海豚调度器：http://42.194.186.185:12345/dolphinscheduler

HBase：http://42.194.186.185:16010/master-status

脚本

```bash
1  #/bin/bash
2  #nohup不挂起，&后台运行 ，2>1& > path:不管是正常还是报错的信息都写入
3  #zookeeper启动服务
4  echo '正在启动zookeeper'
5  hosts=(node1 node2 node3)
6  for host in ${hosts[*]}
7  do
8    ssh $host "source /etc/profile;/export/server/zookeeper/bin/zkServer.sh start"
9  done
10
11 #HBase启动服务--只需要启动一个
12 echo '正在启动HBase'
13 /export/server/hbase/bin/start-hbase.sh
14
15 #hive的metastore和hiveserver2的服务
16 echo '正在启动hive-metastore'
17 nohup /export/server/hive/bin/hive --service metastore 2>&1 > /tmp/hive-metastore.log &
18 echo '正在启动hive-hiveserver2'
19 nohup /export/server/hive/bin/hive --service hiveserver2 2>&1 > /tmp/hive-
   hiveserver2.log &
20 #nohup /export/server/hive/bin/hive --service metastore &
21 #nohup /export/server/hive/bin/hive --service hiveserver2 &
22 #启动MR和spark的历史服务
23 echo '正在启动MR-history'
24 /export/server/hadoop/sbin/mr-jobhistory-daemon.sh start historyserver
25 echo '正在启动spark-history'
26 /export/server/spark/sbin/start-history-server.sh
27 #启动spark连接客户端相关参数
28 echo '正在启动spark客户端连接'
29 /export/server/spark/sbin/start-thriftserver.sh \
30 --hiveconf hive.server2.thrift.port=10001 \
31 --hiveconf hive.server2.thrift.bind.host=node1 \
32 --master local[*]
33
34 #明细定义
35 start-thriftserver.sh \
```

```
36   --name sparksql-thrift-server \
37   --master yarn \
38   --deploy-mode client \
39   --driver-cores 4 \
40
     --driver-memory 12g \
41   --hiveconf hive.server2.thrift.http.port=10001 \
42   --num-executors 8 \
43   --executor-memory 4g \
44   --conf spark.sql.shuffle.partitions=2
45   #jps -m
46   jps -m
47   #第二种方案(企业)
48   #/export/server/spark/sbin/start-thriftserver.sh \
49   #--hiveconf hive.server2.thrift.port=10001 \
50   #--hiveconf hive.server2.thrift.bind.host=node1 \
51   #--master yarn \
52   #--num-executors 50 \
53   #--executor-cores 4 \
54   #--executor-memory 12G \
55   echo '正在启动海豚调度器'
56   /export/server/dolphinscheduler/bin/start-all.sh
57    echo '查看zookeeper运行状态'
58    /export/server/zookeeper/bin/zkServer.sh status
59
60    echo "正在启动kafka"
61    nohup /export/server/kafka/bin/kafka-server-start.sh
     /export/server/kafka/config/server.properties 2>&1 > /tmp/kafka.log &
62
63    echo "启动Phoenix客户端"
64    /export/server/phoenixbin/python2 sqlline.py node1:2181
65
```

```
1   kafka
2   #kafka创建主题
3   kafka-topics.sh --create --bootstrap-server node1:9092,node2:9092,node3:9092 --
    partitions 3 --replication-factor 1 --topic ODS_BASE_LOG_1018
4   #kafka生产者
5   kafka-console-producer.sh --broker-list node1:9092,node2:9092,node3:9092 --topic test001
```

```
6   #kafka消费者
7   kafka-console-consumer.sh  --bootstrap-server node1:9092,node2:9092,node3:9092 --topic
    ODS_BASE_LOG_1018
8   #kafka查看主题
9   kafka-topics.sh --bootstrap-server node1:9092,node2:9092,node3:9092 --list
10  #kafka查看主题详情
11  kafka-topics.sh --bootstrap-server node1:9092 --describe --topic test001
12
13  flume
14  #启动flume
15  flume-ng agent -c conf -f /export/server/flume/conf/momo_taildir_source_kafka_sink.conf
    -n a1 -Dflume.root.logger=INFO,console
16
17  phoneix
18  #启动phoneix客户端
19  /export/server/phoenix/bin
    /python2 sqlline.py node1:2181
20
21  启动陌陌jar包
22  java -jar /export/data/momo_init/MoMo_DataGen.jar MoMo_Data.xlsx /export/data/momo_data/
    1000
23
```

## 下载

```
1   #dolphinscheduler调度器
2   wget https://archive.apache.org/dist/dolphinscheduler/1.3.5/apache-dolphinscheduler-
    incubating-1.3.5-dolphinscheduler-bin.tar.gz
3   #MySQL
4   wget https://dev.mysql.com/get/Downloads/MySQL-8.0/mysql-8.0.27-1.el8.x86_64.rpm-
    bundle.tar
5   #spark
6   wget https://archive.apache.org/dist/spark/spark-3.1.2/spark-3.1.2-bin-hadoop3.2.tgz
7   #anaconda3
8   wget https://repo.anaconda.com/archive/Anaconda3-2021.05-Linux-x86_64.sh
9   #hadoop未编译
10  wget https://archive.apache.org/dist/hadoop/common/hadoop-3.3.0/hadoop-3.3.0-src.tar.gz
11  #zookeeper
12  wget https://archive.apache.org/dist/zookeeper/zookeeper-3.7.0/apache-zookeeper-3.7.0-
    bin.tar.gz
13  #hive
14  wget https://archive.apache.org/dist/hive/hive-3.1.2/apache-hive-3.1.2-bin.tar.gz
```

```
15  #sqoop
16  wget https://archive.apache.org/dist/sqoop/1.4.7/sqoop-1.4.7.bin__hadoop-2.6.0.tar.gz
17  #sqoop-commns-lang.jar包
18  wget https://repo1.maven.org/maven2/commons-lang/commons-lang/2.6/commons-lang-2.6.jar
19  #sqoop-hive.jar包
20  wget https://repo1.maven.org/maven2/org/apache/hive/hive-common/3.1.2/hive-common-
    3.1.2.jar
21  #MySQL8.0jar包
22  wget https://repo1.maven.org/maven2/mysql/mysql-connector-java/8.0.27/mysql-connector-
    java-8.0.27.jar
23  #MySQL5.0jar包
24  wget https://repo1.maven.org/maven2/mysql/mysql-connector-java/5.1.49/mysql-connector-
    java-5.1.49.jar
25  #HBase
26  wget https://archive.apache.org/dist/hbase/2.1.0/hbase-2.1.0-bin.tar.gz
27  #Phoenix
28  wget https://archive.apache.org/dist/phoenix/apache-phoenix-5.0.0-HBase-2.0/bin/apache-
    phoenix-5.0.0-HBase-2.0-bin.tar.gz
29  #kafka
30  wget https://archive.apache.org/dist/kafka/2.4.1/kafka_2.12-2.4.1.tgz
31  #flume
32  wget https://archive.apache.org/dist/flume/1.9.0/apache-flume-1.9.0-bin.tar.gz
```

python安装库

```
pip install pyarrow,pyspark
```

python中创建使用pyspark

```
1   import string,time,os
2   from numpy import double
3   from pyspark import StorageLevel
4   from pyspark.sql import SparkSession, Row,functions
5   from pyspark.sql.functions import udf,pandas_udf
6   from pyspark.sql.types import *
7   import pandas as pd
8
9   # 这里选择本地pyspark环境执行spark代码
10  os.environ['JAVA_HOME']='/export/server/jdk'
11  os.environ['SPARK_HONE'] = '/export/server/spark'
12  PYSPARK_PYTHON = '/root/anaconda3/bin/python'
13  # 当存在多个版本时,不指定很可能会导致出错
```

```python
14  os.environ['PYSPARK_PYTHON'] = PYSPARK_PYTHON
15  os.environ['PYSPARK_DRIVER_PYTHON'] = PYSPARK_PYTHON
16  if __name__ == '__main__':
17      # 1.创建SparkSession上下文对象
18      spark = SparkSession.builder.appName('spark_no_hive').master('local[*]')\
19          .config('hive.metastore.uris','thrift://node1:9083')\
20          .config('hive.metastore.warehouse.dir','/user/hive/warehouse')\
21          .config('spark.sql.shuffle.partitions','4')\
22          .enableHiveSupport()\
23          .getOrCreate()
24      spark.sparkContext.setLogLevel('WARN')
25      #Arrow 是一种内存中的列式数据格式，用于spark中以jvm和python进程之间有效的传输数据
26      #需要安装arrow pip install pyarrow
27      #2.开启pyarrow，能加快计算速度，原理2个：1.基于内存减少了序列化和反序列化开销，2.基于向量
    计算vectorize
28      spark.conf.set('spark.sql.execution.arrow.pyspark.enabled','true')
```

spark写出

```python
1   #3.将DataFrame数据保存为text格式
2   #保存为text格式，只能写出一列，不能写出多列
3   df.select(functions.concat_ws('_','user_id','movieID','rating','timestamp')).write.mode(
    'overwrite').format('text').save('file://./out1/')
4   #4.将DataFrame数据保存为csv格式 -- header 不忽略字段
5   df.coalesce(1).write.mode('overwrite').format('csv').option('sep',',').option('header',T
    rue).save('file://./csv')
6   #5.将DataFrame数据保存为json格式
7   df.coalesce(1).write.mode('overwrite').format('json').save('file://./json')
8   #6.将DataFrame数据保存为parquet格式
9   df.coalesce(1).write.mode('overwrite').format('parquet').save('file://./parquet')
10  #parquet 语法2 --spark默认支持的格式就是parquet
11  df.coalesce(1).write.mode('overwrite').save('file://./parquet2')
12  #mysql连接
13  df.write.format('jdbc').mode('overwrite').option('url',
    'jdbc:mysql://node1:3306/insurance') \
14      .option('dbtable', 'policy_actuary') \
15      .option('user', 'root').option('password', '123456').save()
16
17  #读取mysql
18  df=spark.read.format('jdbc') \
19      .option('url', 'jdbc:mysql://42.194.186.185:3306/?
    serverTimezone=UTC&characterEncoding=utf8&useUnicode=true') \
```

```python
        .option('dbtable', 'bigdata.tb_top10_movies') \
        .option('user', 'root').option('password', '123456') \
        .load()
spark=SparkSession.builder.appName("word").mater("local[*]").getOrcreate()
#2-加载text文件形成DataFrame
df=spark.read.format('text').load('path')'
#简化写法
df1=spark.read.text('path')

#3-加载csv文件形成DataFrame
df2=spark.read.format('csv')\
    .option('sep',';')\
    .option('header',True)\
    .option('encoding','utf-8')\
    .option('inferSchema',True)\
    .load(path)
#4-加载json文件形成DataFrame
df3=spark.read.format('json').load('path')

#5-加载parquet文件形成DataFrame
df4=spark.read.format('parquet')\
    .load(path)
 #6-添加schema加载文件形成DataFrame
df=spark.read.schema('id int,name string, score int')\
    .csv('path')
```

全局

```
#JAVA_HOME
JAVA_HOME=/export/server/jdk
CLASSPATH=.:$JAVA_HOME/lib
PATH=$JAVA_HOME/bin:$PATH
export JAVA_HOME CLASSPATH PATH
#HADOOP_HOME
export HADOOP_HOME=/export/server/hadoop
export PATH=$PATH:$HADOOP_HOME/bin:$HADOOP_HOME/sbin
#hive和spark的免账号密码登录
alias beelinehive="/export/server/apache-hive/bin/beeline -u jdbc:hive2://node1:10000 -n root -p 123456"
alias beelinespark="/export/server/spark/bin/beeline -u jdbc:hive2://node1:10001 -n root -p 123456"
```

```bash
12  #SPARK_HOME
13  export SPARK_HOME=/export/server/spark
14  export PATH=$PATH:$SPARK_HOME/bin
15  #ANACONDA_HOME
16  export ANACONDA_HOME=/root/anaconda3
17  export PATH=$ANACONDA_HOME/bin:$PATH:/root
18  #HIVE
19  export HIVE_HOME=/export/server/hive
20  export PATH=$PATH:$HIVE_HOME/bin
21  #SQOOP_HOME
22  export SQOOP_HOME=/export/server/sqoop
23  export PATH=$PATH:$SQOOP_HOME/bin
24  #ZOOKEEPER_HOME
25  export ZOOKEEPER_HOME=/export/server/zookeeper/
26  export PATH=$PATH:$ZOOKEEPER_HOME/bin
27  #HBASE_HOME
28  export HBASE_HOME=/export/server/hbas
29  export PATH=$PATH:$HBASE_HOME/bin
30  #KAFKA_HOME
31  export KAFKA_HOME=/export/server/kafka
32  export PATH=:$PATH:${KAFKA_HOME}/bin
33  #phoenix_HOME
34  export PHONEIX_HOME=/export/server/phoenix
35  #FLUME
36  export FLUME_HOME=/export/server/flume
37  export PATH=$PATH:$FLUME_HOME/bin
38
```

```sql
1  【计算链条比较长，每一步骤如果有精度损失，那么叠加到最后，可能会导致结果不准确,则设置如下】
2  set spark.sql.decimalOperations.allowPrecisionLoss=false
3
4  【spark和vine导出打印表头，有数据库名】
5  set hive.cli.print.header=true;
6
7  【spark和vine导出打印表头，去数据库名】
8  set hive.cli.print.header=true;set hive.resultset.use.unique.column.names=false;
9
10  【spark设置shuffle并行度】
```

```
11  spark.sql('set spark.sql.shuffle.partitions=4')
12  set spark.sql.shuffle.partitions=4
13
14  【hdfs 查看总文件大小】
15  hdfs dfs -du -h path
16
17  【hdfs副本数】
18  dfs.replication
19
20  【修复MySQL元数据，和分区表的实际目录一致】
21  msck repair table 库.表;
22
23  【开启hive中的mapjoin】--大表的数据会分成多个task和整张小表去join
24  set hive.auto.convert.join=true [hive2默认为true]
25
26  【hive和spark的mapjoin小表大小】
27  set hive.auto.convert.join.noconditionaltask.size=10000000;【默认为10M单位字节】
28  set spark.sql.autoBroadcastJoinThreshold=51200000;
29
30  【hive的bucket mapjoin 分桶】
31  set hive.optimize.bucketmapjoin = true;【默认为true】
32
33  【sortmerge bucket mapjoin（大表join大表smb）】
34  set hive.auto.convert.join=true;
35  set hive.optimize.bucketmapjoin = true;   -- 开启 bucket map join
36  set hive.auto.convert.sortmerge.join=true; -- 开启 SBM join支持--归并
37  set hive.optimize.bucketmapjoin.sortedmerge = true;   -- 自动尝试开启 SMB join
38  --set hive.enforce.sorting=true;   -- 开启强制排序
39  --set hive.enforce.bucketing=true;-- 注释的只能在配置文件写死或者不执行
40  --set hive.input.format=org.apache.hadoop.hive.ql.io.BucketizedHiveInputFormat;
41  set hive.auto.convert.join.noconditionaltask = true;
42  --下面可以不设置
43  set hive.auto.convert.join.noconditionaltask.size = 50000000;
44  set hive.auto.convert.sortmerge.join.bigtable.selection.policy
45     = org.apache.hadoop.hive.ql.optimizer.TableSizeBasedBigTableSelectorForAutoSMJ;
46
47  【MR的task内存】
48  set mapreduce.map.java.opts=-Xmx6000m;
49  set mapreduce.map.memory.mb=6096;
```

```
50  set mapreduce.reduce.java.opts=-Xmx6000m;
51  set mapreduce.reduce.memory.mb=6096;
52
53  【MR缓冲区大小】
54  set mapreduce.task.io.sort.mb=100
55
56  【MR环形缓冲区溢出的阈值】
57  mapreduce.map.sort.spill.percent=0.8
58
59  【reduce拉取并行度】
60  set mapreduce.reduce.shuffle.parallelcopies=8
61  set mapreduce.reduce.shuffle.read.timeout=180000
62
63  【MR小文件处理】
64  #设置Hive中底层MapReduce读取数据的输入类：将所有文件合并为一个大文件作为输入
65  set hive.input.format=org.apache.hadoop.hive.ql.io.CombineHiveInputFormat;
66  #如果hive的程序，只有maptask，将MapTask产生的所有小文件进行合并
67  set hive.merge.mapfiles=true;
68  #如果hive的程序，有Map和ReduceTask,将ReduceTask产生的所有小文件进行合并
69
    set hive.merge.mapredfiles=true;
70  #每一个合并的文件的大小
71
    set hive.merge.size.per.task=256000000;
72  #平均每个文件的大小，如果小于这个值就会进行合并
73
    set hive.merge.smallfiles.avgsize=16000000;
74
75  【查看hive和spark的执行计划】
76  explain select...
77
78  【markdown代码折叠】
79  <details>
80  <summary><b>点击查看完整代码</b></summary>
81  <pre><code>
82  </code></pre>
83  </details>
84  【Linux关闭邮件提醒】
85  echo "unset MAILCHECK" >> /etc/profile
86  source /etc/profile
87
```

```
88  【Linux ifconfig没有ip】
89  [方案1推荐]
90  [root@localhost~]# ifup ens33
91  Error: Connection activation failed: No suitable device found for this
    connection(device lo not availa ble because device is strictly unmanaged).
92
93  [root@localhost~]# chkconfig NetworkManager off
94  Note: Forwarding request to 'systemctI disable NetworkManager.service'.
95  Removed symlink /etc/systemd/system/multi-user.target.wants/NetworkManager.service.
96  Removed symlink /etc/systemd/system/dbus-org.freedesktop.NetworkManager.service.
97  Removed symlink /etc/systemd/system/dbus-org.freedesktop.nm-dispatcher.service.
98  Removed symlink /etc/systemd/system/network-online.target.wants/NetworkManager-wait-
    online.service.
99
100 [root@localhost~]# chkconfig network on
101 [root@localhost~]# service NetworkManager stop
102 Redirecting to /bin/systemctl stop NetworkManager.service
103 [root@localhost~]# service network start
104 Starting network (via systemct1):[ ok]
105
106 [方案2]
107 ip add show
108 将link/ether后的MAC地址写到ens33中的HWADDR
109 vim /etc/sysconfig/network-scripts/ifcfg-ens33
110
111 【Linux搜索指定文件类型的指定内容】
112 grep -rin "9820" --include "*.xml" /export/server/
113 find / -name "*.log" | grep "error"
114
115 【SparkSQL报错Filesystem Closed】
116 修改hadoop配置文件core-site.xml：
117 <property>
118         <name>fs.hdfs.impl.disable.cache</name>
119         <value>true</value>
120 </property>
121
122 【启动hive的metastore、hiveserver2服务】
123 nohup hive --service metastore  2>&1 > /tmp/hive-metastore.log &
124 nohup hive --service hiveserver2 2>&1 > /tmp/hive-hiveserver2.log &
125
```

```
126
127  【hive元数据初始化和更新】
128  schematool -dbType mysql -initSchema
129  schematool -dbType mysql -upgradeSchema
130
131  【mysql的授权语句】
132  GRANT ALL PRIVILEGES ON hive.* TO 'root'@'%' IDENTIFIED BY '123456';
133  flush privileges;
134
135  【hive的jdbc连接地址】
136  jdbc:hive2://node1:10000
137
138
139  /export/server/zookeeper/bin/zkServer.sh stop
140  /export/server/zookeeper/bin/zkServer.sh start
141  /export/server/zookeeper/bin/zkServer.sh status
142
143  /export/server/dolphinscheduler/bin/stop-all.sh
144  /export/server/dolphinscheduler/bin/start-all.sh
145
146
147
148  【开启Hive的本地模式】
149  set hive.exec.mode.local.auto=true;(默认为false)
150  【hive开启负载均衡】平均分区，让每个分区都是同样的数据，如果多个去重则会报错，这个针对group by
151  hive.groupby.skewindata=true
152
153
154
155  【hadoop退出安全模式】
156  hdfs dfsadmin -safemode leave
157  hdfs dfsadmin -safemode forceExit
158
159  【重新启动zk和kafka】
160  zkServer.sh stop
161  kafka-server-stop.sh
162  zkServer.sh start
163  nohup /export/server/kafka/bin/kafka-server-start.sh
     /export/server/kafka/config/server.properties &
164
```

```
165  【彻底删除kafka主题】

166  kafka-topics.sh --zookeeper node3:2181 --list

167  kafka-topics.sh --zookeeper node3:2181 --delete --topic spark_kafka

168

169  【zkCli.sh】

170  ls /config/topics

171  rmr /config/topics/spark_kafka

172  rmr /brokers/topics/spark_kafka

173  rmr /admin/delete_topics/spark_kafka

174

175  【创建主题】

176  kafka-topics.sh --zookeeper node3:2181 --create --topic spark_kafka --partitions 3 --
     replication-factor 1

177  kafka-topics.sh --zookeeper node3:2181 --list

178

179  【启动生产者和消费者】

180  kafka-console-producer.sh --broker-list node3:9092 --topic spark_kafka

181  kafka-console-consumer.sh --from-beginning --bootstrap-server node3:9092 --topic
     spark_kafka

182  kafka-console-consumer.sh --from-beginning --bootstrap-server node3:9092 --topic
     __consumer_offsets

183

184

185  【完整的删除再重建主题步骤start

186  关闭生产者和消费者

187

188  kafka-topics.sh --zookeeper node3:2181 --list

189  kafka-topics.sh --zookeeper node3:2181 --delete --topic stationTopic

190

191  kafka-topics.sh --zookeeper node3:2181 --list

192

193  zkCli.sh

194  ls /config/topics

195  rmr /config/topics/stationTopic

196  rmr /brokers/topics/stationTopic

197  rmr /admin/delete_topics/stationTopic

198

199  删除log文件目录

200

201  kafka-topics.sh --zookeeper node3:2181 --create --topic stationTopic --partitions 3 --
     replication-factor 1
```

```
202  kafka-topics.sh --zookeeper node3:2181 --list
203  再次重启kafka
204  kafka-server-stop.sh
205  nohup /export/server/kafka/bin/kafka-server-start.sh
     /export/server/kafka/config/server.properties &
206
207  启动生产者
208  kafka-console-producer.sh --broker-list node3:9092 --topic stationTopic
209  启动消费者
210  kafka-console-consumer.sh --bootstrap-server node3:9092 --topic stationTopic --from-
     beginning
211  完整的删除再重建主题步骤end】
212
213
214  【安装redis】
215  tar
216  yum install gcc
217  cd /redis3.2.8
218  make
219  make PREFIX=/usr/local/src/redis install
220  cp /redis3.2.8/redis.conf /redis/bin/redis.conf
221  启动redis
222  redis-server redis.conf
223  redis-cli -h 192.168.88.163
224
225  【启动spark-thriftserver】
226  /export/server/spark/sbin/start-thriftserver.sh \
227    --hiveconf hive.server2.thrift.port=10001 \
228    --hiveconf hive.server2.thrift.bind.host=node1 \
229    --master local[*]
230
231  /export/server/spark/sbin/start-thriftserver.sh \
232  --hiveconf hive.server2.thrift.port=10001 \
233  --hiveconf hive.server2.thrift.bind.host=p1 \
234  --master local[*]
235
236
237  /export/server/spark/sbin/start-thriftserver.sh \
238  --hiveconf hive.server2.thrift.port=10001 \
239  --hiveconf hive.server2.thrift.bind.host=p1 \
```

```
240  --master yarn \
241  --deploy-mode client \
242  --driver-memory 1g \
243  --executor-memory 2g \
244  --executor-cores 4 \
245  --num-executors 25 \
246
247  !connect jdbc:hive2://192.168.88.166:10001
248
249  【kafka】
250  启动Zookeeper 服务
251  zookeeper-daemon.sh start
252
253  启动Kafka 服务
254  kafka-daemon.sh start
255
256
257  【免秘钥登录】
258  ssh-keygen -t rsa
259  ssh-copy-id node1
260  scp /root/.ssh/authorized_keys node2:/root/.ssh
261  scp /root/.ssh/authorized_keys node3:/root/.ssh
262
263  【sqoop】
264  sqoop import \
265  -D mapred.job.name=sqoop_import_dd_table  \
266  --connect "jdbc:mysql://192.168.88.163:3306/insurance"  \
267  --username root  \
268  --password "123456"  \
269  --table dd_table  \
270  --hive-import  \
271  --hive-database insurance  \
272  --hive-table dd_table  \
273  --hive-overwrite  \
274  -m 1  \
275  --fields-terminated-by ','  \
276  --delete-target-dir
277
278  【mysql】
```

```
279  mysqldump insurance --add-drop-trigger --result-file=/opt/insurance/insurance.sql --
     user=root --host=192.168.88.163 --port=3306
280  mysql -uusername -ppassword -h ip_address -P 3306 --default-character-set = utf8mb4 $
     {i} </ data / mysql_back / $ {i} .sql
281  mysql -uroot -p  --default-character-set=utf8mb4 insurance </opt/insurance/insurance.sql
282
283
284  【dolphinscheduler】
285  CREATE DATABASE dolphinscheduler DEFAULT CHARACTER SET utf8 DEFAULT COLLATE
     utf8_general_ci;
286  set global validate_password.policy=0;
287  set global validate_password.length=1;
288
289  GRANT ALL PRIVILEGES ON dolphinscheduler.* TO '{user}'@'%' IDENTIFIED BY '{password}';
290  flush privileges;
291
292  安装dolphinscheduler 要启动zookeeper
293  一定要配在ds的env中配置sqoop的环境变量
294  需要root给租户授权最大。
295  谁将脚本上传到Linux，谁就记得授执行权限。
296
297  安装dos2unix命令
298   dos2unix /opt/insurance/sqoop/sqoop_import_mort_10_13.sh
299
300  修改工作流，要先下线
301
302
303  【structured Streaming】
304  --memory sink
305  CREATE TABLE db_spark.tb_word_count (
306    id int NOT NULL AUTO_INCREMENT,
307    word varchar(255) NOT NULL,
308    count int NOT NULL,
309    PRIMARY KEY (id),
310    UNIQUE KEY word (word)
311  ) ENGINE=InnoDB DEFAULT CHARSET=utf8mb4 COLLATE=utf8mb4_0900_ai_ci;
312
313  REPLACE INTO  tb_word_count (id, word, count) VALUES (NULL, ?, ?);
314
315
316  【spark yarn Pi】
```

```
317  /export/server/spark/bin/spark-submit \
318  --master yarn \
319  --class org.apache.spark.examples.SparkPi \
320  ${SPARK_HOME}/examples/jars/spark-examples_2.11-2.4.5.jar \
321  10
322
323  【WordCount yarn】
324  /export/server/spark/bin/spark-submit \
325  --master yarn \
326  --driver-memory 512m \
327  --executor-memory 512m \
328  --executor-cores 1 \
329  --num-executors 2 \
330  --queue default \
331  --class cn.itcast.spark._2SparkWordCount \
332  /opt/spark-chapter01-1.0-SNAPSHOT.jar
333
334
335
336  【 Run application locally on 8 cores】
337  /export/server/spark/bin/spark-submit \
338    --class org.apache.spark.examples.SparkPi \
339    --master local[8] \
340  ${SPARK_HOME}/examples/jars/spark-examples_2.11-2.4.5.jar \
341    100
342
343  # Run on a Spark standalone cluster in client deploy mode
344  ./bin/spark-submit \
345    --class org.apache.spark.examples.SparkPi \
346    --master spark://207.184.161.138:7077 \
347    --executor-memory 20G \
348    --total-executor-cores 100 \
349  ${SPARK_HOME}/examples/jars/spark-examples_2.11-2.4.5.jar \
350    1000
351
352  # Run on a Spark standalone cluster in cluster deploy mode with supervise
353  ./bin/spark-submit \
354    --class org.apache.spark.examples.SparkPi \
355    --master spark://207.184.161.138:7077 \
356    --deploy-mode cluster \
```

```
357    --supervise \
358    --executor-memory 20G \
359    --total-executor-cores 100 \
360    /path/to/examples.jar \
361    1000
362
363  # Run on a YARN cluster
364  export HADOOP_CONF_DIR=XXX
365  ./bin/spark-submit \
366    --class org.apache.spark.examples.SparkPi \
367    --master yarn \
368    --deploy-mode cluster \  # can be client for client mode
369    --executor-memory 20G \
370    --num-executors 50 \
371    /path/to/examples.jar \
372    1000
373
374  # Run a Python application on a Spark standalone cluster
375  ./bin/spark-submit \
376    --master spark://207.184.161.138:7077 \
377    examples/src/main/python/pi.py \
378    1000
379
380  # Run on a Mesos cluster in cluster deploy mode with supervise
381  ./bin/spark-submit \
382    --class org.apache.spark.examples.SparkPi \
383    --master mesos://207.184.161.138:7077 \
384    --deploy-mode cluster \
385    --supervise \
386    --executor-memory 20G \
387    --total-executor-cores 100 \
388    http://path/to/examples.jar \
389    1000
390
391  # Run on a Kubernetes cluster in cluster deploy mode
392  ./bin/spark-submit \
393    --class org.apache.spark.examples.SparkPi \
394    --master k8s://xx.yy.zz.ww:443 \
395    --deploy-mode cluster \
396    --executor-memory 20G \
```

```
397    --num-executors 50 \
398    http://path/to/examples.jar \
399    1000
400
```

```
1   hive动态分区
2   -----------------因为采用了动态分区插入技术 因此需要设置相关参数--------------
3   --分区
4   SET hive.exec.dynamic.partition=true;
5   SET hive.exec.dynamic.partition.mode=nonstrict;
6   set hive.exec.max.dynamic.partitions.pernode=10000;
7   set hive.exec.max.dynamic.partitions=100000;
8   set hive.exec.max.created.files=150000;
9   --hive压缩
10  set hive.exec.compress.intermediate=true;
11  set hive.exec.compress.output=true;
12  --写入时压缩生效
13  set hive.exec.orc.compression.strategy=COMPRESSION;
```

```
1   --开启分桶，如果不开启，是不会启动多个reduce分桶的
2   set hive.enforce.bucketing=true;
3   --2.x版本可以通过以下参数禁止使用load语句加载数据到表中
4   set hive.strict.checks.bucketing = true;
```

常见的文件类型：textfile/orc/parquet

```
1   hiveorc索引
2   stored as orc ('orc.create.index'='true')
3   查询数据时，开启row group index 过滤查询
4   hive.optimize.index.filter=true
```

```
1   hive优化
2   矢量化查询：开启了：Hive每个批次读取1024条，处理1024条
3   set hive.vectorized.execution.enabled = true;
4   set hive.vectorized.execution.reduce.enabled = true;
5
6   零拷贝：
7   功能：数据在操作系统中从内存中不用经过多次拷贝直接读取
```

```
 8   不开启：必须在内存中经过多次交换才能读取到数据
 9   开启了：数据可以直接从内存中读取
10   set hive.exec.orc.zerocopy=true;
11
12   关联优化器
13   功能：Hive在解析SQL语句，转换为MapReduce时候，可以将相关联的部分合并在一起执行
14   自动判断所有执行过程语法数是否有重合的过程，放在一起执行
15   set hive.optimize.correlation=true;
16
17   动态生成分区的线程数
18   说明：在执行动态分区过程中，可以运行多少个线程来生产分区数据，线程数量越多，执行效率越高，前提
     是有资源
19   hive.load.dynamic.partitions.thread  默认值为 15
```

## hive建表

```
 1   hive建表语句
 2   create [external | temporary] table [dbname.]tbname(
 3           colName1 type1 comment '',
 4           colName2 type1 comment '',
 5           colName3 type1 comment '',
 6           ……
 7           colNameN type1 comment ''
 8
 9   ) comment '表的注释'
10   partitioned by (col type)
11   clustered by (col) [sorted by col desc] into N buckets
12   row format delimited fields terminated by
13   stored as orc [tblpropertied ("orc.compress"="[ZLIB|SNAPPY]")]
14   location '指定路径'
15
16   jdbc:mysql://hadoop01:3306/yipin?useUnicode=true&characterEncoding=UTF-
     8&autoReconnect=true'
```

## mapreduce内存溢出优化

```
 1   MapReduce的内存配置
 2   map阶段:mapreduce.map.memory.mb ： 每一个map默认申请的内存大小，默认值为 0  表示自动获取；
     mapreduce.map.java.opts ：  每一个map的jvm的内存大小；
 3   注意：   mapreduce.map.java.opts一定要小于mapreduce.map.memory.mb；
 4
```

```
5   reduce阶段:mapreduce.reduce.memory.mb : 每一个reduce默认申请的内存大小，默认值为 0   表示自动
    获取；mapreduce.reduce.java.opts :   每一个reduce的jvm的内存大小；
6   注意：    mapreduce.reduce.java.opts 一定要小于mapreduce.reduce.memory.mb；
7
8   注意：MR中所有的内存配置，都不能大于nodemanager的内存大小
9
10  jvm-java的格式为：-Xmx4096m
11  其他：4096
```

## 压缩配置

```
1   yarn配置：   在CM中直接配置
2   mapreduce.map.output.compress   是否开启map端压缩配置       默认开启的
3   mapreduce.map.output.compress.codec   map端采用何种压缩方案
4        建议配置为：org.apache.hadoop.io.compress.SnappyCodec
5
6   mapreduce.output.fileoutputformat.compress 是否开启reduce端压缩配置   默认不开启的 对最终结
    果有影响
7   mapreduce.output.fileoutputformat.compress.codec reduce端需要采用何种压缩操作
8        建议配置为：org.apache.hadoop.io.compress.SnappyCodec
9
10  mapreduce.output.fileoutputformat.compress.type 采用压缩的方式
11        建议：block 块压缩
12
13  hive配置：   此配置需要在会话中执行 MR开了压缩，开启以下才会进行压缩
14  set hive.exec.compress.intermediate=true; 开启中间结果的压缩
15  set hive.exec.compress.output=true; 是否开启最终结果压缩
16
```