

下载spark

```
1 wget https://archive.apache.org/dist/spark/spark-3.1.2/spark-3.1.2-bin-hadoop3.2.tgz
```

下载python环境变量 anaconda3

```
1 wget https://repo.anaconda.com/archive/Anaconda3-2021.05-Linux-x86_64.sh
```

将下载的spark解压和anaconda3(直接运行) 安装

```
1 #spark解压
2 #解压
3 tar -zxf /export/software/spark-3.1.2-bin-hadoop3.2.tgz -C /export/server/
4 #配置软连接
5 ln -s /export/server/spark-3.1.2-bin-hadoop3.2 /export/server/spark
6 #anaconda3安装
7 bash /export/software/Anaconda3-2021.05-Linux-x86_64.sh
```

配置环境变量-hadoop的全部变量, profile是对全局不包括其他客户, .bashrc是对其他用户, 或者可以解决远程无权限问题

/etc/profile和/root/.bashrc 两个文件都配置一下

```
1 #JAVA_HOME
2 JAVA_HOME=/export/server/jdk1.8.0_241
3 CLASSPATH=.:$JAVA_HOME/lib
4 PATH=$JAVA_HOME/bin:$PATH
5 export JAVA_HOME CLASSPATH PATH
6 #HADOOP_HOME
7 export HADOOP_HOME=/export/server/hadoop
8 export PATH=$PATH:$HADOOP_HOME/bin:$HADOOP_HOME/sbin
9 alias beeline="/export/server/apache-hive/bin/beeline -u jdbc:hive2://node1:10000 -n root -p 123456"
10 export SPARK_HOME=/export/server/spark
11 export PATH=$PATH:$SPARK_HOME/bin
12 export ANACONDA_HOME=/root/anaconda3
13 export PATH=$ANACONDA_HOME/bin:$PATH:/root
```

生效文件

```
1 source /etc/profile
2 source /root/.bashrc
```

查看python是否生效

```
1 #直接输入python如果是3.8则成功了
2 python
```

安装pyspark

```
1 pip install pyspark
```

在pycharm运行-- 配置连接省略了 -- 直接在pycharm运行

```
1 from pyspark import SparkConf,SparkContext
2 import os
3 os.environ['SPARK_HOME']='/export/server/spark'
4 PYSPARK_PYTHON='/root/anaconda3/bin/python'
5 os.environ['PYSPARK_PYTHON']=PYSPARK_PYTHON
6 os.environ['PYSPARK_DRIVER_PYTHON']=PYSPARK_PYTHON
7
8 if __name__ == '__main__':
9     conf= SparkConf().setAppName('wordcount').setMaster('local[*]')
10    sc=SparkContext(conf=conf)
11    rdd1=sc.textFile('file:///export/pyworkspace_sz28/pyspark_sz28/pyspark-
sparkbase_3.1.2/data/words.txt')
12    rdd2=rdd1.flatMap(lambda line:line.split(' '))
13    rdd3=rdd2.map(lambda word: (word ,1))
14    rdd4=rdd3.reduceByKey(lambda x,y:x+y)
15    list1=rdd4.collect()
16    for x in list1 : print(x)
```

读取HDFS输入输出文件-- 可以直接在pycharm运行

```
1 from pyspark import SparkConf, SparkContext
2 import os
3 import sys
4
5 # 这里选择本地pyspark环境执行spark代码
6 os.environ['SPARK_HONE'] = '/export/server/spark'
7 PYSPARK_PYTHON = '/root/anaconda3/bin/python'
8 #当存在多个版本时,不指定很可能会导致出错
9 os.environ['PYSPARK_PYTHON'] = PYSPARK_PYTHON
10 os.environ['PYSPARK_DRIVER_PYTHON'] = PYSPARK_PYTHON
11
```

```

12 if __name__ == '__main__':
13     conf = SparkConf().setAppName('wordcount').setMaster('local[*]')
14     sc = SparkContext(conf=conf)
15     rdd1 = sc.textFile('hdfs://node1:8020/pydata/words.txt')
16     rdd2 = rdd1.flatMap(lambda line: line.split(' '))
17     rdd3 = rdd2.map(lambda word: (word, 1))
18     rdd4 = rdd3.reduceByKey(lambda x, y: x + y)
19     rdd4.saveAsTextFile('hdfs://node1:8020/pydata/output2.txt')

```

上传文件

```

1 hdfs dfs -mkdir /pydata
2 hdfs dfs -put /export/pyworkspace/pyspark_sz28/pyspark-sparkbase_3.1.2/data/words.txt
  /pydata

```

使用spark-submit提交运行

```

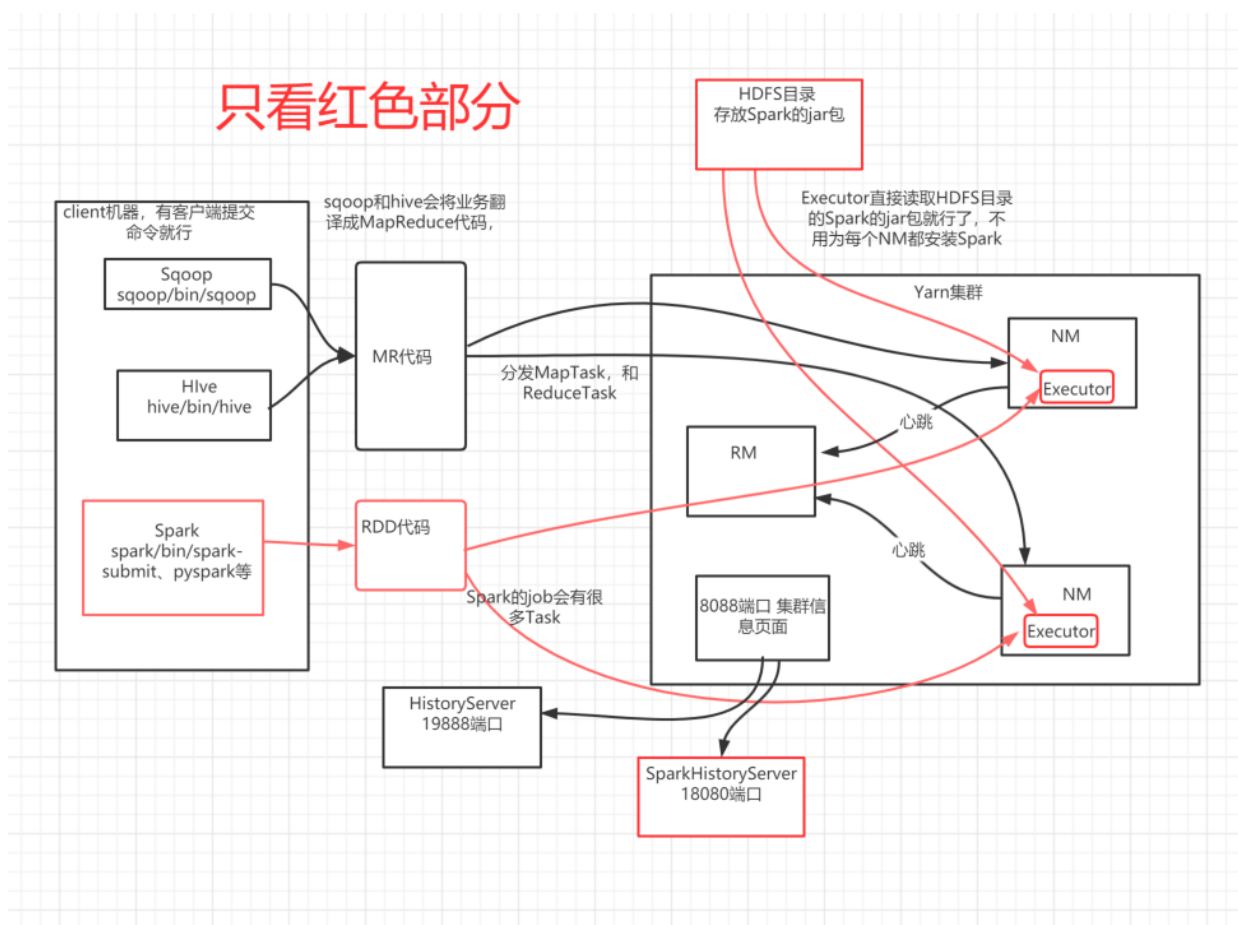
1 import sys
2
3 from pyspark import SparkConf, SparkContext
4 import os
5 os.environ['SPARK_HOME'] = '/export/server/spark'
6 PYSPARK_PYTHON = '/root/anaconda3/bin/python'
7 os.environ['PYSPARK_PYTHON'] = PYSPARK_PYTHON
8 os.environ['PYSPARK_DRIVER_PYTHON'] = PYSPARK_PYTHON
9
10 if __name__ == '__main__':
11     #1. 首先创建sparkcontext上下文环境,
12     conf = SparkConf().setAppName('wordcount').setMaster('local[*]')
13     sc = SparkContext(conf=conf)
14     #2. 从外部文件数据源读取数据
15     rdd1 = sc.textFile(sys.argv[1])
16     #3. 执行flatMap执行扁平化操作
17     rdd2 = rdd1.flatMap(lambda line: line.split(' '))
18     #4. 执行map转化操作, 得到(word, 1)
19     rdd3 = rdd2.map(lambda word: (word, 1))
20     #5. reduceByKey将相同的key的value数据累加操作
21     rdd4 = rdd3.reduceByKey(lambda x, y: x + y)
22     #6. 将结果打印
23     list1 = rdd4.collect()
24     for x in list1: print(x)

```

在linux客户端运行

```
1 spark-submit --master local[*] /export/pyworkspace/pyspark_sz28/pyspark-  
  sparkbase_3.1.2/main/wordcount_spark-submit.py /export/pyworkspace/pyspark_sz28/pyspark-  
  sparkbase_3.1.2/data/words.txt  
2
```

sparkOnYarn安装配置



每个节点都需要python3环境, 所有每台机器都装anaconda3, 可以不装pyspark

```
1 scp /export/software/Anaconda3-2021.05-Linux-x86_64.sh root@node2:$PWD  
2 scp /export/software/Anaconda3-2021.05-Linux-x86_64.sh root@node3:$PWD
```

复制hive的/bin/hive-size.xml到spark的同目录下--metastore地址

```
1 cp /export/server/hive/conf/hive-size.xml /export/server/spark/conf/
```

创建spark.env.sh

```
1 cd /export/server/spark/conf  
2 cp spark-env.sh.template spark-env.sh
```

在spark.env.sh添加

```
1 JAVA_HOME=/export/server/jdk1.8.0_241
2 ## HADOOP软件配置文件目录，读取HDFS上文件和运行YARN集群
3 HADOOP_CONF_DIR=/export/server/hadoop/etc/hadoop
4 YARN_CONF_DIR=/export/server/hadoop/etc/hadoop
5 SPARK_HISTORY_OPTS="-Dspark.history.fs.logDirectory=hdfs://node1:8020/sparklog/ -
  Dspark.history.fs.cleaner.enabled=true"
```

整合历史服务器并关闭资源检查

使用虚拟机运行服务,默认会检查,如果内存不足无法运行

在\$HADOOP_HOME/etc/hadoop/yarn-site.xml,添加:

```
1 <!-- 关闭yarn内存检查 -->
2 <property>
3     <name>yarn.nodemanager.pmem-check-enabled</name>
4     <value>>false</value>
5 </property>
6 <property>
7     <name>yarn.nodemanager.vmem-check-enabled</name>
8     <value>>false</value>
9 </property>
```

然后分发到其他节点

```
1 scp /export/server/hadoop/etc/hadoop/yarn-site.xml
  node2:///export/server/hadoop/etc/hadoop/
2 scp /export/server/hadoop/etc/hadoop/yarn-site.xml
  node3:///export/server/hadoop/etc/hadoop/
```

配置spark历史服务器spark

```
1 #手动创建HDFS目录
2 hdfs dfs -mkdir /sparklog
3 cd /export/server/spark/conf
4 cp spark-defaults.conf.template spark-defaults.conf
```

创建spark.defaults.conf

```
1 cd /export/server/spark/conf
2 cp spark.defaults.conf.template spark.defaults.conf
3 #在spark-defaults.conf添加内容
4 spark.eventLog.enabled true
```

```
5 spark.eventLog.dir                hdfs://node1:8020/sparklog/
6 spark.eventLog.compress           true
7 spark.yarn.historyServer.address   node1:18080
8 spark.yarn.jars hdfs://node1:8020/spark/jars/*
```

设置日志级别

```
1 cp log4j.properties.template log4j.properties
2 #修改log4j.properties,将INFO改成WARN 大概在19行
```

配置依赖spark jar包

```
1 ## hdfs上创建存储spark相关jar包目录
2 hadoop fs -mkdir -p /spark/jars/
3 ## 上传$SPARK_HOME/jars所有jar包
4 hadoop fs -put /export/server/spark/jars/* /spark/jars/
```

启动

启动服务:HDFS、YARN、MRHistoryServer和Spark HistoryServer

```
1 ## 启动HDFS和YARN服务, 在node1执行命令
2 start-dfs.sh
3 start-yarn.sh
4 或
5 start-all.sh
6 注意: 在onyarn模式下不需要启动spark/bin/start-all.sh (jps查看一下看到worker和master)
7 ## 启动MRHistoryServer服务, 在node1执行命令
8 mr-jobhistory-daemon.sh start historyserver
9 ## 启动Spark HistoryServer服务, , 在node1执行命令
10 /export/server/spark/sbin/start-history-server.sh
11
```

测试

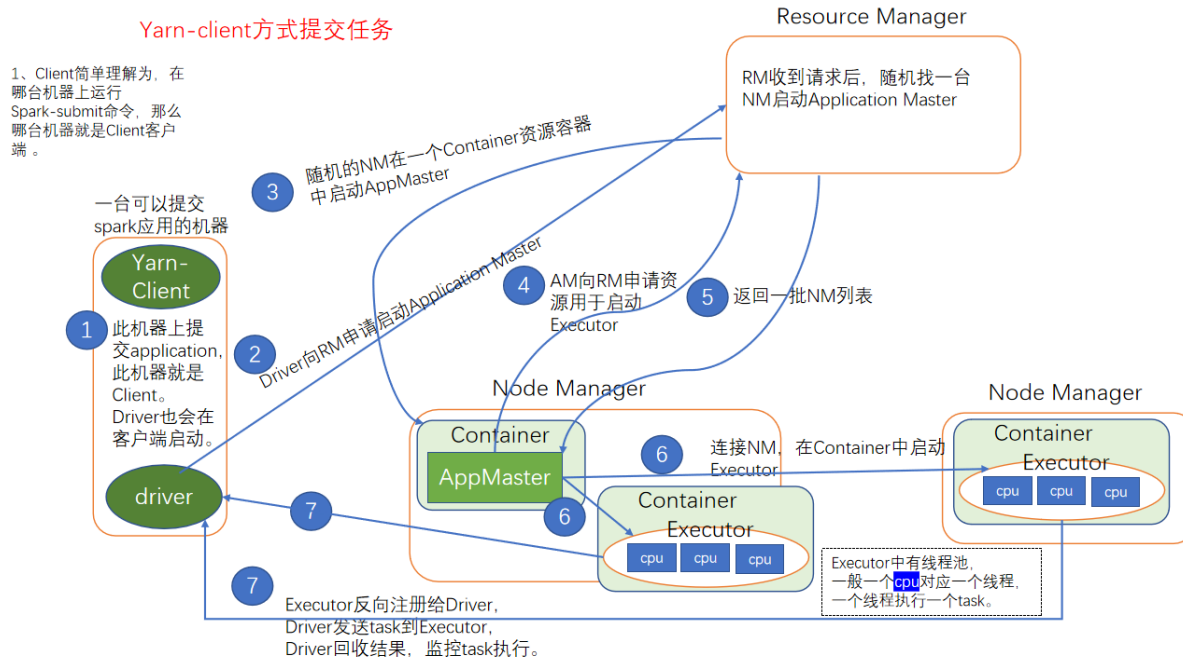
client模式

```
1 SPARK_HOME=/export/server/spark
2 ${SPARK_HOME}/bin/spark-submit \
3 --master yarn \
4 --deploy-mode client \
5 --driver-memory 512m \
6 --executor-memory 512m \
7 --executor-cores 1 \
```

```

8  --num-executors 2 \
9  --queue default \
10 --conf "spark.pyspark.driver.python=/root/anaconda3/bin/python3" \
11 --conf "spark.pyspark.python=/root/anaconda3/bin/python3" \
12 ${SPARK_HOME}/examples/src/main/python/pi.py \
13 10

```



cluster 模式

```

1  ${SPARK_HOME}/bin/spark-submit \
2  --master yarn \
3  --deploy-mode cluster \
4  --driver-memory 512m \
5  --executor-memory 512m \
6  --executor-cores 1 \
7  --num-executors 2 \
8  --queue default \
9  --conf "spark.pyspark.driver.python=/root/anaconda3/bin/python3" \
10 --conf "spark.pyspark.python=/root/anaconda3/bin/python3" \
11 ${SPARK_HOME}/examples/src/main/python/pi.py \
12 10

```

```

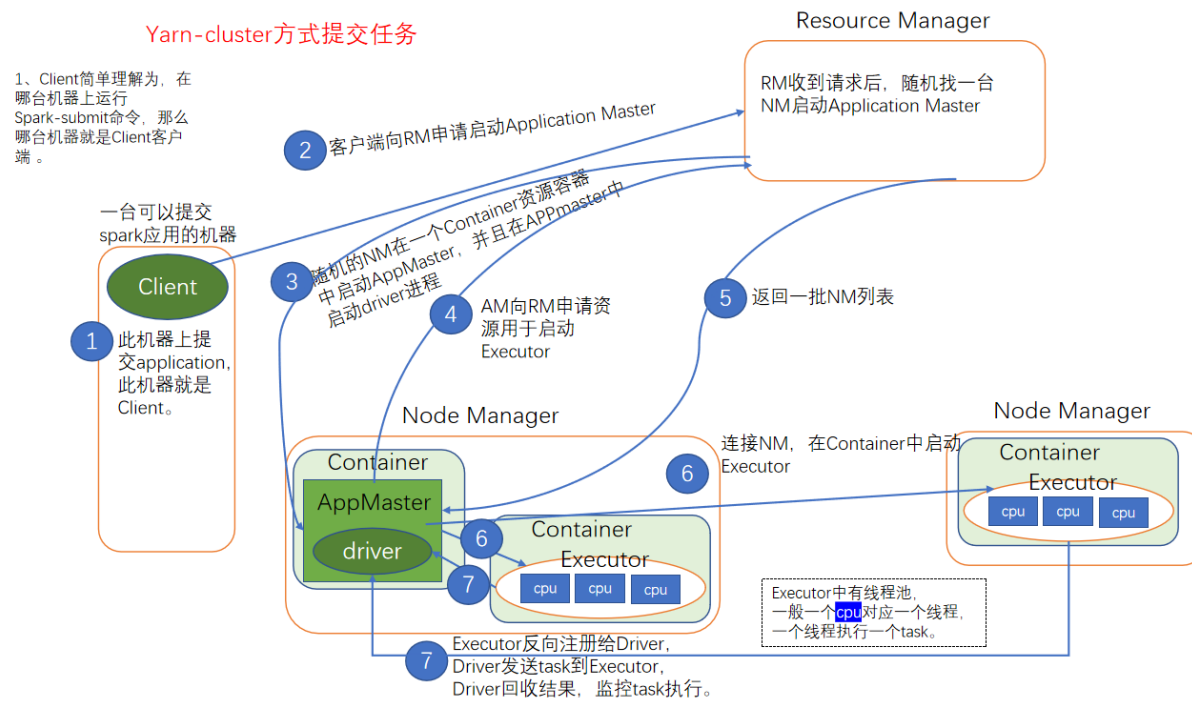
1  spark-submit \
2  --master yarn \
3  --deploy-mode cluster \

```

```

4 --driver-memory 512m \
5 --executor-memory 512m \
6 --executor-cores 1 \
7 --num-executors 2 \
8 --queue default \
9 --class bigdata.spark.sql.Spark06_SparkSQL \
10 root/spark-1.0-SNAPSHOT.jar
11

```



参数说明

```

1 --master spark使用什么资源管理器运行,比如yarn
2 --deploy-mode 模式:client和cluster
3 --name 给一个程序起一个名字[可选]
4 --class 包名.包类 [可选]
5 --driver-core 给driver申请几个核
6 --driver-memory 申请driver进程的内存
7 --executor-memory 申请每个executor内存
8 --executor-core 申请每个executor几个核
9 --num-executor 一共申请几个executor
10 --queue 用来隔离CPU和内存资源,每个队列中都包含指定容量的CPU和内存

```

启动spark端口

```

1 nohup /export/server/hive/bin/hive --service metastore 2>&1 > /tmp/hive-metastore.log &
2 /export/server/spark/sbin/start-thriftserver.sh \

```



```
3  --hiveconf hive.server2.thrift.port=10001 \  
4  --hiveconf hive.server2.thrift.bind.host=node1 \  
5  --master local[*]
```

spark和hive乱码（查询）

```
1  -- 下面修改是在MySQL中修改，因为MySQL记录维护着元数据  
2  
3  use hive3;  
4  
5  -- （1）修改表字段注解和表注解  
6  
7  alter table COLUMNS_V2 modify column COMMENT varchar(256) character set utf8;  
8  alter table TABLE_PARAMS modify column PARAM_VALUE varchar(4000) character set utf8;  
9  
10 -- （2）修改分区字段注解  
11  
12 alter table PARTITION_PARAMS modify column PARAM_VALUE varchar(4000) character set utf8  
13 ;  
14 alter table PARTITION_KEYS modify column PKEY_COMMENT varchar(4000) character set utf8;  
15  
16 # （3）修改索引注解  
17  
18 alter table INDEX_PARAMS modify column PARAM_VALUE varchar(4000) character set utf8;  
19  
20 -- 下面这个在Hive的配置文件中修改  
21  
22 # <!-- 存储元数据mysql相关配置 -->  
23 # <property>  
24 #   <name>javax.jdo.option.ConnectionURL</name>  
25 #   <value>jdbc:mysql://node1:3306/hive3?  
26     createDatabaseIfNotExist=true&useSSL=false&useUnicode=true&characterEncoding  
27     =UTF-8</value>  
28 # </property>  
29 #  
30 # 在IDEA/pycharm中修改  
31 # tools（工具）-部署-advanced（配置）-把GBK改成utf-8
```