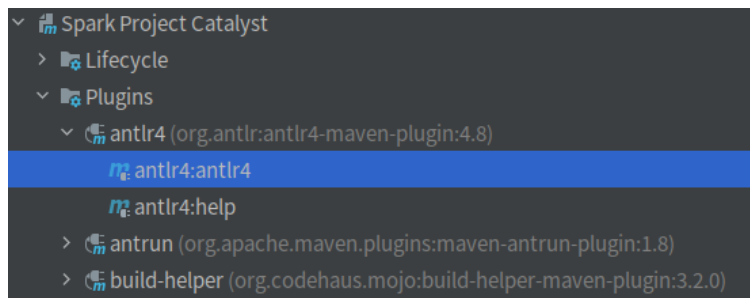


1.在SqlBase.g4中添加语法规则，路径为sql/catalyst/src/main/antlr4/org/apache/spark/sql/catalyst/parser/SqlBase.g4

```
1 statement
2 | COMPACT TABLE target=tableIdentifier partitionSpec?
3   (INTO fileNum=INTEGER_VALUE FILES)?                                #compactTable
4 ansiNonReserved
5 | FILES
6 nonReserved
7 | FILES
8 //--SPARK-KEYWORD-LIST-START
9 FILES: 'FILES';
```

2.执行antlr4:antlr4插件，自动生成代码



3. SparkSqlParser.scala中添加visitCompactTable方法：

```
1 override def visitCompactTable(ctx: CompactTableContext): LogicalPlan = withOrigin(ctx) {
2   val table: TableIdentifier = visitTableIdentifier(ctx.tableIdentifier())
3   // 解析获得文件数
4   val fileNum: Option[Int] = if (ctx.INTEGER_VALUE() != null) {
5     Some(ctx.INTEGER_VALUE().getText.toInt)
6   } else {
7     None
8   }
9   // 解析获得partitionSpec, 格式partition(key1=value1,key2=value2)
10  val partition: Option[String] = if (ctx.partitionSpec() != null) {
11    Some(ctx.partitionSpec().getText)
12  } else {
13    None
14  }
15  CompactTableCommand(table, fileNum, partition);
16 }
```

4.添加 CompactTableCommand 类的实现 在org.apache.spark.sql.execution.command包下新建CompactTableCommand类

```
1 package org.apache.spark.sql.execution.command
2
3 import org.apache.spark.sql.{Row, SaveMode, SparkSession}
4 import org.apache.spark.sql.catalyst.TableIdentifier
5 import org.apache.spark.sql.catalyst.expressions.{Attribute, AttributeReference}
6 import org.apache.spark.sql.types.StringType
7
8 case class CompactTableCommand(table: TableIdentifier,
9                                fileNum: Option[Int],
10                                partitionSpec: Option[String]) extends LeafRunnableCommand {
11
12   private val defaultSize = 128 * 1024 * 1024
```

```

13
14 override def output: Seq[Attribute] = Seq(
15     AttributeReference("compact", StringType, nullable = false)()
16 )
17
18 override def run(sparkSession: SparkSession): Seq[Row] = {
19     // 设置当前数据库
20     sparkSession.catalog.setCurrentDatabase(table.database.getOrElse("default"))
21     // 临时表格式: curTable_timestamp
22     val tempTableName = "`" + table.identifier + "_" + System.currentTimeMillis() + "`"
23
24     val originDataFrame = sparkSession.table(table.identifier)
25     // 计算分区数, 如果fileNum有效, 则为fileNum, 否则使用默认分区大小128M计算分区数
26     val partitions = fileNum match {
27         case Some(num) => num
28         case None => (sparkSession.sessionState
29             .executePlan(originDataFrame.queryExecution.logical)
30             .optimizedPlan.stats.sizeInBytes / defaultSize).toInt + 1
31     }
32
33     if (partitionSpec.nonEmpty) {
34         // 如果不更改这个参数, 使用默认的static, 在动态插入时, 不管插入的分区是否存在, 都会导致所有的分区被覆盖, 数据丢失
35         sparkSession.conf.set("spark.sql.sources.partitionOverwriteMode", "dynamic")
36         // Dynamic partition strict mode requires at least one static partition column.
37         // To turn this off set hive.exec.dynamic.partition.mode=nonstrict
38         sparkSession.conf.set("hive.exec.dynamic.partition.mode", "nonstrict")
39
40         // 当partitionSpec有值时, partition(key1=value1,key2=value2)转换为key1=value1 AND key2=value2的格式
41         val conditionExpr = partitionSpec.get.trim
42             .stripPrefix("partition(").dropRight(1)
43             .replace(",", " AND ")
44
45         // 设置分区数及where条件将数据写入临时表中
46         originDataFrame
47             .where(conditionExpr)
48             .repartition(partitions)
49             .write.mode(SaveMode.Overwrite)
50             .saveAsTable(tempTableName)
51
52         // 将临时表中的数据重新插入到原表中
53         sparkSession
54             .table(tempTableName)
55             .write.mode(SaveMode.Overwrite)
56             .insertInto(table.identifier)
57     } else {
58         // 当partitionSpec不存在时, 设置分区数将数据读取到临时表中
59         originDataFrame
60             .repartition(partitions)
61             .write
62             .mode(SaveMode.Overwrite)
63             .saveAsTable(tempTableName)
64
65         // 读取临时表, 对原表进行覆盖

```

```
66     sparkSession.table(tempTableName)
67     .write
68     .mode(SaveMode.Overwrite)
69     .saveAsTable(table.identifier)
70 }
71
72 // 删除临时表，此处不删除用于查看分区效果
73 // sparkSession.sql(s"DROP TABLE ${tempTableName}")
74
75 Seq(Row(s"compact table ${table.identifier} finished."))
76 }
77 }
```

5, 编译

```
1 build/mvn clean package -DskipTests -Phive -Phive-thriftserver
```

6, 运行

```
1 bin/spark-sql, COMPACT TABLE test INTO fileNum=100;
```

参考: <https://gitee.com/joyjung/bigdata-homework-compact>