如何解读四种 Oracle 优化器直方图

作者:包光磊(Todd Bao)

http://www.shoug.info/archives/todd-bao.html



本文的适应读者: 初步理解直方图对优化器的作用。

预计阅读时间: 10分钟。

数据库一共有四种直方图,它们是:频度型直方、顶端频度型直方、高度平衡型直方与混合型直方。

示例中的 DEPARTMENT_ID 字段相异值为 10、20、30、40、50、60、70、80、90、100、110, 一共十一个部门。

-- Type 1: 频度型直方

测试 (表内行数 82):

```
begin
 dbms_stats.gather_table_stats(
  'HR',
  'EMPLOYEES',
  method opt=>'for columns department id size 100',
  estimate percent=>dbms stats.auto sample size);
end;
/
define owner='hr'
define tab='employees'
define col='department id'
select column_name, num_distinct, num_buckets, histogram
from dba tab col statistics
where owner=upper('&owner') and table_name=upper('&tab') and
column name=upper('&col')
select column name, endpoint number, endpoint value from
dba tab histograms
where owner=upper('&owner') and table_name=upper('&tab') and
column name=upper('&col')
order by 1, 2
undefine owner
undefine tab
undefine col
```

字段名字段相异值直方类型COLUMN_NAMENUM_DISTINCT NUM_BUCKETSHISTOGRAM

DEPARTMENT ID 11

由于计划收集直方数(100)大于实际字段相异值数(11)所以实际直方数等于字段相异值数。在这种情况下,每个相异值分配一个直方(一柱)统计其所占行数,统计的结果在直方收集时是精确到每一个相异值的,其他的类型不具备这一特点。

11

FREQUENCY

字段名	行数	字段值
COLUMN_NAME	ENDPOINT_NUMBER	ENDPOINT_VALUE
DEPARTMENT_ID	1	<mark>10</mark>
DEPARTMENT_ID	<mark>3</mark>	<mark>20</mark>
DEPARTMENT_ID	5	30
DEPARTMENT_ID	6	40
DEPARTMENT_ID	31	50
DEPARTMENT_ID	36	60
DEPARTMENT_ID	37	70
DEPARTMENT_ID	71	80
DEPARTMENT_ID	74	90
DEPARTMENT_ID	80	100
DEPARTMENT_ID	82	110

由于行数是累积的,所以 DEPARTMENT_ID 等于 10 的行数是 1、为 20 的是 (3-1) 等于 2、以此类推。

-- Type 2:顶端频度直方

测试(表内行数82):

```
begin
  dbms_stats.gather_table_stats(
   'HR','EMPLOYEES',
   method_opt=>'for columns department_id size 5',
   estimate_percent=>dbms_stats.auto_sample_size);
end;
/
define owner='hr'
define tab='employees'
define col='department_id'
```

```
select column name, num distinct, num buckets, histogram
from dba tab col statistics
where owner=upper('&owner') and table name=upper('&tab') and
column name=upper('&col')
select column_name, endpoint_number, endpoint_value
from dba tab histograms
where owner=upper('&owner') and table name=upper('&tab') and
column name=upper('&col')
order by 1, 2
undefine owner
undefine tab
undefine col
字段名
                       字段相异值
                                    直方数
                                                直方类型
COLUMN NAME
                      NUM DISTINCT NUM BUCKETS HISTOGRAM
                                                TOP-FREQUENCY
DEPARTMENT ID
                       11
                                    5
字段名
                  行数
                                  字段值
COLUMN NAME
               ENDPOINT NUMBER ENDPOINT VALUE
DEPARTMENT ID
               1
DEPARTMENT ID
               26
DEPARTMENT ID
               60
DEPARTMENT ID
               66
```

在 11 个部门里选择 5 个值作为直方的统计有三个原因: a. 10 和 110 是两个极值,必须出现; b. 只能再加三个部门是因为直方数限制为 5; c. 50、80 和 100 选择这三个部门的原因是属于这些部门的行数总和最接近内部阈值:1-(1/2) 为 100 数)与行基数的积,本例中为: 1-(1/5),即 100 80%的行数。

DEPARTMENT ID

67

除两个极值以外,大约占总行数 20%左右的 DEPARTMENT_ID 字段值在此例中没有展现,它们被优化器认为选择度偏高。

与频度直方比较很容易发现这种直方把"矮楼"(选择度高)铲平了,只对"高楼"(选择度低)感兴趣,自带强拆队,有做房地产商的潜质。

×注意: 是否测试不成功? 是否收集不到这种直方? 是否收到的是 HEIGHT BALANCED-高度平衡型的? 那么请期待下个版本的 Oracle 数据库。如果下个版本



还是没有这种直方,请期待下下个版本的 Oracle 的数据库。

-- Type 3: 高度平衡型直方

一旦刚才提到的房地产商出现,我们就可以和高度平衡型直方说再见了。为了忘却的纪念让我再介绍一下它。

测试 (表内行数 107)

省略收集…

字段名	字段相异值	直方数	直方类型
COLUMN_NAME	NUM_DISTINCT	NUM_BUCKETS	HISTOGRAM
DEPARTMENT_ID	11	5	HEIGHT BALANCED
字段名	行单位数	字段值	VALUE
COLUMN_NAME	ENDPOINT_NUMBI	ER ENDPOINT_	
DEPARTMENT_ID DEPARTMENT_ID DEPARTMENT_ID DEPARTMENT_ID	0 2 4 5	10 50 80 110	

查看此类直方必须了解表的行基数,此例为107。

一个行单位等于行数处以直方数=107/5 约 21.

所以,估计各行数是

10 号部门: 21*<mark>0</mark> = 大约 0, 当然不可能是 0 行, 极少的意思

20-50 号部门: 21*(2-0) = 大约 42 60-80 号部门: 21*(4-2) = 大约 42 90-110 号部门: 21*(5-4) = 大约 21

与实际比较一下:	与实	际	比车	☆—	下	:
----------	----	---	----	----	---	---

SYS@fmw//Scripts > select count(*) from hr. employees where department_id between 11 and 50;

COUNT(*)

54

SYS@fmw//Scripts> select count(*) from hr.employees where department_id between 51 and 80;

COUNT(*)

40

SYS@fmw//Scripts> select count(*) from hr.employees where department_id between 81 and 110;

COUNT(*)

11

-- Type 4: 混合型直方

Oracle 多加一个字段来描述这种直方: ENDPOINT_REPEAT_COUNT。

测试 (表内行数 82)

字段名	字段相异值	直方数	直方类型
COLUMN_NAME	NUM_DISTINCT	NUM_BUCKETS	HISTOGRAM
DEPARTMENT_ID	11	4	HYBRID

字段名	VME	行数	MIMDED	字段值	WALLE	字段值重		COLINT
COLUMN_NA	AME 	ENDPOINI	_NUMBER	ENDPOINI_	_VALUE	ENDPOINT_	_KEPEAT_ 	_COUNT
DEPARTMEN	NT_ID	1		10		1		
DEPARTMEN	NT_ID	<mark>31</mark>		50		<mark>25</mark>		
DEPARTMEN	NT_ID	71		80		34		
DEPARTMEN	NT_ID	<mark>82</mark>		110		2		
10 早並7	的怎粉	4						
10 号部门		1;	00/=		□ 2	a - 1 -		
20-50 号音				,其中 50				
60-80 号音	邓门的行数	数: <mark>71</mark> -	- <mark>31</mark> =40 行	F,其中 80	号部门] 34 行;		
90-110 号	部门的行	·数: 82-	- <mark>71</mark> =9 行	其中 110	号部门] 2 行。		

什么,没有这种直方?

← 请参考 Type 2.

四种直方的解读至此介绍完毕。

参考资料:

www.oracle-base.com/blog/2012/10/06/oracle-openworld-2012-day-5/

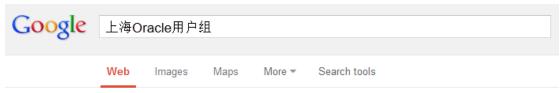
郑重声明:本文仅代表个人观点,不代表 Oracle 公司。

包光磊

新浪微博 http://weibo.com/toddbao

我主持的国内 ADF/SOA 论坛: http://www.databi.cn/forum-57-1.html

SHOUG 的网址: www.shoug.info 找我们:



About 234,000 results (0.34 seconds)

<u>上海Oracle用户组| SHOUG,走近全系Oracle技术和数据库专家</u>

www.shoug.info/ Translate this page

5 days ago – SHOUG的全称是ShangHai Oracle Users Group,中文为**上海Oracle用户组。** SHOUG的成员仅仅局限于上海地区吗? 上海是国际化大都市,我们 ...

SHOUG成员| 上海Oracle用户组

www.shoug.info/archives/category/shoug成员 ▼ Translate this page 2 days ago – 在Oracle数据库领域具有比较扎实的理论基础、较丰富的实践经验、处理过较多的Oracle性能优化和数据恢复案例,乐于分享,学习Oracle技术。