

Spark SQL 1.1.0和Hive的兼容说明

Spark SQL也是可以部署在当前的Hive warehouse。

Spark SQL 1.1.0的 Thrift JDBC server 被设计成兼容当前的Hive数据仓库。你不需要修改你的Hive元数据，或者是改变表的数据存放目录以及分区。



微信扫一扫，加关注
即可及时了解Spark、Hadoop或者Hbase
等相关的文章
欢迎关注微信公共帐号：iteblog_hadoop

过往记忆博客 (<http://www.iteblog.com>)
专注于Hadoop、Spark、Flume、Hbase等
技术的博客，欢迎关注。

Hadoop、Hive、Hbase、Flume等交流群：138615359和149892483

如果想及时了解Spark、Hadoop或者Hbase相关的文章，欢迎关注微信公共帐号：iteblog_hadoop

以下列出来的是当前Spark SQL (1.1.0) 对Hive特性的支持

一、Hive查询语句，包括：

- SELECT
- GROUP BY
- ORDER BY
- CLUSTER BY
- SORT BY

二、所有的Hive操作：

- Relational operators (=, <=>, ==, <>, <, >, >=, <=, etc)
- Arithmetic operators (+, -, *, /, %, etc)
- Logical operators (AND, &&, OR, ||, etc)
- Complex type constructors
- Mathematical functions (sign, ln, cos, etc)
- String functions (instr, length, printf, etc)
- User defined functions (UDF)
- User defined aggregation functions (UDAF)
- User defined serialization formats (SerDe's)
- Joins
- JOIN
- {LEFT | RIGHT | FULL} OUTER JOIN
- LEFT SEMI JOIN
- CROSS JOIN

- Unions
- Sub queries
- SELECT col FROM (SELECT a + b AS col from t1) t2
- Sampling
- Explain
- Partitioned tables

三、所有的Hive DDL函数，包括：

- CREATE TABLE
- CREATE TABLE AS SELECT
- ALTER TABLE

四、大部分的Hive数据类型，包括：

- TINYINT
- SMALLINT
- INT
- BIGINT
- BOOLEAN
- FLOAT
- DOUBLE
- STRING
- BINARY
- TIMESTAMP
- ARRAY<>
- MAP<>
- STRUCT<>

以下是不支持Hive的：

一、Major Hive Features

带有buckets的table，目前Spark SQL还不支持。

二、Esoteric Hive Features

- 1、带有不同输入格式的分区表，在Spark SQL中，所有的表分区的输入格式必须相同；
- 2、不等值的outer join ("key

三、Hive的输入 输出格式

1、CLI的文件格式: 对于返回到CLI界面的结果信息，Spark SQL目前只支持TextOutputFormat

2、Hadoop archive

四、Hive优化

有一大部分的Hive优化在当前的Spark SQL是不支持的，这些优化中（包括了Indexs）在Spark SQL是不重要的，因为Spark

SQL是基于内存的计算模型。其他的优化将会在Spark SQL以后的版本得到支持。（下面我就不翻译了，太要时间了。下面的英文都简单易懂）

- 1、Block level bitmap indexes and virtual columns (used to build indexes)
- 2、Automatically convert a join to map join: For joining a large table with multiple small tables, Hive automatically converts the join into a map join. We are adding this auto conversion in the next release.
- 3、Automatically determine the number of reducers for joins and groupbys: Currently in Spark SQL, you need to control the degree of parallelism post-shuffle using "SET spark.sql.shuffle.partitions=[num_tasks];". We are going to add auto-setting of parallelism in the next release.
- 4、Meta-data only query: For queries that can be answered by using only meta data, Spark SQL still launches tasks to compute the result.
- 5、Skew data flag: Spark SQL does not follow the skew data flags in Hive.
- 6、STREAMTABLE hint in join: Spark SQL does not follow the STREAMTABLE hint.
- 7、Merge multiple small files for query results: if the result output contains multiple small files, Hive can optionally merge the small files into fewer large files to avoid overflowing the HDFS metadata. Spark SQL does not support that.

仔细看下就知道这些其实都是Shark对Hive的兼容，[《Shark对Hive的兼容性总结》](#)，其实就是把Shark迁移进Spark SQL了。期待Spark SQL更强大的功能了。

本博客文章除特别声明，全部都是原创！
原创文章版权归过往记忆大数据（[过往记忆](#)）所有，未经许可不得转载。
本文链接：[【】（）](#)