

Multi column correlation

I have encountered this situation many times before: Cost based optimizer assumes no correlation between two columns (until 10g) and this has the effect of reducing cardinality of a row source erroneously. Consider following example.

```
create table t_vc as select mod(n, 100) n1, mod(n, 100) n2 , mod(n, 50) n3 ,
mod(n, 20) n4
from (select level n from dual connect by level <= 10001);
```

There is a strong correlation between n1 and n2 above. $N1 = N2$. There is some correlation among other columns. Let's collect statistics with histograms on all columns.

```
exec dbms_stats.gather_Table_stats( user, 'T_VC', estimate_percent => null,
method_opt => 'for all columns size 254');
```

```
explain plan for select count(*) from t_vc where n1=10 and n2=10;
```

```
-----
| Id | Operation          | Name | Rows | Bytes | Cost (%CPU)| Time     |
-----
|  0 | SELECT STATEMENT   |      |    1 |    6 |    9  (0)| 00:00:01 |
|  1 |  SORT AGGREGATE    |      |    1 |    6 |           |         |
|*  2 |   TABLE ACCESS FULL| T_VC |    1 |    6 |    9  (0)| 00:00:01 |
-----
```

Notice the # of rows column in the table above. It is just 1. But, there are 100 rows returned for that SQL.

CBO estimate is simple (example without considering histograms):

$$\begin{aligned} \# \text{ of rows} &\sim \text{total \# of rows} * (1/\text{NDV for } n1) * (1/\text{NDV for } n2) \\ &= 10000 * (1/100) * (1/100) = 1 \text{ row.} \end{aligned}$$

Far from the truth!

Extended stats

Oracle 11g introduces extended stats to relieve some pain. In 11g, an extended stats can be added between columns, enabling CBO to consider correlation between these column values.

```
SELECT dbms_stats.create_extended_stats(
    ownname => user, tabname => 'T_VC',
    extension => '(n1, n2)' ) AS n1_n2_correlation
```

FROM dual;

N1_n2_correlation

```
-----
SYS_STUBZH0IHA7K$KEBJVXO5LOHAS
```

Let's collect stats again on this table and check the SQL plan.

```
exec dbms_stats.gather_Table_stats( user, 'T_VC', estimate_percent => null,
method_opt => 'for all columns size 254');
```

```
explain plan for select count(*) from t_vc where n1=10 and n2=10;
```

```
-----
```

Id	Operation	Name	Rows	Bytes	Cost (%CPU)	Time
0	SELECT STATEMENT		100	1200	9 (0)	00:00:01
* 1	TABLE ACCESS FULL	T_VC	100	1200	9 (0)	00:00:01

```
-----
```

Estimated # of rows correctly at 100.

Under the wrap

Adding an extended stats adds a new virtual column to the table. Here is the line from sqltrace. Virtual column name is cryptic and seems to have been derived from table_name, column name combinations.

```
alter table "CBQT"."T_VC" add (SYS_STUBZH0IHA7K$KEBJVXO5LOHAS as
(sys_op_combined_hash(n1, n2)) virtual BY USER for statistics);
```


Cardinality estimates are exactly matching with reality. Cardinality calculations are quite important for performance.

In the next section, we will discuss this further ;-)