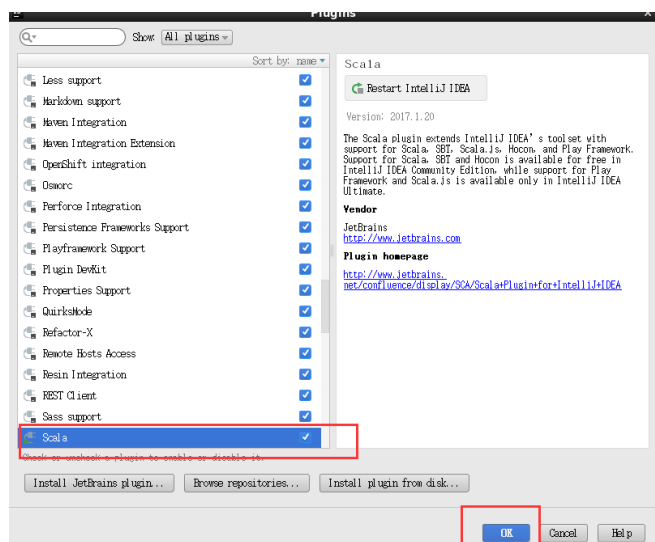
选择Plugins：



选择前面上传的scala插件压缩包：



点击OK：



重启：

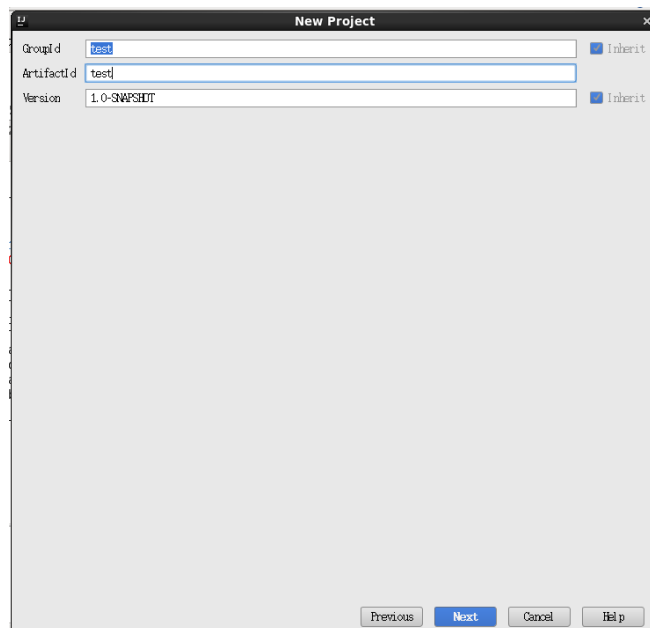
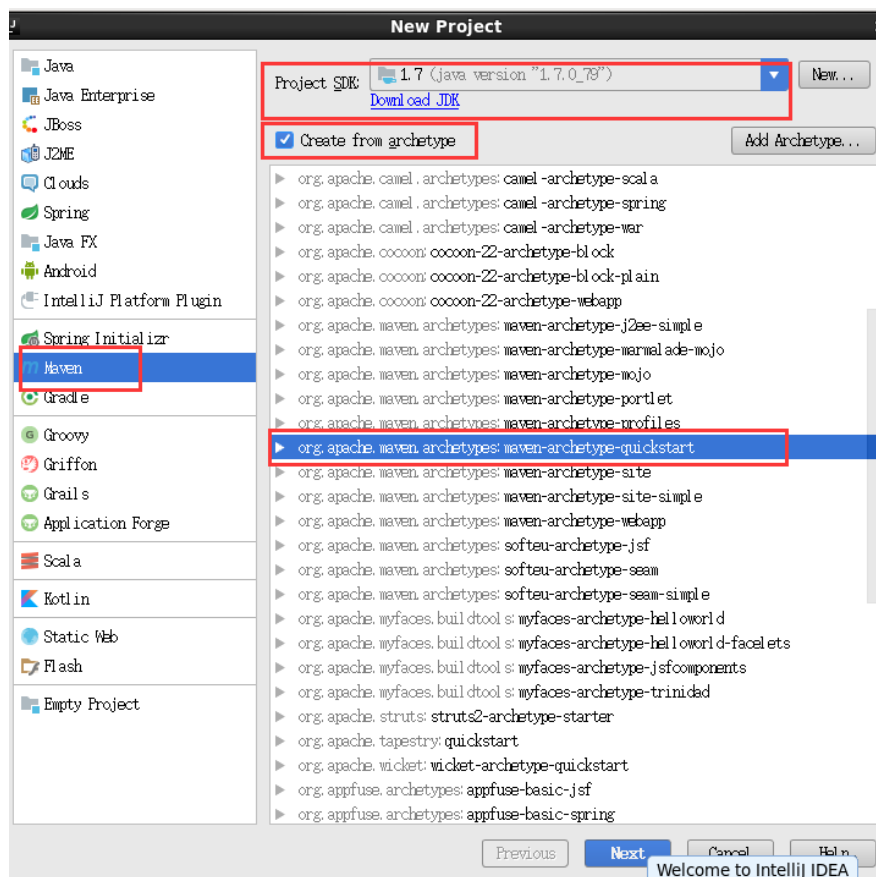


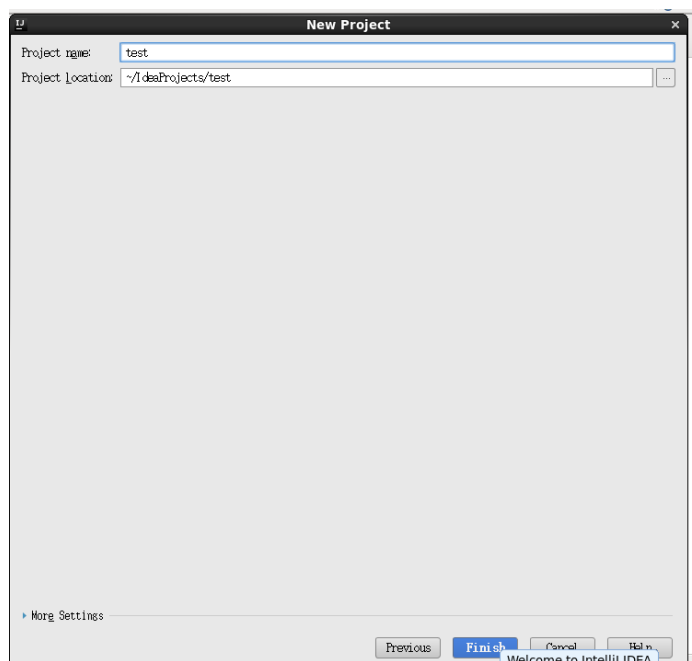
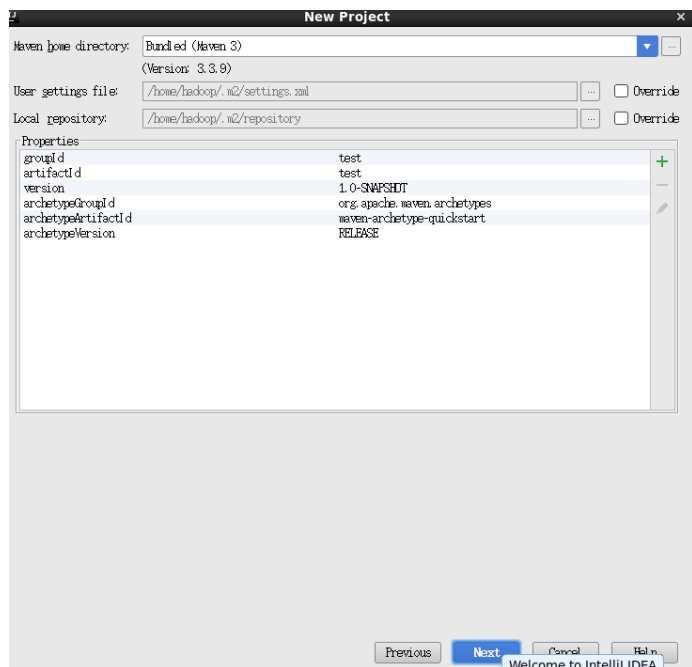
3.3 创建、配置测试的Spark项目

点击第一个，创建测试项目



创建Maven工程，第一次创建功能需要选择JDK的安装路径：





(5) 修改pom文件，增加项目需要的依赖

```
1 <project xmlns="http://maven.apache.org/POM/4.0.0"
  xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
2       xsi:schemaLocation="http://maven.apache.org/POM/4.0.0
  http://maven.apache.org/xsd/maven-4.0.0.xsd">
3   <modelVersion>4.0.0</modelVersion>
4
5   <groupId>Test</groupId>
6   <artifactId>Test</artifactId>
7   <version>1.0-SNAPSHOT</version>
```

```
8     <packaging>jar</packaging>
9
10    <name>JSDPI</name>
11    <url>http://maven.apache.org</url>
12    <repositories>
13        <repository>
14            <id>cloudera</id>
15            <url>https://repository.cloudera.com/artifactory/cloudera-
repos/</url>
16        </repository>
17    </repositories>
18
19    <properties>
20        <project.build.sourceEncoding>UTF-8</project.build.sourceEncoding>
21        <spark.version>1.6.0-cdh5.7.0</spark.version>
22        <jedis.version>2.8.2</jedis.version>
23        <hadoop.version>2.6.0-cdh5.7.0</hadoop.version>
24        <fastjson.version>1.2.14</fastjson.version>
25        <jetty.version>9.2.5.v20141112</jetty.version>
26        <container.version>2.17</container.version>
27        <java.version>1.8</java.version>
28        <scala.version>2.10.6</scala.version>
29    </properties>
30
31
32
33    <dependencies>
34        <dependency>
35            <groupId>redis.clients</groupId>
36            <artifactId>jedis</artifactId>
37            <version>${jedis.version}</version>
38        </dependency>
39        <dependency>
40            <groupId>org.apache.hadoop</groupId>
41            <artifactId>hadoop-common</artifactId>
42            <version>${hadoop.version}</version>
43            <exclusions>
44                <exclusion>
45                    <groupId>javax.servlet</groupId>
46                    <artifactId>*</artifactId>
47                </exclusion>
48            </exclusions>
49        </dependency>
50        <dependency>
51            <groupId>org.apache.hadoop</groupId>
52            <artifactId>hadoop-hdfs</artifactId>
```



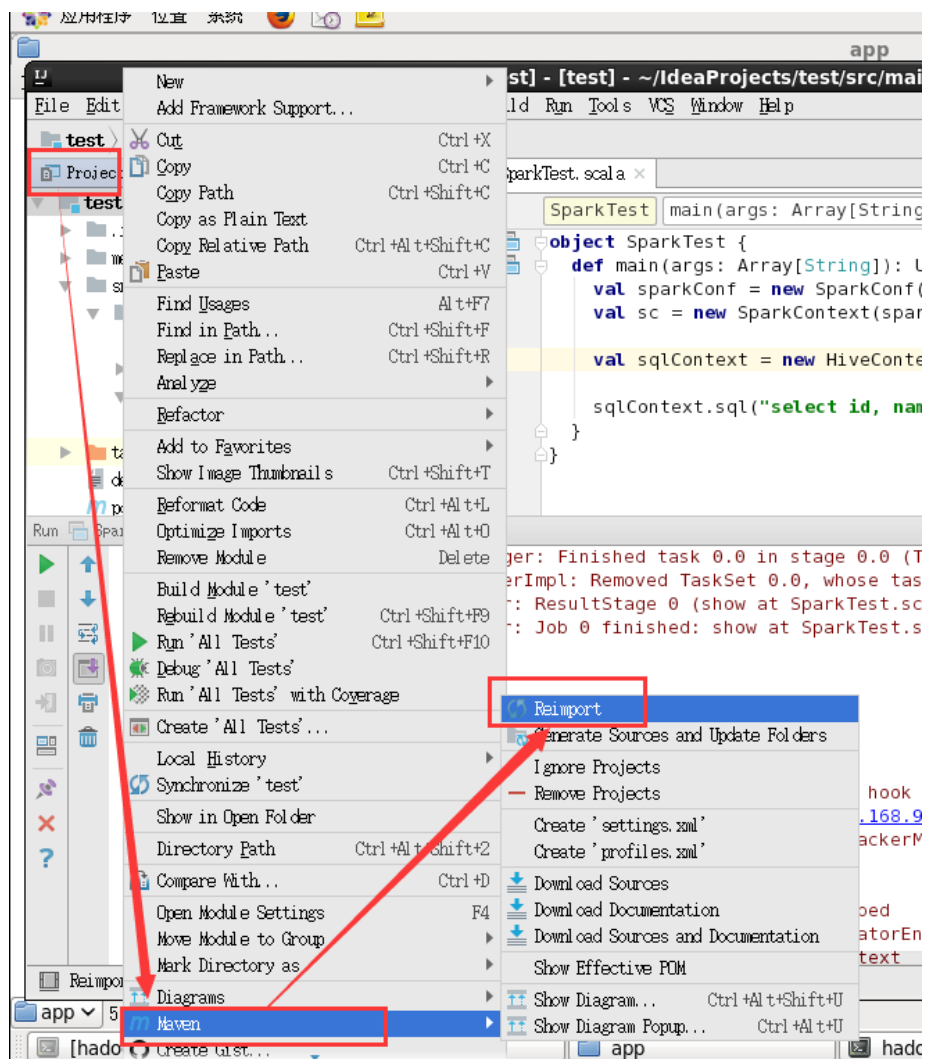
```
53         <version>${hadoop.version}</version>
54     <exclusions>
55         <exclusion>
56             <groupId>javax.servlet</groupId>
57             <artifactId>*</artifactId>
58         </exclusion>
59     </exclusions>
60 </dependency>
61 <dependency>
62     <groupId>org.apache.hadoop</groupId>
63     <artifactId>hadoop-client</artifactId>
64     <version>${hadoop.version}</version>
65     <exclusions>
66         <exclusion>
67             <groupId>javax.servlet</groupId>
68             <artifactId>*</artifactId>
69         </exclusion>
70     </exclusions>
71 </dependency>
72 <dependency>
73     <groupId>org.apache.spark</groupId>
74     <artifactId>spark-core_2.10</artifactId>
75     <version>${spark.version}</version>
76 </dependency>
77 <dependency>
78     <groupId>org.apache.spark</groupId>
79     <artifactId>spark-sql_2.10</artifactId>
80     <version>${spark.version}</version>
81 </dependency>
82 <dependency>
83     <groupId>org.apache.spark</groupId>
84     <artifactId>spark-streaming_2.10</artifactId>
85     <version>${spark.version}</version>
86 </dependency>
87 <dependency>
88     <groupId>org.apache.spark</groupId>
89     <artifactId>spark-yarn_2.10</artifactId>
90     <version>${spark.version}</version>
91 </dependency>
92 <dependency>
93     <groupId>org.apache.spark</groupId>
94     <artifactId>spark-hive_2.10</artifactId>
95     <version>${spark.version}</version>
96 </dependency>
97 <dependency>
98     <groupId>com.google.guava</groupId>
```

```
99         <artifactId>guava</artifactId>
100         <version>18.0</version>
101     </dependency>
102     <dependency>
103         <groupId>com.alibaba</groupId>
104         <artifactId>fastjson</artifactId>
105         <version>${fastjson.version}</version>
106     </dependency>
107     <dependency>
108         <groupId>junit</groupId>
109         <artifactId>junit</artifactId>
110         <version>3.8.1</version>
111         <scope>test</scope>
112     </dependency>
113     <dependency>
114         <groupId>org.spark-project.hive</groupId>
115         <artifactId>hive-jdbc</artifactId>
116         <version>0.12.0</version>
117     </dependency>
118     <dependency>
119         <groupId>log4j</groupId>
120         <artifactId>log4j</artifactId>
121         <version>1.2.14</version>
122     </dependency>
123     <dependency>
124         <groupId>com.fasterxml.jackson.core</groupId>
125         <artifactId>jackson-core</artifactId>
126         <version>2.5.3</version>
127     </dependency>
128     <dependency>
129         <groupId>com.fasterxml.jackson.core</groupId>
130         <artifactId>jackson-annotations</artifactId>
131         <version>2.5.3</version>
132     </dependency>
133     <dependency>
134         <groupId>org.codehaus.jackson</groupId>
135         <artifactId>jackson-mapper-asl</artifactId>
136         <version>1.9.0</version>
137     </dependency>
138     <dependency>
139         <groupId>org.apache.spark</groupId>
140         <artifactId>spark-streaming-kafka_2.10</artifactId>
141         <version>${spark.version}</version>
142     </dependency>
143     <dependency>
144         <groupId>mysql</groupId>
```

```
145         <artifactId>mysql-connector-java</artifactId>
146         <version>5.1.34</version>
147     </dependency>
148 </dependencies>
149 <build>
150     <plugins>
151         <plugin>
152             <artifactId>maven-assembly-plugin</artifactId>
153             <version>2.3</version>
154             <configuration>
155                 <classifier>dist</classifier>
156                 <appendAssemblyId>true</appendAssemblyId>
157                 <descriptorRefs>
158                     <descriptor>jar-with-dependencies</descriptor>
159                 </descriptorRefs>
160             </configuration>
161             <executions>
162                 <execution>
163                     <id>make-assembly</id>
164                     <phase>package</phase>
165                     <goals>
166                         <goal>single</goal>
167                     </goals>
168                 </execution>
169             </executions>
170         </plugin>
171
172         <plugin>
173             <artifactId>maven-compiler-plugin</artifactId>
174             <configuration>
175                 <source>1.7</source>
176                 <target>1.7</target>
177             </configuration>
178         </plugin>
179
180         <plugin>
181             <groupId>net.alchim31.maven</groupId>
182             <artifactId>scala-maven-plugin</artifactId>
183             <version>3.2.2</version>
184             <executions>
185                 <execution>
186                     <id>scala-compile-first</id>
187                     <phase>process-resources</phase>
188                     <goals>
189                         <goal>compile</goal>
190                     </goals>
```

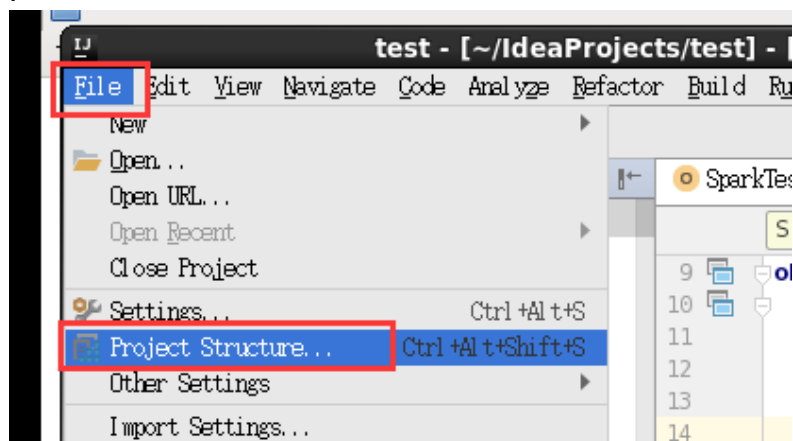
```
191         </execution>
192     </executions>
193     <configuration>
194         <scalaVersion>${scala.version}</scalaVersion>
195         <recompileMode>incremental</recompileMode>
196         <useZincServer>true</useZincServer>
197         <args>
198             <arg>-unchecked</arg>
199             <arg>-deprecation</arg>
200             <arg>-feature</arg>
201         </args>
202         <jvmArgs>
203             <jvmArg>-Xms1024m</jvmArg>
204             <jvmArg>-Xmx1024m</jvmArg>
205         </jvmArgs>
206         <javacArgs>
207             <javacArg>-source</javacArg>
208             <javacArg>${java.version}</javacArg>
209             <javacArg>-target</javacArg>
210             <javacArg>${java.version}</javacArg>
211             <javacArg>-Xlint:all,-serial,-path</javacArg>
212         </javacArgs>
213     </configuration>
214 </plugin>
215
216 </plugins>
217 </build>
218 </project>
219
```

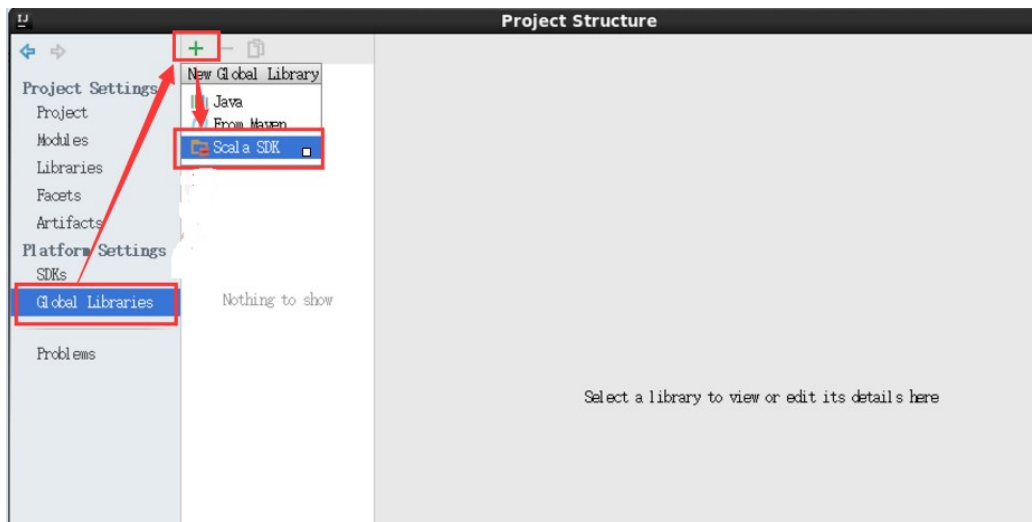
在项目上右键，选择Maven-》Reimport, idea会自动从资源库下载需要的依赖jar包：整个过程可能需要很久（几个小时或者更久）。



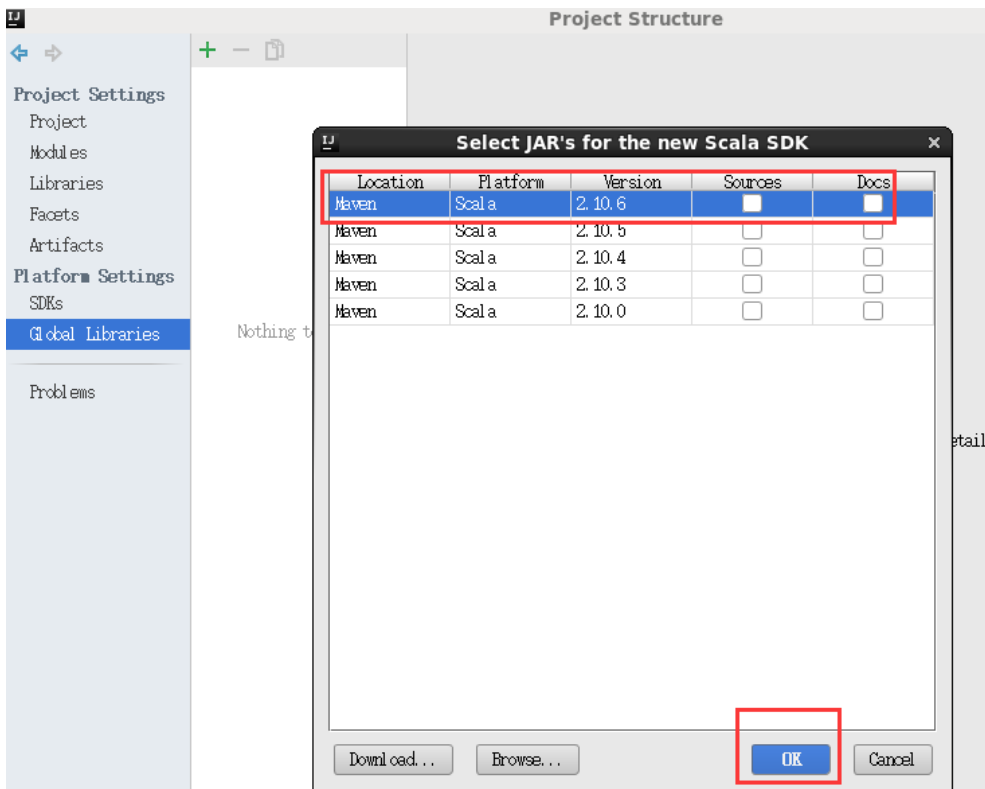
（6）增加项目的Scala支持

pom文件的依赖包下载完毕后，增加项目的Scala支持。





选择Scala 2.10.6：

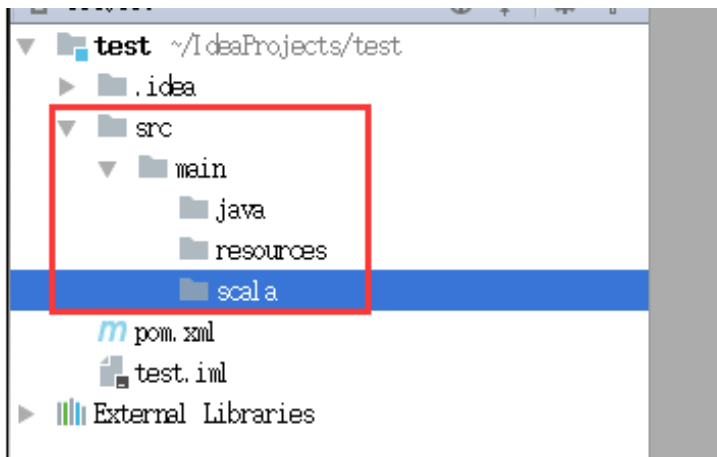


(7) 创建项目目录结构

在项目下面，创建如下层次的目录：

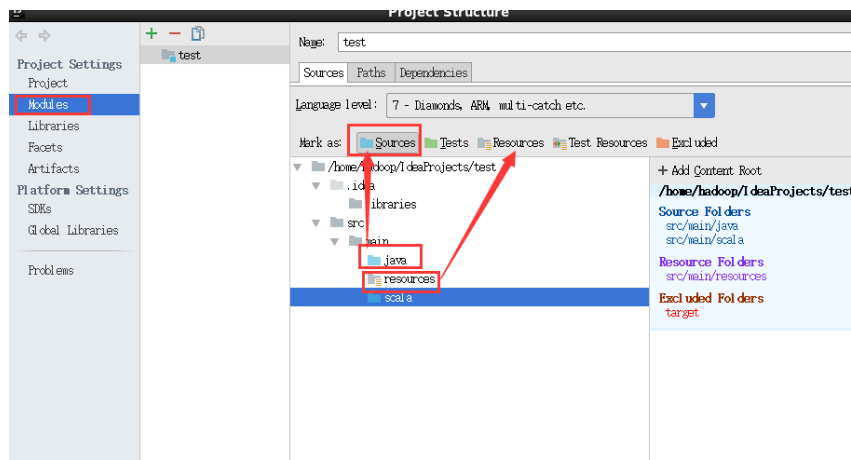
```

-----test/src/main
-----| java
-----| resources
-----| scala
  
```



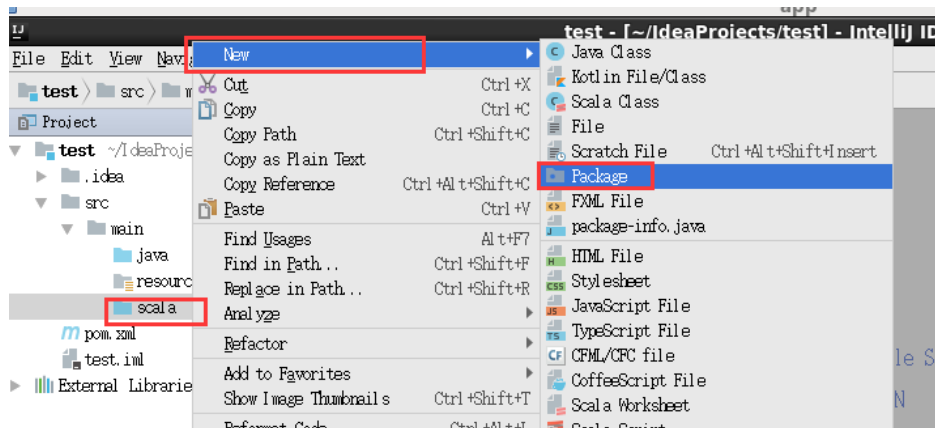
设置每个目录的模块功能：

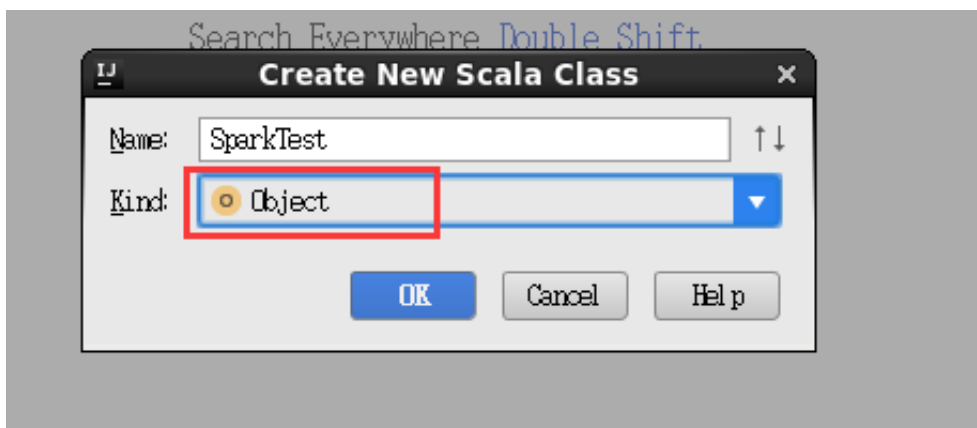
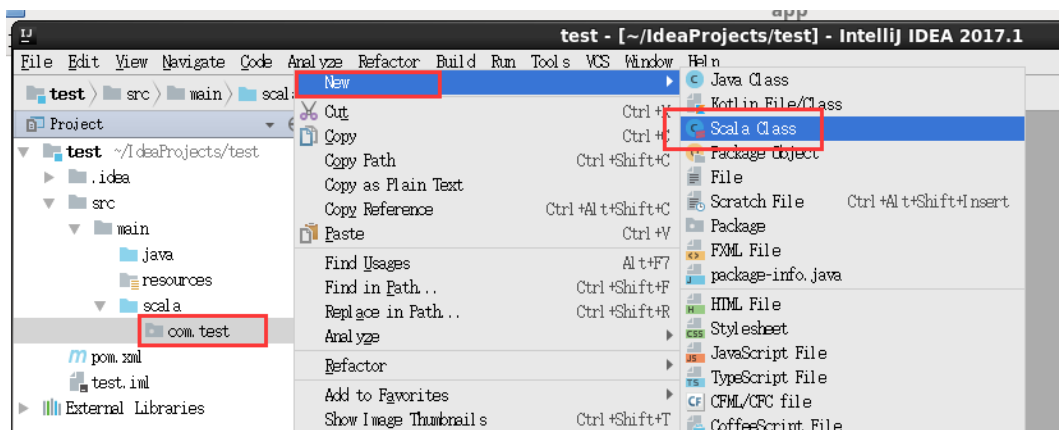
选择java， 设置为Sources。 选择resources目录， 设置为Resources。 如下：



（8）创建scala测试程序

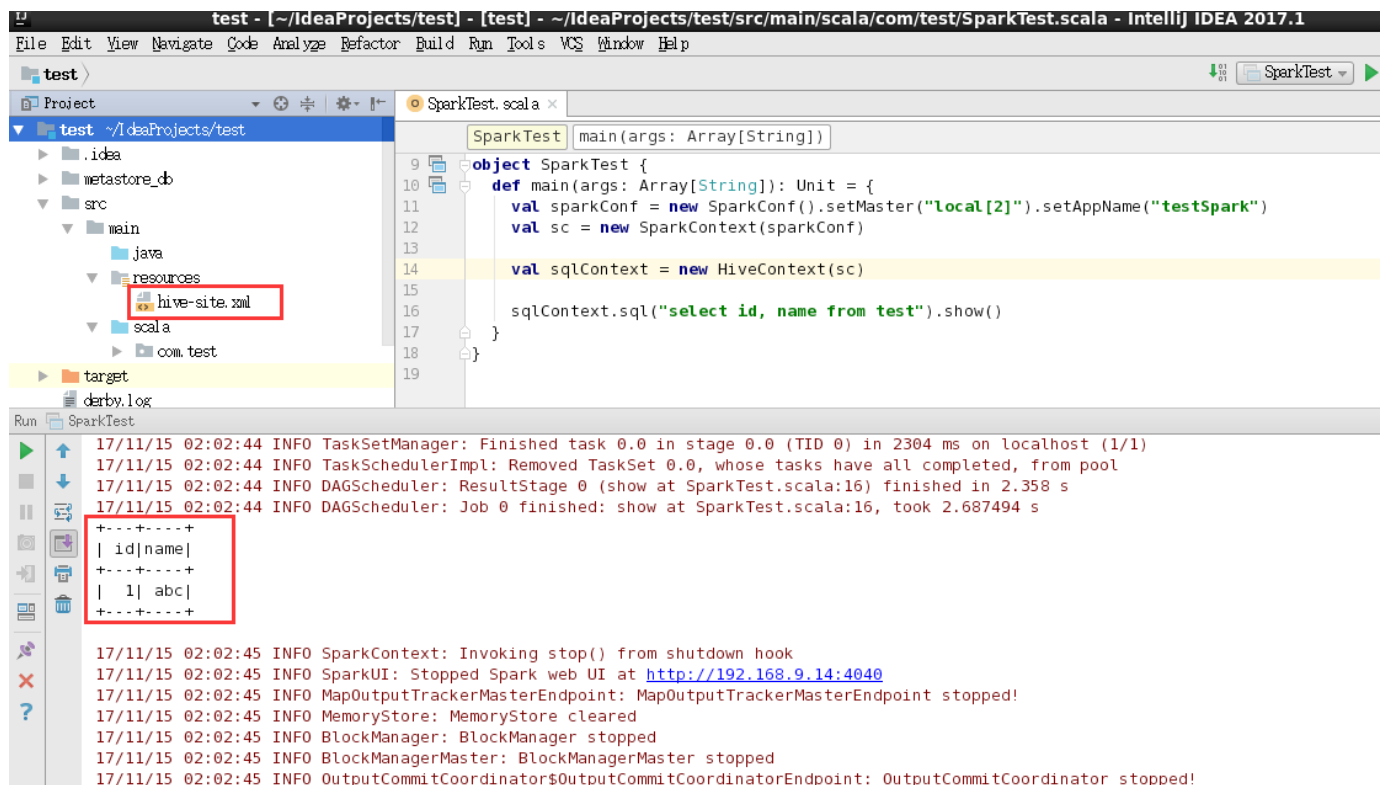
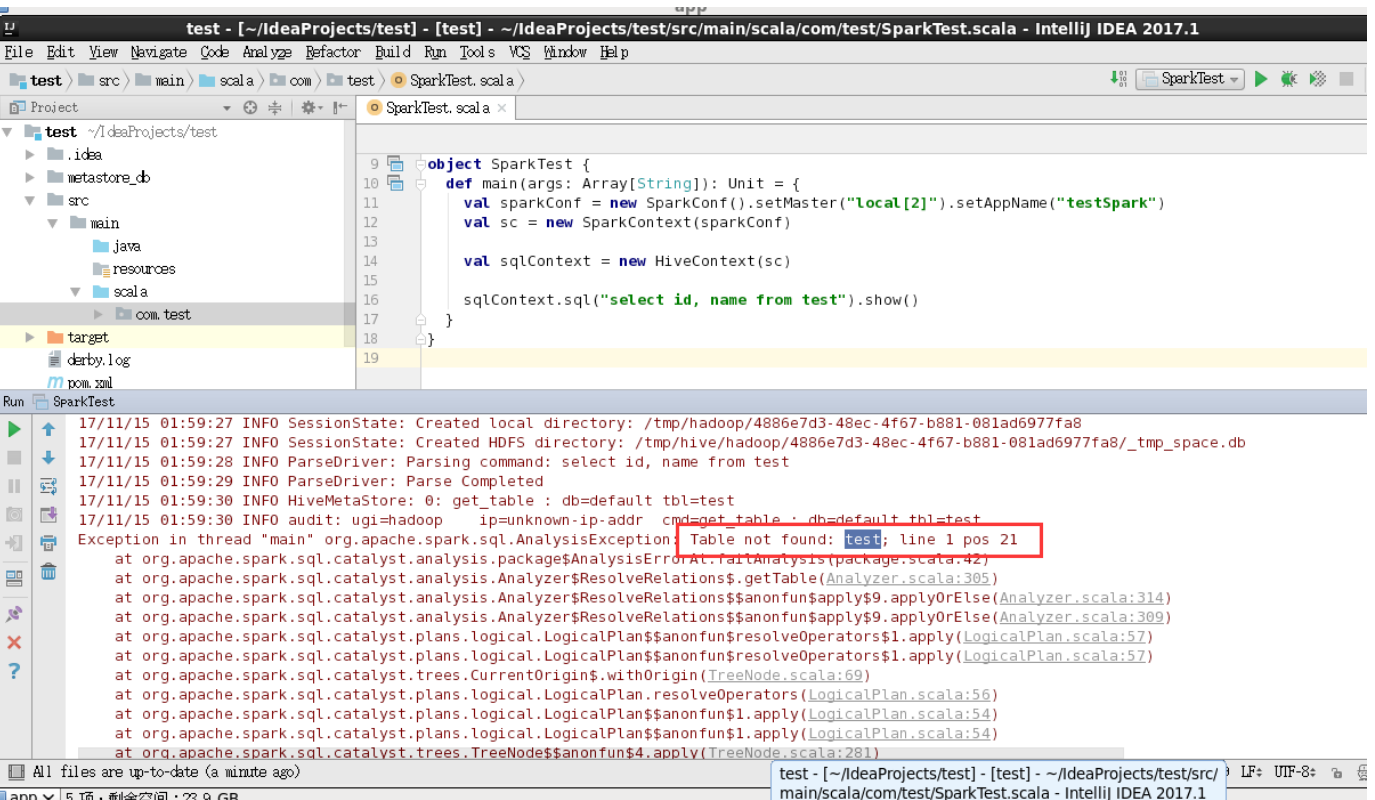
先创建包：





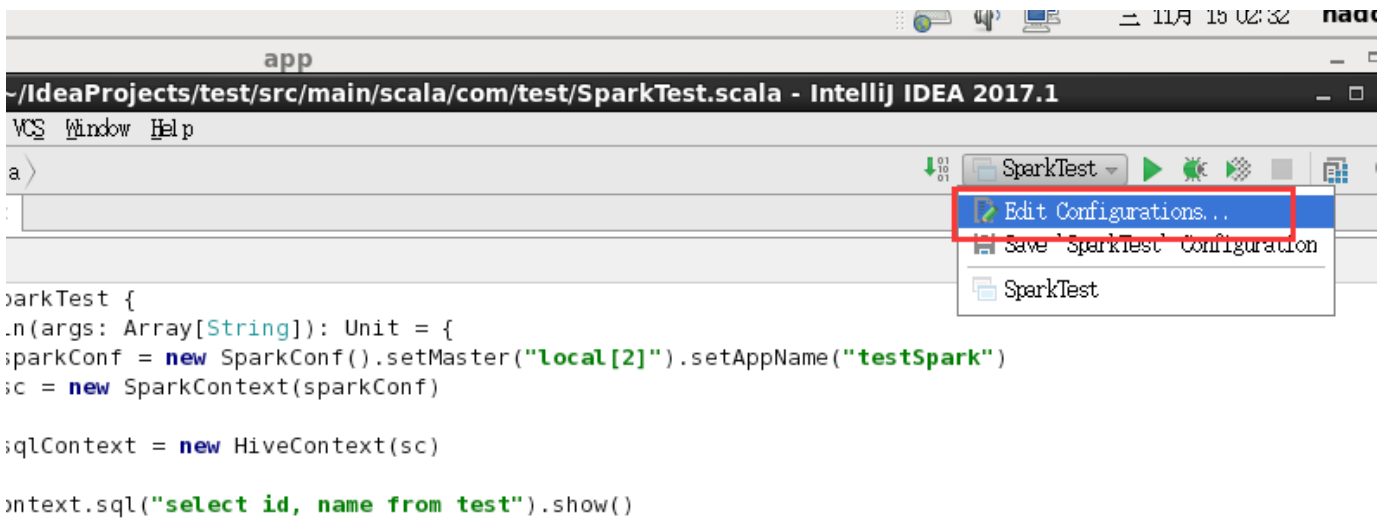
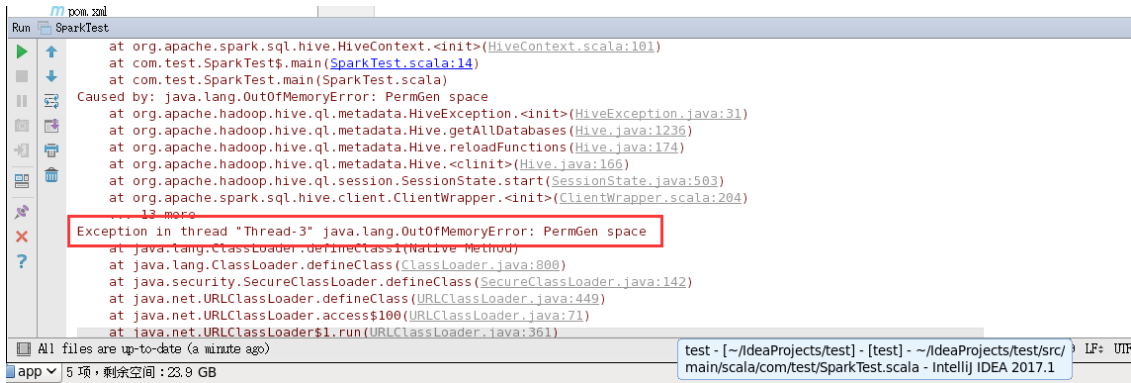
测试代码：

```
1 package com.test
2
3 import org.apache.spark.sql.hive.HiveContext
4 import org.apache.spark.{SparkConf, SparkContext}
5
6 /**
7  * Created by hadoop on 17-11-15.
8  */
9 object SparkTest {
10   def main(args: Array[String]): Unit = {
11     val sparkConf = new SparkConf().setMaster("local[2]").setAppName("testSpark")
12     val sc = new SparkContext(sparkConf)
13
14     val sqlContext = new HiveContext(sc)
15
16     sqlContext.sql("select id, name from test").show()
17   }
18 }
19
```

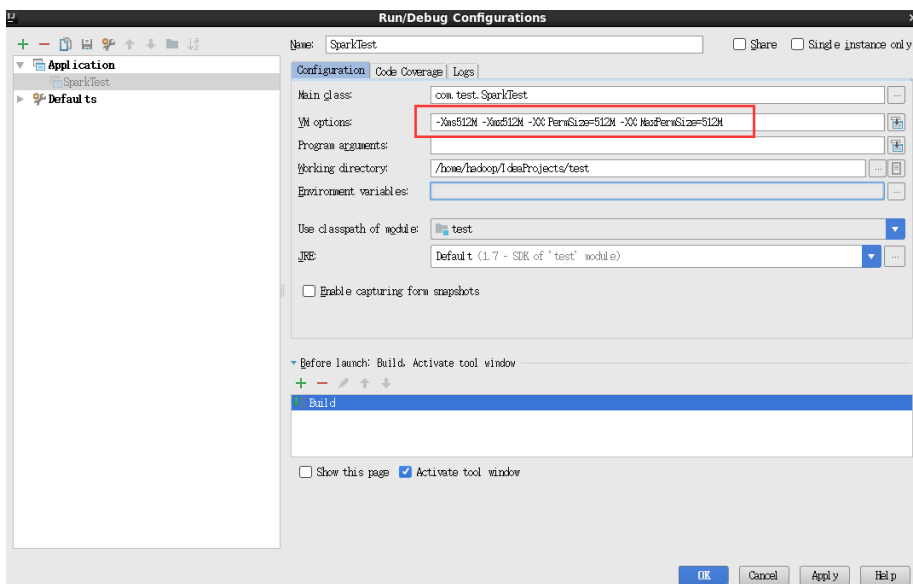



如果在idea操作hdfs，需要将hadoop的core-site.xml复制到resources目录。

(9) 如果报内存溢出错误，需要修改下IDEA的配置：



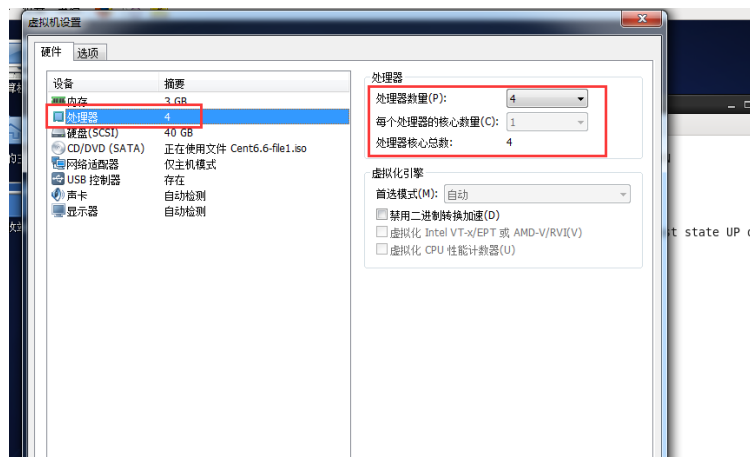
在VM options设置内存：-Xms512M -Xmx512M -XX:PermSize=512M -
XX:MaxPermSize=512M



4. 其他

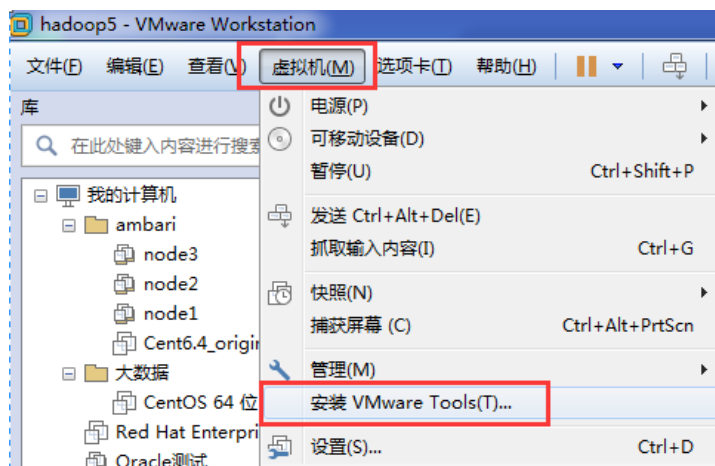
4.1 设置较大的虚拟机内存和CPU

内存和CPU尽量多分点，内存过小会导致spark作业提交失败。
CPU核心数太小会导致两个线程提交spark作业无法运行。



4.2 虚拟机和物理机无缝相互复制/粘贴

(1) 选择安装VMware Tools



(2) 挂载光盘

```
[root@spark1234 ~]#  
[root@spark1234 ~]# mount /dev/cdrom /mnt
```

(3) 进入挂载目录，解压到/root/目录

```
[root@spark1234 /]# cd /mnt  
[root@spark1234 mnt]#  
[root@spark1234 mnt]# tar -xzf VMwareTools-10.0.6-3595377.tar.gz /root/
```

(4) 进入vmware-tools的解压目录，一路回车即可：

```
[root@spark1234 ~]#  
[root@spark1234 ~]# cd vmware-tools-distrib/  
[root@spark1234 vmware-tools-distrib]#  
[root@spark1234 vmware-tools-distrib]# ls  
bin doc FILES installer vgauth vmware-install.real  
caf etc INSTALL lib vmware-install.pl  
[root@spark1234 vmware-tools-distrib]#  
[root@spark1234 vmware-tools-distrib]# perl vmware-install.pl
```

安装完成后，重启虚拟机，物理机和虚拟机可以无缝复制粘贴。