# 搭建Hadoop+Spark分布式集群

## 实验目的

搭建真实的分布式计算环境。
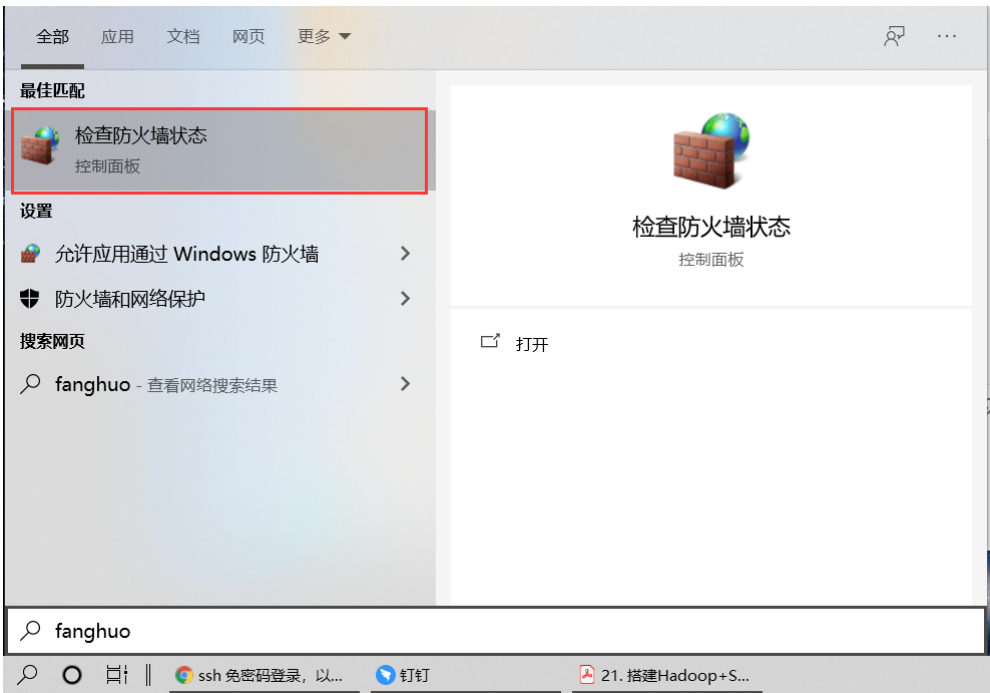
## 实验内容
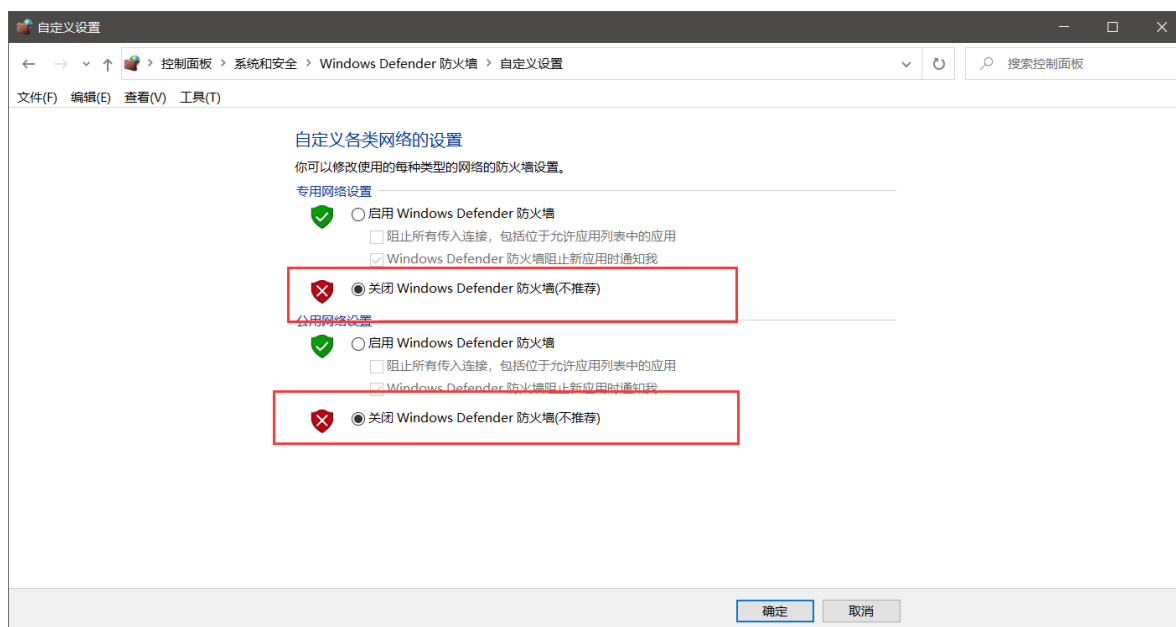
### 1.准备机器



### 2.配置桥接网卡

**以master为例说明网络配置：**

> 1.首先关闭宿主机（windows10）的防火墙

## 2.在虚拟机软件上添加一张桥接网卡

虚拟机设置

| 硬件 | 选项 |

设备状态

设备 | 摘要
内存 | 4 GB
处理器 | 2
硬盘 (SCSI) | 30 GB
CD/DVD (SATA) | 自动检测
网络适配器 2 | 桥接模式 (自动)
USB 控制器 | 存在
声卡 | 自动检测
打印机 | 存在
显示器 | 自动检测

☐ 已连接(C)
☑ 启动时连接(O)

网络连接
◉ 桥接模式(B): 直接连接物理网络
　☐ 复制物理网络连接状态(P)
○ NAT 模式(N): 用于共享主机的 IP 地址
○ 仅主机模式(H): 与主机共享的专用网络
○ 自定义(U): 特定虚拟网络
　VMnet0
○ LAN 区段(L):

LAN 区段(S)...　高级(V)...

添加(A)...　移除(R)

确定　取消　帮助

2.编辑网卡设置

文件(F)　编辑(E)　查看(V)　虚拟机(M)　选项卡(T)　帮助(H)

库　×

在此处键入内容...

我的计算机
　kali-xuegod53
　chen
　controller
　compute1
共享的虚拟机

主页 ×　chen ×

chen

▶ 开启此虚拟机
编辑虚拟机设置

▼ 设备
内存　4 GB
处理器

| 名称 | 类型 | 外部连接 | 主机连接 | DHCP | 子网地址 |
|------|------|----------|----------|------|----------|
| VMnet0 | 自定义... | - | - | - | 192.168.58.0 |
| VMnet1 | 仅主机... | - | 已连接 | 已启用 | 192.168.44.0 |
| VMnet8 | NAT 模式 | NAT 模式 | 已连接 | 已启用 | 10.0.0.0 |

添加网络(E)... 　移除网络(O) 　重命名网络(W)...

VMnet 信息

○ 桥接模式(将虚拟机直接连接到外部网络)(B)

已桥接至(G): ⌄ 　自动设置(U)...

○ NAT 模式(与虚拟机共享主机的 IP 地址)(N) 　NAT 设置(S)...

⦿ 仅主机模式(在专用网络内连接虚拟机)(H)

☐ 将主机虚拟适配器连接到此网络(V)

主机虚拟适配器名称: VMware 网络适配器 VMnet0

☐ 使用本地 DHCP 服务将 IP 地址分配给虚拟机(D) 　DHCP 设置(P)...

子网 IP (I): 192 . 168 . 58 . 0 　　子网掩码(M): 255 . 255 . 255 . 0

⚠ 需要具备管理员特权才能修改网络配置。 🛡更改设置(C)
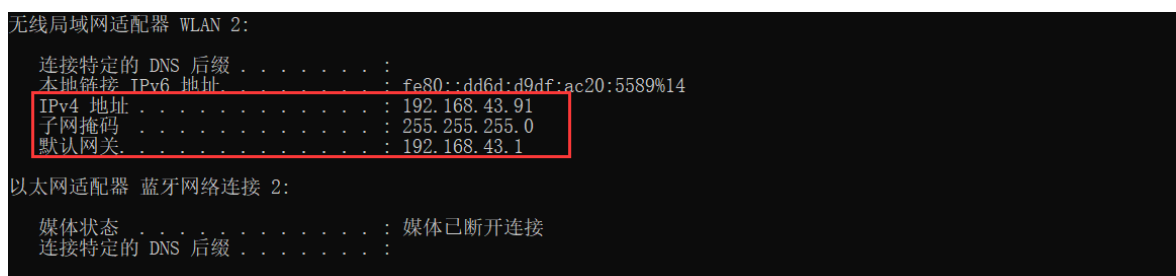
还原默认设置(R) 　导入(T)... 　导出(X)... 　确定 　取消 　应用(A) 　帮助

| 名称 | 类型 | 外部连接 | 主机连接 | DHCP | 子网地址 |
|------|------|---------|---------|------|---------|
| VMnet0 | 桥接模式 | Intel(R) Wireless-AC 9462 #2 | - | - | - |
| VMnet1 | 仅主机... | - | 已连接 | 已启用 | 192.168.44.0 |
| VMnet8 | NAT 模式 | NAT 模式 | 已连接 | 已启用 | 10.0.0.0 |

添加网络(E)...　移除网络(O)　重命名网络(W)...

VMnet 信息

◉ 桥接模式(将虚拟机直接连接到外部网络)(B)

　已桥接至(G): Intel(R) Wireless-AC 9462 #2 ▽　自动设置(U)...

◯ NAT 模式(与虚拟机共享主机的 IP 地址)(N)　NAT 设置(S)...

◯ 仅主机模式(在专用网络内连接虚拟机)(H)

☐ 将主机虚拟适配器连接到此网络(V)

　主机虚拟适配器名称: VMware 网络适配器 VMnet0

☐ 使用本地 DHCP 服务将 IP 地址分配给虚拟机(D)　DHCP 设置(P)...

子网 IP (I): [　.　.　.　]　子网掩码(M): [　.　.　.　]

还原默认设置(R)　导入(T)...　导出(X)...　确定　取消　应用(A)　帮助

3.开启桥接网卡，同时关闭第一张网卡，也可以不关闭



### chen

▶ 开启此虚拟机
　编辑虚拟机设置

▼ 设备

| 内存 | 4 GB |
|------|------|
| 处理器 | 2 |
| 硬盘 (SCSI) | 30 GB |
| CD/DVD (SATA) | 自动检测 |
| 网络适配器 2 | 桥接模式 (自动) |
| USB 控制器 | 存在 |
| 声卡 | 自动检测 |
| 打印机 | 存在 |
| 显示器 | 自动检测 |

桥接模式 (自动)

虚拟机设置

硬件 | 选项

| 设备 | 摘要 |
|---|---|
| 内存 | 4 GB |
| 处理器 | 2 |
| 硬盘 (SCSI) | 30 GB |
| CD/DVD (SATA) | 自动检测 |
| 网络适配器 2 | 桥接模式 (自动) |
| USB 控制器 | 存在 |
| 声卡 | 自动检测 |
| 打印机 | 存在 |
| 显示器 | 自动检测 |

设备状态
☐ 已连接(C)
☑ 启动时连接(O)

网络连接
◉ 桥接模式(B): 直接连接物理网络
☐ 复制物理网络连接状态(P)
◯ NAT 模式(N): 用于共享主机的 IP 地址
◯ 仅主机模式(H): 与主机共享的专用网络
◯ 自定义(U): 特定虚拟网络
VMnet0
◯ LAN 区段(L):

LAN 区段(S)... | 高级(V)...

添加(A)... | 移除(R)

确定 | 取消 | 帮助

4.然后令宿主机（windows10）连接到手机热点。打开cmd，使用命令【ipconfig】查看ip地址



无线局域网适配器 WLAN 2:

连接特定的 DNS 后缀 . . . . . . . . :
本地链接 IPv6 地址. . . . . . . . . . : fe80::dd6d:d9df:ac20:5589%14
IPv4 地址 . . . . . . . . . . . . : 192.168.43.91
子网掩码 . . . . . . . . . . . . : 255.255.255.0
默认网关. . . . . . . . . . . . . : 192.168.43.1

以太网适配器 蓝牙网络连接 2:

媒体状态 . . . . . . . . . . . . : 媒体已断开连接
连接特定的 DNS 后缀 . . . . . . . . :

再开启虚拟机，编辑配置文件【/etc/network/interfaces】

添加以下内容:

```
auto ens38
iface ens38 inet static
# IP地址和宿主机保持在同一网段，这里指定92
address 192.168.43.92
# 和宿主机一致
netmask 255.255.255.0
# 和宿主机一致
gateway 192.168.43.1
```



保存退出重启【reboot】

其他的机器和master配置方法相同。配置文件写入相对自己宿主机的内容

验证互相ping，可以通。

## 3.设置免密登录

每台机器上生成密钥

```
ssh-keygen -t rsa
# 选择覆盖原来的文件，其他都直接回车
# 私钥id_rsa 公钥 id_rsa.pub
```

将每台机器的公钥（id_rsa.pub里的内容）都放到authorized_keys文件里，authorized_keys中有所有机器的公钥，所有机器都有authorized_keys这个文件。

因为不同的机器用户名不一样，所以如果想通过【ssh worker1】实现连接，还需要修改（创建）【~/.ssh/config】写入如下内容：然后将其复制到所有机器。

```
Host master
user chen
Host worker1
user huang
Host worker2
user chen
Host worker3
user guo
```



分发给其他机器

```
scp ~/.ssh/config chen@worker2:~/.ssh/
scp ~/.ssh/config huang@worker1:~/.ssh/
scp ~/.ssh/config guo@worker3:~/.ssh/
```

测试

# 4.配置hadoop

首先将Hadoop改为集群模式。在master主机中修改下面四个文件。

> 1.修改 core-site.xml

将配置文件 /apps/hadoop/etc/hadoop/core-site.xml中fs.defaultFS的值由hdfs://localhost:9000改为hdfs://master:9000，修改以后，如下图所示



> 2.修改hdfs-site.xml文件

将配置文件/apps/hadoop/etc/hadoop/hdfs-site.xml中dfs.replication的值由1改为4,修改以后如下图所示

3.修改workers文件

将配置文件/apps/hadoop/etc/hadoop/workers中localhost改为

```
master
worker1
worker2
worker3
```

修改后如下图所示



4.修改yarn-site.xml文件

将以下内容添加到配置文件/apps/hadoop/etc/hadoop/ yarn-site.xml中

```
<property>
<name>yarn.resourcemanager.hostname</name>
<value>master</value>
</property>
```

修改以后的文件内容如下图所示



> 5.将上面修改的四个文件复制到worker1和worker2两个节点，覆盖原来的文件

```
cd /apps/hadoop/etc/hadoop/
scp core-site.xml hdfs-site.xml workers yarn-site.xml \
worker1:/apps/hadoop/etc/hadoop/
scp core-site.xml hdfs-site.xml workers yarn-site.xml \
worker2:/apps/hadoop/etc/hadoop/
scp core-site.xml hdfs-site.xml workers yarn-site.xml \
worker3:/apps/hadoop/etc/hadoop/
```

## 删除伪分布式namenode文件

重新对分布式文件系统进行格式化前，需要删除三台主机中/data/tmp/hadoop/hdfs/目录下的文件和文件夹。首先删除master上/data/tmp/hadoop/hdfs/目录下的文件和文件夹

```
rm -rf /data/tmp/hadoop/hdfs/*
```

删除另外三台主机上相应的文件.

## 格式化分布式文件系统

在主节点master执行以下命令

```
hadoop namenode -format
```

```
2020-12-20 22:48:54,154 INFO util.GSet: 0.029999999329447746% max memory 869.5
MB = 267.1 KB
2020-12-20 22:48:54,154 INFO util.GSet: capacity     = 2^15 = 32768 entries
2020-12-20 22:48:54,193 INFO namenode.FSImage: Allocated new BlockPoolId: BP-10
31313234-192.168.43.92-1608533334179
2020-12-20 22:48:54,212 INFO common.Storage: Storage directory /data/tmp/hadoop
/hdfs/name has been successfully formatted.
2020-12-20 22:48:54,238 INFO namenode.FSImageFormatProtobuf: Saving image file
/data/tmp/hadoop/hdfs/name/current/fsimage.ckpt_0000000000000000000 using no co
mpression
2020-12-20 22:48:54,384 INFO namenode.FSImageFormatProtobuf: Image file /data/t
mp/hadoop/hdfs/name/current/fsimage.ckpt_0000000000000000000 of size 389 bytes
saved in 0 seconds.
2020-12-20 22:48:54,415 INFO namenode.NNStorageRetentionManager: Going to retai
n 1 images with txid >= 0
2020-12-20 22:48:54,432 INFO namenode.NameNode: SHUTDOWN_MSG:
/************************************************************
SHUTDOWN_MSG: Shutting down NameNode at master/192.168.43.92
************************************************************/
chen@master:~$
```

至此，Hadoop分布式集群就设置好了，下面进行测试。

# 5.测试hadoop

### 启动Hadoop

在master节点执行

```
/apps/hadoop/sbin/start-all.sh
```

### 查看Hadoop进程

查看主节点和其他节点的Hadoop进程



可以看到HDFS的NameNode和SecondaryNameNode，以及Yarn的ResourceManager只运行在主节点；HDFS的DataNode 和MapReduce的NodeManager只运行在从节点。

### 测试HDFS

在HDFS上创建目录/input

```
hadoop fs -mkdir /input
```

查看是否创建成功

```
hadoop fs -ls /
```



**将文件传到HDFS**

```
hadoop fs -put /data/testfile /input
```

**运行wordcount**

```
cd /apps/hadoop/share/hadoop/mapreduce/ hadoop jar hadoop-mapreduce-examples-3.0.0.jar wordcount /input/testfile /output
```

**查看结果**

```
hadoop fs -cat /output/*
```



**webUI**

[http://master:8088/](http://master:8088/) 可以查看 Hadoop 集群，节点及任务相关信息。可以看到现在活跃的节点数是4。



**HDFS Web 界面**

在浏览器中访问 http://master:9870，可以查看 HDFS 相关信息，浏览 HDFS 上的文件系统



# 6.配置spark

## 修改配置文件

在master节点修改下面三个文件

> 1.修改spark-env.sh，将SPARK_MASTER_IP的值改为master，修改后，如下图所示



> 2.修改slaves文件，将localhost改为

```
master
worker1
worker2
worker3
```

3.修改spark-defaults.conf，将spark.master 改为spark://master:7077，spark.eventLog.dir 改为 hdfs://master:9000/spark/eventLog。修改后 如下图所示



eventLog 用来存放日志，需要手动创建

```
hadoop fs -mkdir -p /spark/eventLog
```



将修改的三个文件复制到worker1和worker2两个节点，覆盖原来的文件

```
cd /apps/spark/conf
scp spark-env.sh slaves spark-defaults.conf worker1:/apps/spark/conf
scp spark-env.sh slaves spark-defaults.conf worker2:/apps/spark/conf
scp spark-env.sh slaves spark-defaults.conf worker3:/apps/spark/conf
```

至此，配置文件就修改好了，下面进行测试。

# 7.测试spark

### 启动spark集群

```
/apps/spark/sbin/start-all.sh
```

### 查看进程

主节点多了两个进程Master和Worker。从节点多了一个进程Worker



**Web UI**

查看spark管理界面，在浏览器中输入 http://master:8080，可以看到Worker有四个。

**Spark** 2.4.3  **Spark Master at spark://master:7077**

**URL:** spark://master:7077
**Alive Workers:** 4
**Cores in use:** 4 Total, 0 Used
**Memory in use:** 4.0 GB Total, 0.0 B Used
**Applications:** 0 Running, 0 Completed
**Drivers:** 0 Running, 0 Completed
**Status:** ALIVE

▼ Workers (4)

| Worker Id | Address | State | Cores | Memory |
|---|---|---|---|---|
| worker-20201221014842-192.168.43.92-7078 | 192.168.43.92:7078 | ALIVE | 1 (0 Used) | 1024.0 MB (0.0 B Used) |
| worker-20201221014847-192.168.43.185-7078 | 192.168.43.185:7078 | ALIVE | 1 (0 Used) | 1024.0 MB (0.0 B Used) |
| worker-20201221014940-192.168.43.228-7078 | 192.168.43.228:7078 | ALIVE | 1 (0 Used) | 1024.0 MB (0.0 B Used) |
| worker-20201221174852-192.168.43.117-7078 | 192.168.43.117:7078 | ALIVE | 1 (0 Used) | 1024.0 MB (0.0 B Used) |

▼ Running Applications (0)

| Application ID | Name | Cores | Memory per Executor | Submitted Time | User | State | Duration |
|---|---|---|---|---|---|---|---|

▼ Completed Applications (0)

| Application ID | Name | Cores | Memory per Executor | Submitted Time | User | State | Duration |
|---|---|---|---|---|---|---|---|

## 运行演示实例

计算PI的值会出错

```
java.lang.NoSuchMethodError: net.jpountz.lz4.LZ4BlockInputStream.<init>
(Ljava/io/InputStream;Z)Vat
```

---

原因：

> 应用在执行时对数据解码（反序列化）时，使用了默认的lz4解压缩算法，在spark-core中依赖的lz4版本是1.4，而kafka-client中依赖的lz4版本是1.3版本，在生成解压器时，版本不兼容异常。

解决办法：

> 可参考网上修改源码解决，也可通过设置"spark.io.compression.codec","snappy"或其他压缩算法规避。鉴于修改源码重新打包替换较为繁琐，建议设置其他压缩算法规避

### 读取数据

读取hdfs上的train.tsv文件，并查看数据项

```python
In [1]: from pyspark import SparkContext,SparkConf
        from pyspark.sql import SQLContext
        conf = SparkConf()
        conf.setAppName('Streaming').set('spark.io.compression.codec','snappy')
        conf.setMaster('local[2]')
        sc = SparkContext(conf = conf)
        sqlContext = SQLContext(sc)
```

```python
In [2]: row_df = sqlContext.read.format("csv")\
        .option("header","true")\
        .option("delimiter","\t")\
        .load("/input/mllib/train.tsv")
        print(row_df.count())

        7395
```

参考：https://blog.csdn.net/zhenzi_PeppaPig/article/details/84442296?utm_medium=distribute.pc_relevant_t0.none-task-blog-BlogCommendFromMachineLearnPai2-1.control&depth_1-utm_source=distribute.pc_relevant_t0.none-task-blog-BlogCommendFromMachineLearnPai2-1.control