

Hash join 性能测试报告

修订历史

版本	修订日期	修订描述	作者	备注
Cedar	2017-09-28	Hash Join性能测试报告	胡爽	无

测试环境

使用6台虚拟机组成的集群作为测试环境，每台虚拟机的配置相同，包括4核1.2GHz主频CPU、100GB内存、3000GB磁盘，虚拟机上安装了CentOS release 6.5系统，相互之间通过千兆以太网连接。其中三台虚拟机环境，每台部署一个数据库集群，包括UpdateServer, RootServer,ChunkServer以及MergeServer。剩下三台每个环境搭建一个MergeServer。

测试案例一

1.案例描述

测试join的深度对性能的影响。

2.测试方法

- 2张500w的表inner join、3张500w的表inner join、4张500w的表inner join、5张500w的表inner join、6张500w的表inner join。
- 2张1000w的表inner join、3张1000w的表inner join、4张1000w的表inner join、5张1000w的表inner join、6张1000w的表inner join

3.测试结果

500w数据量的测试结果如下图：



1000w数据量的测试结果如下图：



测试结果：随着join深度增加，Merge Join和Bloomfilter Join所费时间大幅度上升，而Hash Join所耗时间缓慢上升，从图一二可以看出性能最好。

测试案例二

1.案例描述

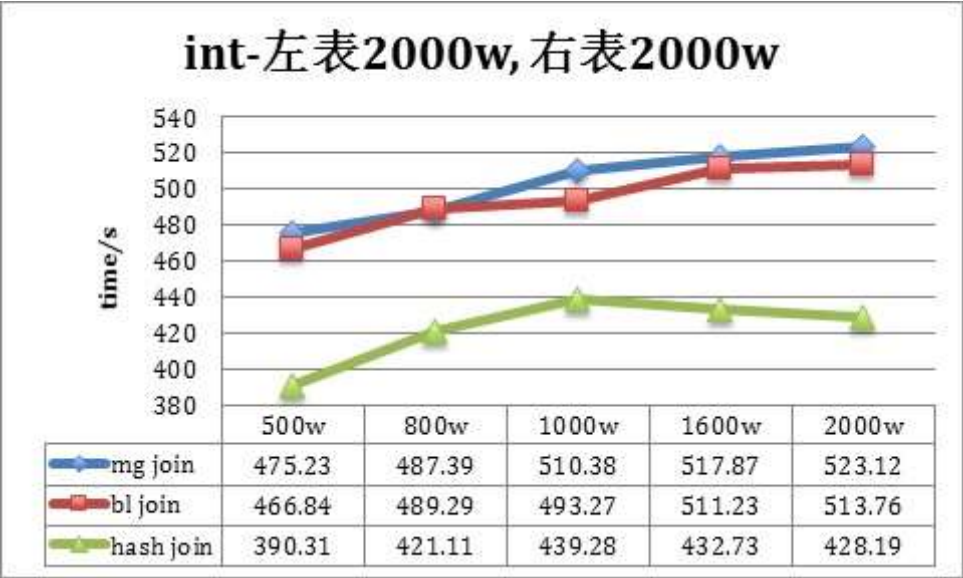
测试Hash Join在选择率为10%的情况下，非重复值个数对性能的影响（连接属性为int类型）。

2.测试方法

左右表都为2000万条数据，左表的非重复值个数为（500w,800w,1000w,1600w,2000w）。

3.测试结果

测试结果如下图：



测试结果：Join表数目为2000w时，随着左表非重复数的增加，Hash Join所耗时间远低于merge join 和bloom filter join，从图三看出hash join性能远好于Bloom filter Join 和Merge Join。

测试案例三

1.案例描述

测试Hash Join在非重复值个数为500万下，选择率对性能的影响（连接属性为int类型）。

2.测试方法

左表1000万条数据（非重复数500万，选择率0.1%，1%，10%，30%，60%，90%）

3.测试结果



在非重复值个数为500w时，Join左右表数目为1000万时，Hash Join所耗时间低于merge join和bloom filter join，从图四看出hash join性能远好于Bloom filter Join 和Merge Join

测试案例四

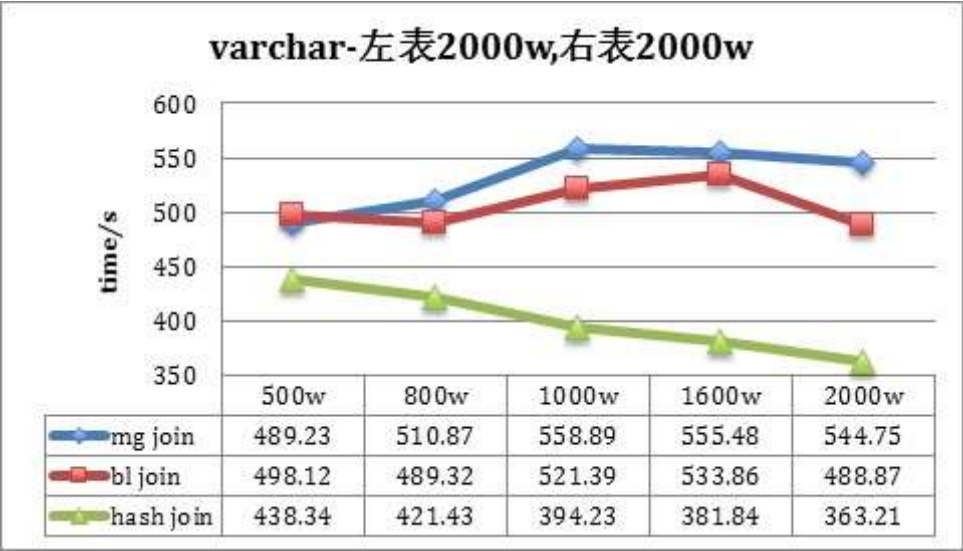
1.案例描述

测试Hash Join在选择率为10%的情况下，非重复值个数对性能的影响（连接属性为varchar类型）。

2.测试方法

左右表2000万条数据（非重复数500w,800w,1000w,1600w,2000w，连接属性为varchar类型）

3.测试结果



在将左表的非重复值个数增加到（500w,800w,1000w,1600w,2000w）时，Join表数目为2000w时，Hash Join性能远好于Bloom filter Join 和Merge Join。即从图看出，当连接属性为varchar类型时，相较于Bloom filter Join 和Merge Join，Hash Join性能更优。

测试案例五

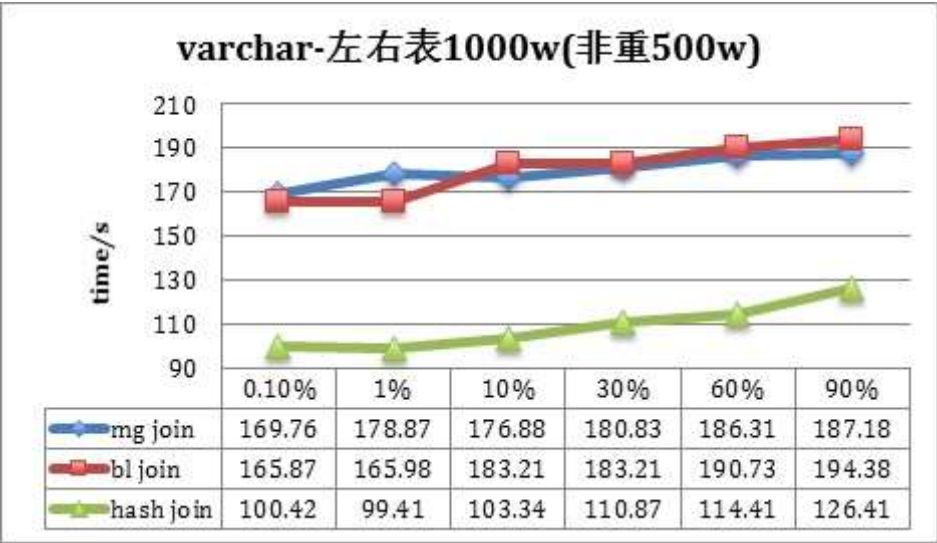
1.案例描述

测试Hash Join在非重复值个数为500万下，选择率对性能的影响（连接属性为varchar类型）。

2.测试方法

左右表1000万条数据（非重复数500万，选择率0.1%，1%，10%，30%，60%，90%）

3.测试结果



连接属性为varchar类型时，随着选择率的升高，相较于Bloom filter Join 和Merge Join，Hash Join性能更优。从图六看出，Hash Join性能远好于Bloom filter Join 和Merge Join。

测试结论

由上述实验分析可知，在大数据量的情况下，hash Join与merge Join，bloom filter join相比性能有明显的优势。