



## Cardinality Estimator

Herbert Albert  
Halbert@SolidQ.com

### Agenda

- Cardinality Estimation Basics
- Statistics Basics
  - Definition
  - Creating and Updating Statistics
  - Best Practices
- Common Cardinality Estimation Problems
- The New Cardinality Estimator
  - Changes
  - Regression Problems and possible Solutions

© 2014 SolidQ

2

## Cardinality Estimation

- Cardinality Estimation is the process of estimating the
  - number of rows processed by the different operators in a query plan
  - distribution of values
  - distinct value counts
  - duplicate counts
- Essential for
  - Query plan generation
    - Finding the right indexes, joins, groups,...
  - Query execution
    - Memory Grants
    - DOP
- Largely based on SQL Server 7.0
  - Fixes in QFEs under trace flag to prevent regressions
  - Some problems needed a major redesign

© 2014 SolidQ

3

## Calculation per Operator

```
Query 1: Query cost (relative to the batch): 100%
```

[illegible]

© 2014 SolidQ

4

# Agenda

- Cardinality Estimation Basics
- **Statistics Basics**
  - Definition
  - Creating and Updating Statistics
  - Best Practices
- Common Cardinality Estimation Problems
- The New Cardinality Estimator
  - Changes
  - Regression Problems and possible Solutions

© 2014 SolidQ

5

# Statistics

- Cardinality estimation is based on statistics
- In the absence of statistics or if statistics cannot be used for the estimation, SQL Server uses heuristics
  - Can lead to horrible plans
- Statistics are created for one or a combination of columns
- Statistics can be created and updated automatically and manually
  - Each index has a statistic

© 2014 SolidQ

6

## Statistics

- Each statistic consists of
  - Header
    - Rows sampled, rows in table, last update timestamp,...
  - Density vector
    - Measure for the uniqueness of a column
    - Calculated by  $1/(\text{number of distinct values})$
    - 1 x the number of rows returns the estimate for one value
    - Histogram is used when possible
  - Histogram
    - Data distribution of the column(s) with up to 200 steps
  - String summary
    - A statistic for the frequency distribution of substrings in a string column which is used for like predicates.

© 2014 SolidQ

7

## DBCC SHOW\_STATISTICS

- General Information on statistics
- Histogram

AAADKHZRMF is the smallest value in the sample  
There are about 61 rows with the value APQYHJBZRA  
The range between ALPAAAMGEBU and APQYHJBZRA has

- ~2213 rows
- 322 distinct values
- Average of 6.876072 row per values

	RANGE_HI_KEY	RANGE_ROWS	EQ_ROWS	DISTINCT_RANGE_ROWS	AVG_RANGE_ROWS
1	AAADKHZRMF	0	1	0	1
2	ABRWPRDERS	899.0234	51.66787	123	7.293088
3	ADUPPCHEBH	2218.063	91.55035	298	5.851113
4	ACQVATWING	1248.25	77.55181	171	7.296851
5	ALPAAAMGEBU	1038.242	75.15327	185	5.613848
6	APQYHJBZRA	2213.344	61.06203	322	6.876072
7	ATVYDZNLV	1640.389	100.0328	254	5.174891
8	AVVLAITHBP	1138.707	51.66787	204	5.586887

© 2008 Solid Quality Mentors

8

## General Assumptions

- Independence
  - Distribution on different columns are independent except correlation information is available
- Uniformity
  - Within each histogram step, distinct values are evenly spread
- Containment and Inclusion
  - If something is searched for, it is assumed that it exists.
  - For example in equi-joins and filter predicates

© 2014 SolidQ

9

## Demo

- Statistic and Cardinality Estimation Basics

© 2014 SolidQ

10

## Creating Statistics

- Automatically
  - With every DML statement that needs statistics, if allowed to do so
    - Only single-columns statistics
    - No filtered statistics
- Explicitly
  - CREATE STATISTICS
  - CREATE INDEX
  - sp\_createstats

© 2014 SolidQ

11

## Updating Statistics

- Automatically triggered through query compilation
  - Colmodctr exceeds threshold
  - Allowed to do so
    - Database level setting
    - Overruled by table level setting
    - Overruled by statistic level setting
    - Set through ALTER DATABASE, UPDATE STATISTICS and sp\_autostats
- Manually
  - UPDATE STATISTICS
  - Index REBUILD
  - sp\_updatestats

© 2014 SolidQ

12

## Update Statistics threshold

- Cardinality goes from 0 to 1
- Old statistics was created with less than 500 rows and colmodctr is greater than 500
- Old statistics with more than 500 rows and colmodctr > 500 + 20% of the rows
  - For filtered statistics the threshold is modified by the selectivity
- Temporary tables have an additional threshold with cardinality of 6

© 2014 SolidQ

13

## Managing Statistics

### Auto Create Statistics

- SQL Server can create new statistics if needed by Query Optimizer

### Auto Update Statistics

- SQL Server can update outdated statistics
- Based on a threshold on data changes
- Tracked per column since SQL Server 2005

### Auto Update Statistics Asynchronously

- Statistic is asynchronously update when query optimizer finds a outdated statistic during optimization
- Uses old statistic for the optimization that triggers the update

© 2010 Solid Quality Mentors

14

## Manual Statistics

- CREATE / UPDATE STATISTICS
  - Ability to turn of automatic update
  - Can be used to create Multi-column statistics
  - Custom sample percentage
- sp\_createstats, sp\_updatestats

```
CREATE STATISTICS statistics_name ON {
table_or_indexed_view_name } ( column [ ,...n ] )
[ WHERE <filter_predicate> ]
[ WITH
[ [ FULLSCAN | SAMPLE number { PERCENT | ROWS } ]
STATS_STREAM = stats_stream ]
[ [ , ] NORECOMPUTE ]
[ [ [ , ] INCREMENTAL = { ON | OFF } ] ] ] ;
```

© 2010 Solid Quality Mentors

15

## General Best Practices

- Keep Auto Options ON
  - Use opt-out for exceptions
    - Statistics that shouldn't be updated
    - Statistics that need a different sampling rate
- Update statistics regularly
  - Don't let the colmodctr trigger an update too often
  - 20% of changes might be too late

© 2014 SolidQ

16

## Demo

- Auto Options
- Querying Statistics

© 2014 SolidQ

17

## Agenda

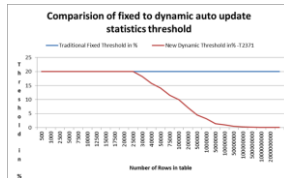
- Cardinality Estimation Basics
- Statistics Basics
  - Definition
  - Creating and Updating Statistics
  - Best Practices
- Common Cardinality Estimation Challenges
- The New Cardinality Estimator
  - Changes
  - Regression Problems and possible Solutions

© 2014 SolidQ

18

## Challenges - Large Tables

- 20% might be too much as threshold for an update
- Trace Flag 2371
  - lower threshold based on # rows
  - <https://support.microsoft.com/en-us/kb/2754171>



Source: <http://blogs.msdn.com/b/asp/sqlserver/archive/2011/09/07/changes-to-automatic-update-statistics-in-sql-server-traceflag-2371.aspx>

© 2014 SolidQ

19

## Challenges – Increasing Values

- Not an issue for highly selective queries
- Update more frequently
- Use Trace Flags for auto-quick-statistics
  - Trace Flag 2389
    - Enable update for known ascending keys
  - Trace Flag 2390
    - Enable also for keys marked as unknown
  - Trace Flag 4139
    - Enable also for keys marked as stationary
    - Since SQL 2012 SP1 CU10

© 2014 SolidQ

20

## Challenges - Multi-column filter

- Assume a query on a table with 10 Mio rows and 2 filters
  - Manufacturer = 'VW' estimated 5 Mio rows
  - Model = 'Golf' estimated 4 Mio rows
- How many rows should SQL Server estimate?
  - 4 Mio rows -> fully dependent
  - 0 rows -> excluding
  - 2 Mio rows -> independent (10x0,5x0,4)
  - Or anything else between 0 and 4 Mio rows?
- SQL Server considers the filters to be independent
- If that isn't the truth
  - Create a multi-column statistic
    - Multi columns statistics are not created automatically!
  - Filtered statistics might also be solution

© 2014 SolidQ

21

## Challenges - Heuristics

- Missing statistics
- Use of local variables in predicates
  - Recompile might be solution (be careful)
- Multistatement table-valued function
  - Use inline table-valued function when possible
- Table variables
  - No statistics
  - No idea of number of rows in the table
  - Use temporary tables
- Non-constant-foldable expressions
  - `YEAR(orderdate) = 2015`
  - Rewrite the expression
  - Use statistics on computed columns
- Comparison of columns
  - `Shipdate > Duedate`
  - Rewrite as expression and use computed column

© 2008 Solid Quality Mentors

22

## Agenda

- Cardinality Estimation Basics
- Statistics Basics
  - Definition
  - Creating and Updating Statistics
  - Best Practices
- Common Cardinality Estimation Problems
- **The New Cardinality Estimator**
  - Changes
  - Regression Problems and possible Solutions

© 2014 SolidQ

23

## New Cardinality Estimator

- Old Cardinality Estimator
  - Largely based on SQL Server 7.0
  - Fixes in QFEs under trace flag to prevent regressions
  - Some problems needed a major redesign
- SQL Server 2014 CE
  - Statistics stays the same “only” the estimation changes
  - Can cause regressions
  - However most query will benefit from changes

© 2014 SolidQ

24



## New Cardinality Estimator

### ► Configuration

- New Cardinality Estimator
  - Compatibility level 120
  - Trace flag 2312 (Query Level)
- Old Cardinality Estimator
  - Compatibility level 110
  - Trace flag 9481 (Query or System Level)
- Test for regression problems
- Extended Events
  - query\_optimizer\_estimate\_cardinality
  - query\_optimizer\_force\_both\_cardinality\_estimation\_behaviors
- CardinalityEstimationModelVersion in Execution Plans

© 2014 SolidQ

25

---

---

---

---

---

---

---

---

## New Cardinality Estimator

### ► Main changes

- Ascending data
  - 120 CE use average cardinality for each value in the column
- Same table filtered predicates
  - Now assumes there's some correlation
- Different tables filtered predicates
  - Now assumes more independence
- Joins
  - More stable algorithms
- Heuristics
  - Multi Statement Table-Valued Function assume 100 rows now

© 2014 SolidQ

26

---

---

---

---

---

---

---

---

## New Cardinality Estimator

### ► Multiple filters on a single table

- Selectivity of a single filter is defined as
  - estimated divided by the total number of rows
- Old CE assumes no dependency with multiple filters
  - $Estimate = Sel1 \times Sel2 \times \dots \times total\ rows$
- New CE estimates some dependency
  - Exponential backoff algorithm
  - $Estimate = Sel1 \times \sqrt{Sel2} \times \sqrt[3]{Sel3} \times \sqrt[4]{Sel4}$
  - Selectivity sorted from highest to lowest
  - Up to 4 filters used
- for highly dependent data also consider other solutions
  - multi-column statistics
  - Filtered statistics

© 2014 SolidQ

27

---

---

---

---

---

---

---

---

## Increasing Values

- Not an issue for highly selective queries
- The new estimator considers the same distribution as within the histogram
  - Trace Flags typically not needed anymore

© 2014 SolidQ

28

## New Cardinality Estimator

- ▶ query\_optimizer\_estimate\_cardinality

- Only for troubleshooting – debug channel

[illegible]

- Correlate with query plan using stats\_collection\_id

NoExpandHint	False
Object	[AdventureWorks2012]
Ordered	True
Output List	Bmk1000; [AdventureW
Parallel	False
Physical Operation	Index Seek
Scan Direction	FORWARD
Seek Predicates	Seek Keys[1]: Prefix: [A
StatsCollectionId	2
Storage	RowStore
TableCardinality	31125

© 2014 SolidQ

29

## New Cardinality Estimator

▶ Resources

- Miloš Radivojević Blog Series
  - <http://milossql.wordpress.com/tag/cardinality-estimator>
- A First Look at the New SQL Server Cardinality Estimator
  - <http://www.sqlperformance.com/2013/12/t-sql-queries/a-first-look-at-the-new-sql-server-cardinality-estimator>
- New functionality in SQL Server 2014 – Part 2 – New Cardinality Estimation
  - <http://blogs.msdn.com/b/sapnsqlserver/archive/2014/01/16/new-functionality-in-sql-server-2014-part-2-new-cardinality-estimation.aspx>

© 2014 SolidQ

30