

Bethune 告警 & MR问题作业优化建议

- 1. map数目过多
 - 问题描述
 - 可能原因
 - 优化建议
- 2. reduce 数据倾斜
 - 问题描述
 - 可能原因
 - 优化建议
- 3. job 总HDFS输出数据量 or 作业总shuffle数据量 过大
 - 问题描述
 - 可能原因
 - 优化建议
- 4. map执行时间倾斜 & map小文件
 - 问题描述
 - 可能原因
 - 优化建议
- 5. Map/Reduce/job GC严重
 - 问题描述
 - 可能原因
 - 优化建议
- 6. map/reduce平均运行时间过长
 - 问题描述
 - 可能原因
 - 优化建议
- 7.job资源申请比例合理性 or Map/Reduce 资源调度响应时间
 - 问题描述
 - 可能原因
 - 优化建议
- 8.job长时间运行
 - 问题描述
 - 可能原因
 - 优化建议
- 9.存在OOM失败task & AM task心跳超时
 - 问题描述
 - 可能原因
 - 优化建议
- 备注

1. map数目过多

问题描述

单个mapreduce作业 map 数目超过20w

可能原因

1. 查询范围过大造成输入数据量大；
2. 分片设置过小；
3. 用户未使用默认的CombineHiveInputFormat

优化建议

当sql的输入数据量太大，导致map task个数特别多，虽然每个task执行时间都不长，但是由于计算资源有限，在资源紧张的时候，一个作业内的多个task只能分批串行执行，导致资源调度开销成为任务执行时间过长的主要因素，花在资源等待上的时间长达几十分钟乃至几小时。另外map task过多也会导致reduce阶段 shuffle时间变长。

1. 减少查询范围以及输入数据量，合理设置sql查询的分区范围，尽量避免全表扫描，考虑生成一些增量的中间表来替代；
2. 通过加大配置项mapreduce.input.fileinputformat.split.maxsize(默认256000000，单位字节)减少分片数，效果是单个map task处理的数据量增加，执行时间可能变长，但是可以降低整个作业资源开销以及调度开销；
3. 如无特殊原因，尽量使用默认的CombineHiveInputFormat，可以合并输入小文件

2. reduce 数据倾斜

问题描述

reduce完成度大于90%，最大reduce执行时间大于30min，task最大shuffle数据量大于6G且大于平均值6倍

可能原因

数据倾斜是一个比较偏业务层并且用户经常遇到的问题，这里主要讲解一下背后的原理以及常见的优化思路，具体如何彻底解决或者优化，还是需要业务结合具体数据和业务场景考虑

什么是数据倾斜

由于数据分布不均匀，大量数据集中在某些特定key上，造成数据热点。

数据倾斜的主要表现

根据mapreduce架构的原理，会按照key把不同的数据hash到不同的reduce task，当存在数据热点时，就会导致某些reduce task处理的数据量远远超过其他task（几倍乃至数十倍），最终表现为少量reduce task执行长尾，任务整体进度长时间卡在99%或者100%。

容易导致数据倾斜的场景

Join/GroupBy/CountDistinct，在存在热点key（例如某个字段存在大量空值）的时候，都会导致一个或少数reduce task处理的数据量远超其他task

优化建议

关于数据倾斜和优化 详见 数据架构组 stack文章 [Hive常见问题和优化手段](#)

1. 过滤掉不符合预期的热点key，例如由于日志信息丢失导致某个字段产生大量空值
2. 加入随机因素，打散热点key
3. 使用map join解决小表关联大表造成的数据倾斜问题

map join是指将做连接的小表全量数据分发到作业的map端进行join，从而避免reduce task产生数据倾斜；公司内map join优化默认已打开（hive.auto.convert.join=true），这个配置跟hive.optimize.skewjoin有冲突，请保证二者只开一个即可；map join需要在内存中加载全部小表数据，容易导致map端OOM，hive.mapjoin.smalltable.filesize这个参数用于设置小表的大小，默认25000000（25M），当小表数据量超过这个大小时，不会走map join优化逻辑，不建议用户把这个参数设置过大

3. job 总HDFS输出数据量 or 作业总shuffle数据量 过大

问题描述

单个作业总HDFS写入量大于20T 或者 单个作业总shuffle数据量大于20T

可能原因

1. 输入数据量过大，查询范围比较大
2. 业务逻辑存在导致数据膨胀的操作

优化建议

1. 减少输入数据量，分批查询；
2. 检查业务逻辑，确认是否存在导致数据膨胀的操作以及是否存在优化空间，减少不必要的输出

4. map执行时间倾斜 & map小文件

问题描述

map完成度大于90%，map最长执行时间大于30min且大于平均值6倍（map task长尾）

单个map最大读操作数大于10000（小文件数过多）

可能原因

出现这种现象的原因通常是：

1. hive查询分多轮stage进行，某些作业在运行过程中产生了大量临时小文件（临时文件数 = task个数 * 动态分区数），在下一轮stage中，由于map task是按照固定数据量大小（默认256M）进行分片，如果生成的文件只有KB级别，则单个task需要处理成千上万个文件，这期间涉及到大量hdfs操作，很容易受到hdfs波动影响导致执行时间拉长。
2. 单个map输入数据量过大，存在复杂计算逻辑 task gc比较严重

优化建议

1. 查看上一轮stage作业是否存在reduce，如果有reduce task，则小文件是reduce生成的，如果单个reduce task执行时间不是特别大，可以适当控制reduce最大并发(hive.exec.reducers.max，默认5120，建议设置为2560/1280/640等)；如果上一轮stage没有reduce，则小文件是map生成的，需要加大split size减少map task (mapreduce.input.fileinputformat.split.maxsize，默认256000000，建议可以设置到1024000000)；
2. 在优化手段1的基础上，还可以使用数据架构组定制开发的根据文件数分片的功能(hadoop默认是按照文件大小分片)，限制单个task处理的文件大小(set mapreduce.split.by.block.num.enable = true; set mapreduce.split.block.number.threshold = 500;)

5. Map/Reduce/job GC严重

问题描述

Map task、Rduce task、job 整体 执行期间gc严重将导致作业执行变慢，task失败，内存超限，甚至导致作业直接失败

可能原因

1. 作业task处理数据量较大
2. 作业存在大量复杂计算逻辑

优化建议

1.加大内存：mapTask gc告警设置mapreduce.map.memory.mb（默认3072），reduceTask gc告警可以设置mapreduce.reduce.memory.mb（默认4096），建议按照512的幅度增加，合理使用避免浪费；

2.如果sql中有join和group by操作，可以调整参数缩小内存buffer检查间隔：

```
set hive.mapjoin.check.memory.rows=10000;
```

```
set hive.groupby.mapaggr.checkinterval=5000;
```

```
set hive.map.aggr.hash.percentmemory=0.3;
```

```
set hive.mapjoin.followby.map.aggr.hash.percentmemory=0.1;
```

```
set hive.map.aggr.hash.force.flush.memory.threshold=0.7;
```

```
set hive.map.aggr.hash.min.reduction=0.3;
```

3.可以选择关闭GBY的map端优化来争取节约内存hive.map.aggr=false;

6. map/reduce平均运行时间过长

问题描述

map平均执行时间大于1h

reduce平均执行时间大于2h

可能原因

1. map并发度不足单个task处理的数据量过大，导致map平均运行时间长；
2. reduce并发度不足单个task处理的数据量过大，导致reduce平均运行时间长；
3. map内存不足，导致map task频繁gc任务运行时间拉长；
4. reduce内存不足，导致reduce task频繁gc任务运行时间拉长；
5. 存在join操作，且record记录较大，读取的数据量大，导致任务频繁gc；
6. 存在group by操作，且存在数据热点的情况；
7. udf或者transform脚本存在耗时操作；

优化建议

这种情况一般是作业整体输入数据量不大，但是由于计算逻辑复杂导致作业执行时间特别长，作业的map/reduce个数通常在个位数或者十位数，这种情况下需要反过来通过增加task并发度，减少单个task处理的数据量来加快任务运行速度

1. map数太少，单个map执行时间过长，可以通过减小mapreduce.input.fileinputformat.split.maxsize(默认256000000，单位字节)增大map并发；
2. reduce数太少，单个reduce执行时间过长，可以通过减小hive.exec.reducers.bytes.per.reducer(默认1073741824，单位字节)来增大reduce并发；
3. map task 存在频繁gc或者OOM问题，增加mapreduce.map.memory.mb（默认3072，单位mb）；
4. reduce task存在频繁gc或者OOM问题，增加mapreduce.reduce.memory.mb（默认4096，单位mb）；
5. 存在join操作，且record记录较大，读取的数据量大，调小join缓存buffer大小，避免因为内存不足导致GC严重（如果调整后问题没有解决可以继续减小阈值），参数调整后可能导致其他原本没有gc问题的task运行时间拉长 set hive.join.emit.interval=500; set hive.join.cache.size=5000;

6. 若有group by操作，且存在数据热点的情况，调小group by的检查阈值，内存buffer大小，记录聚合触发阈值，避免因为group by占用内存过多，导致GC严重，运行延迟（如果调整后问题没有解决可以继续减小阈值），参数调整后可能导致其他原本没有gc问题的task运行时间拉长 set hive.groupby.mapaggr.checkinterval=5000; set hive.map.aggr.hash.percentmemory=0.3; set hive.map.aggr.hash.min.reduction=0.3; set hive.map.aggr.hash.force.flush.memory.threshold=0.7; set hive.mapjoin.followby.map.aggr.hash.percentmemory=0.1;
7. 联系udf或transform脚本owner确认是否存在耗时操作；

7.job资源申请比例合理性 or Map/Reduce 资源调度响应时间

问题描述

内存占用超1T，job运行时长超过2h.(内存申请值/cpu申请值)>5

作业map、reduce task资源满足总耗时过长 作业执行因为队列资源不足执行时间变长

可能原因

作业出现这类问题一般是由于上面某几个问题导致的，另外用户可能对资源正常申请量概念不明确，存在将其他作业的优化参数直接拷贝到新作业执行的现象，导致资源不合理使用，影响部门队列其他作业的资源使用

优化建议

1. 在**bethune** 对应作业页面查看作业是否存在其他健康问题
2. 用户可以通过 [Yarn分析看板](#) [分析看板使用说明](#) 来查看自己部门的资源使用情况，查看部门队列资源是否紧张，资源占用大头的用户和DAG
3. 如果还有问题，需要紧急调整资源和作业优先级，请拉上YARN oncall同学，本组负责人，详细说明需要调整优先级的原因，为什么要优先于其他人跑，需要负责人批准

8.job长时间运行

问题描述

单个作业运行时间超过18h

可能原因

可能造成作业整体运行时间过长的原因通常不是单一因素，可能包括计算资源不足，输入数据量过大，数据倾斜导致长尾task，gc频繁，并发不足等多种原因，需结合其他报警内容一起分析并给出优化建议

优化建议

1. 用户可以通过 [Yarn分析看板](#) [分析看板使用说明](#) 来查看自己部门的资源使用情况，查看部门队列资源是否紧张，资源占用大头的用户和DAG
2. 缩小查询范围，减少输入数据量，分批查询；
3. 结合其他问题，综合给出优化建议

9.存在OOM失败task & AM task心跳超时

问题描述

作业存在由于OutOfMemory，gc超时，心跳超时等原因失败的task

可能原因

1. map/reduce 内存不足
2. 存在join操作，且record记录较大，读取的数据量大，导致任务频繁gc；
3. 存在group by操作，且存在数据热点的情况，导致任务频繁gc；

优化建议

1. 对于作业只有极个别task由于以上原因失败的情况，数据架构组开发了针对这类task失败自动调整内存重试的功能，可以先不用调整作业内存参数
2. map task 存在频繁gc或者OOM问题，增加mapreduce.map.memory.mb（默认3072，单位mb）建议以512mb为单位进行调节，避免资源浪费；
3. reduce task存在频繁gc或者OOM问题，增加mapreduce.reduce.memory.mb（默认4096，单位mb）建议以512mb为单位进行调节，避免资源浪费；

备注

YARN & MR 问题可以优先查看 [YARN用户常见 FAQ](#) 来自助查询

**如果作业有其他问题和疑问，请通过数据服务热线提问

