

Carnegie Mellon University

OPTIMIZE!

Database Query Optimization

Query Cost Models: Cardinality Estimation

ADMINISTRIVIA

Project #1 is due Friday Feb 28th.

Project #2 proposals are due Monday Mar 10th.

LAST CLASS

We discussed the data structures that DBMSs maintain to summarize the contents of tables.

→ Most Common: Equi-Depth, End-Bias Histograms (Heavy Hitters), HyperLogLog

Today's class is about how the optimizer's cost model uses these data structures...

CARDINALITY ESTIMATION

Cardinality estimation is the process of predicting the number of rows that will be returned by a query operation, such as a filter or join, to help the optimizer choose the most efficient execution plan.

The number of tuples that will be generated per operator is computed from its selectivity multiplied by the number of tuples in its input.

TODAY'S AGENDA

Cardinality Estimation Basics

The Germans' Survey

Implementations

CARDINALITY ESTIMATION

There are three cardinality estimations on logical expressions an optimizer must support as the core of its cost model:

- Selection Conditions (filters)
- Join Size Estimation
- Distinct Value Estimation

These will serve as building blocks for more complex expressions:

- Multiple Joins
- Group By

DERIVABLE STATISTICS

For each relation R , the DBMS maintains statistics to approximate the following information:

- N_R : Number of tuples in R .
- $V(A, R)$: Number of distinct values for attribute A .

The selection cardinality $SC(A, R)$ is the average number of records with a value for an attribute A given $N_R / V(A, R)$

SINGLE SELECTION CONDITION

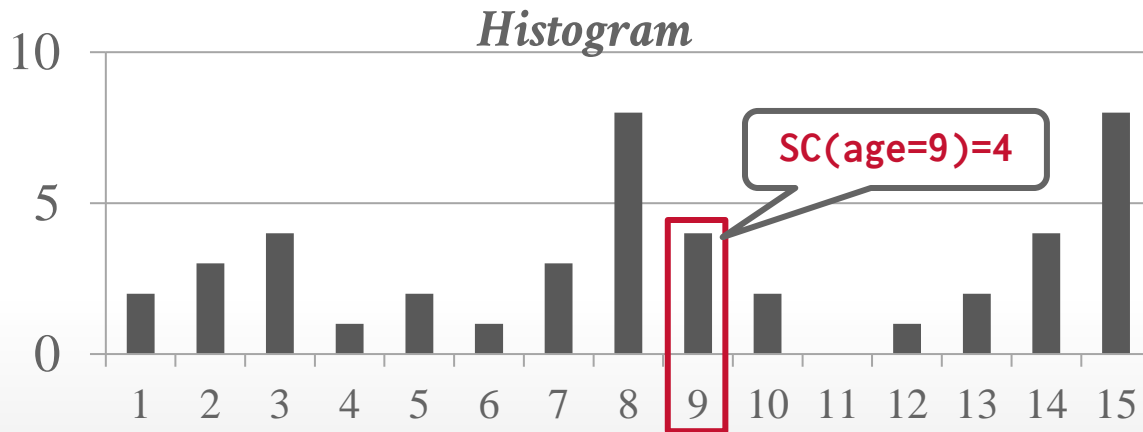
The selectivity (*sel*) of a predicate *P* is the fraction of tuples that qualify.

```
SELECT * FROM people  
WHERE age = 9
```

Equality Predicate: *A=constant*

→ $sel(A=constant) = \#occurrences / |R|$

→ Example: $sel(age=9) = 4 / 45 = 0.088$



SINGLE SELECTION CONDITION

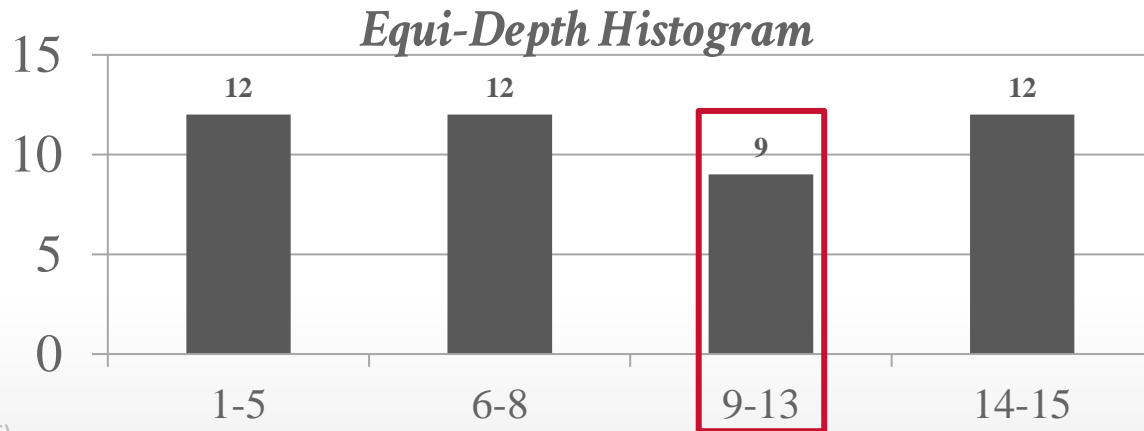
The selectivity (*sel*) of a predicate *P* is the fraction of tuples that qualify.

```
SELECT * FROM people
WHERE age = 9
```

Equality Predicate: $A = \text{constant}$

→ $sel(A = \text{constant}) = \# \text{occurrences} / |R|$

→ Example: $sel(\text{age} = 9) = 4 \cancel{\times} 45 = 0.088$
 $\approx (9/5) / 45 \approx 1.8 / 45 \approx 0.04$



ASSUMPTIONS

Assumption #1: Uniform Data

→ The distribution of values (except for the heavy hitters) is the same within a histogram bucket.

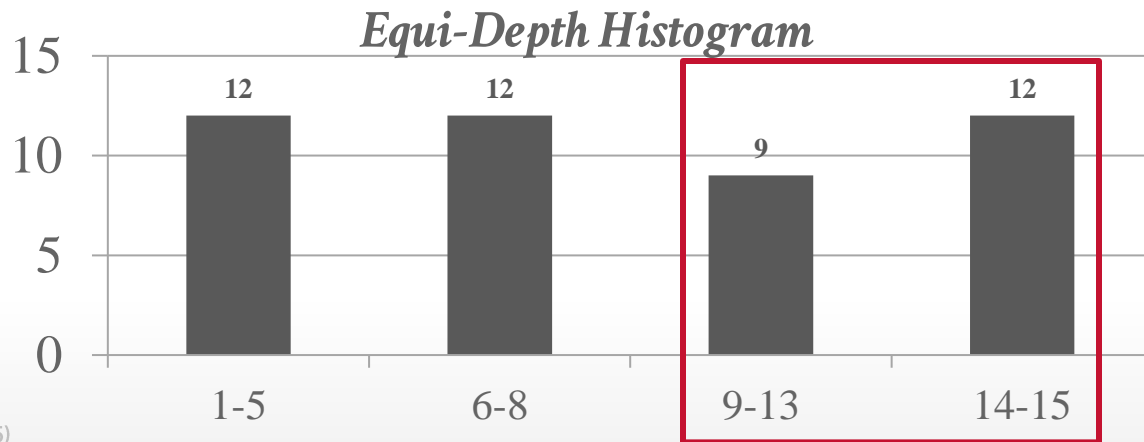
SINGLE SELECTION CONDITION

Range Predicate:

→ $sel(A \geq a) = (\#RANGE-ROWS + \#EQ-ROWS) / |R|$

→ Example: $sel(age \geq 7) \approx ((9+12))$

```
SELECT * FROM people  
WHERE age >= 7
```



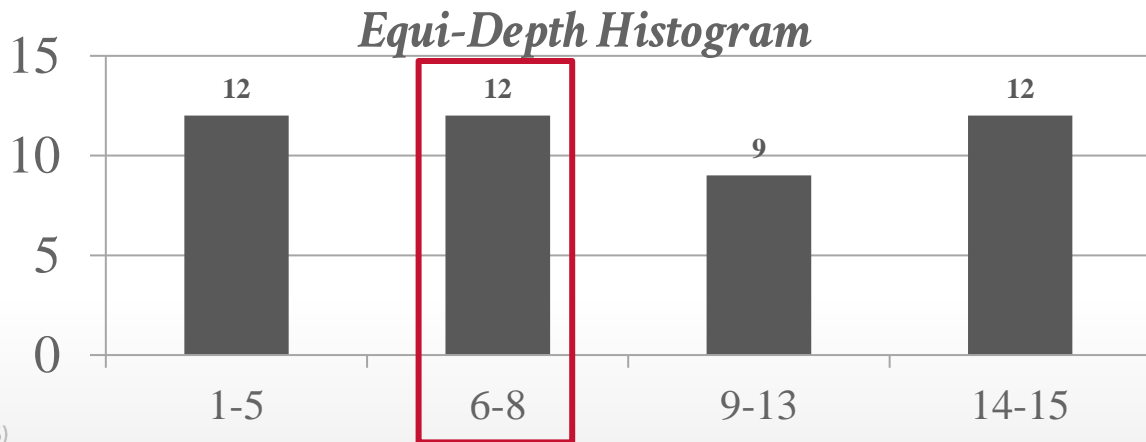
SINGLE SELECTION CONDITION

Range Predicate:

→ $sel(A \geq a) = (\#RANGE-ROWS + \#EQ-ROWS) / |R|$

→ Example: $sel(age \geq 7) \approx ((9+12) + (2 \times (12/3))) / 45$
 $\approx 29 / 45 \approx 0.6444$

```
SELECT * FROM people  
WHERE age >= 7
```



SINGLE SELECTION CONDITION

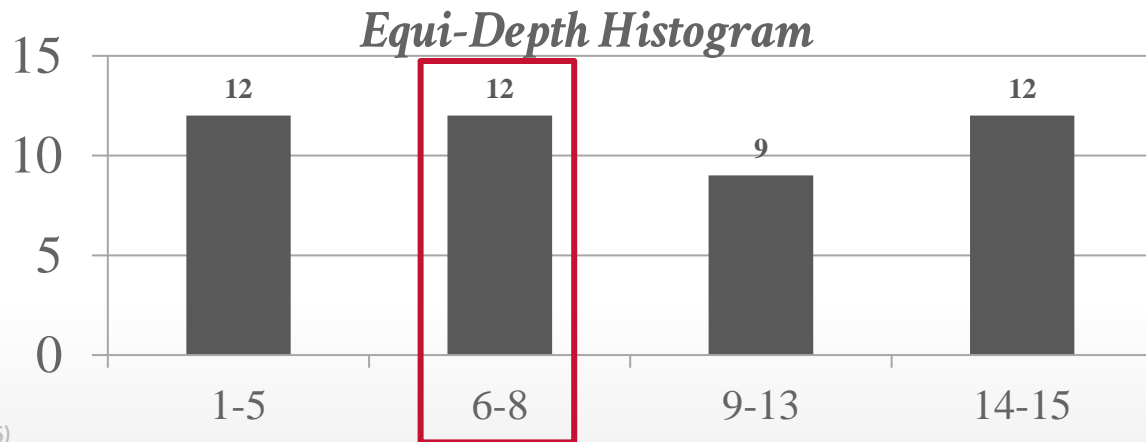
Range Predicate:

→ $sel(A \geq a) = (\#RANGE-ROWS + \#EQ-ROWS) / |R|$

→ Example: $sel(age \geq 7) \approx ((9+12) + (2 \times (12/3))) / 45$
 $\approx 29 / 45 \approx 0.6444$

```
SELECT * FROM people
WHERE age >= 7
```

! *This assumes continuous distribution of values.*



SINGLE SELECTION CONDITION

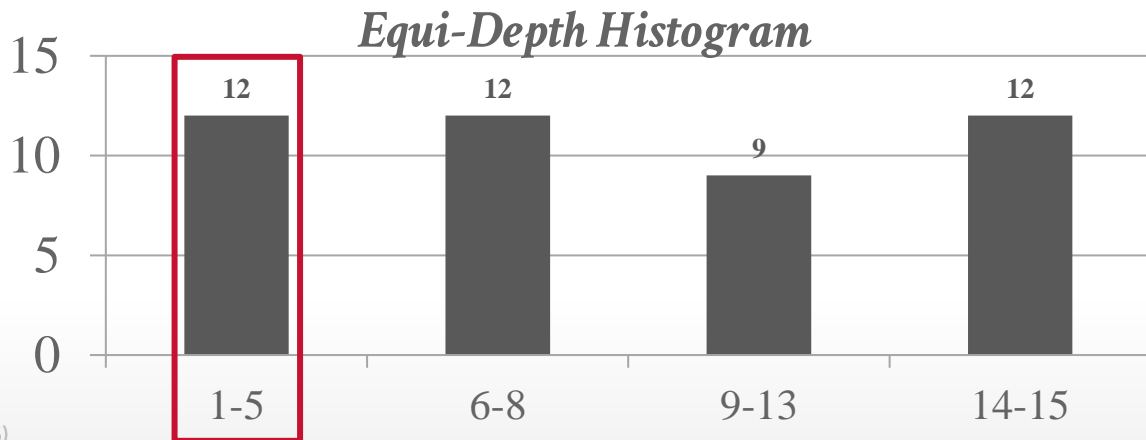
Negation Query:

→ $sel(not\ P) = 1 - sel(P)$

→ Example: $sel(age \neq 2) \approx 1 - ((12/5) / 45)$
 $\approx 1 - (2.4 / 45) \approx 1 - 0.05 \approx 0.95$

```
SELECT * FROM people  
WHERE age != 2
```

⚠ Observation: Selectivity \approx Probability



OBSERVATION

We now can compute selectivities for individual predicates, but what happens if there are multiple predicates in a query?

→ Even though the predicates are on the same table, the attributes may have different distributions.

```
SELECT * FROM people
WHERE age = 2
      AND name LIKE 'A%'
```

OBSERVATION

We now can compute selectivities for individual predicates, but what happens if there are multiple predicates in a query?

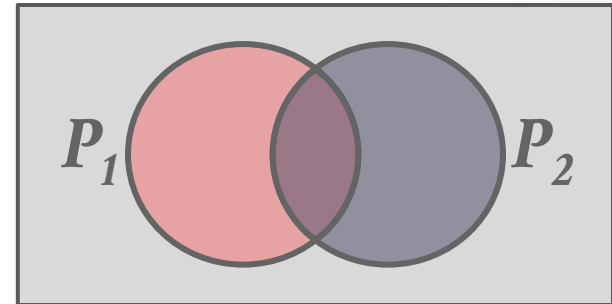
→ Even though the predicates are on the same table, the attributes may have different distributions.

Example:

→ $sel(age = 2) \approx 0.053$

→ $sel(name LIKE 'A\%') \approx 0.1$

```
SELECT * FROM people
WHERE age = 2
      AND name LIKE 'A%'
```



OBSERVATION

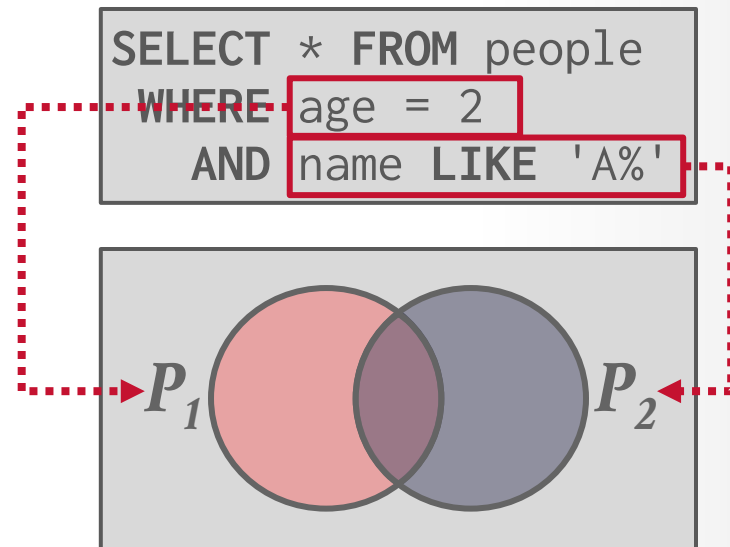
We now can compute selectivities for individual predicates, but what happens if there are multiple predicates in a query?

→ Even though the predicates are on the same table, the attributes may have different distributions.

Example:

→ $sel(age = 2) \approx 0.053$

→ $sel(name LIKE 'A\%') \approx 0.1$



OBSERVATION

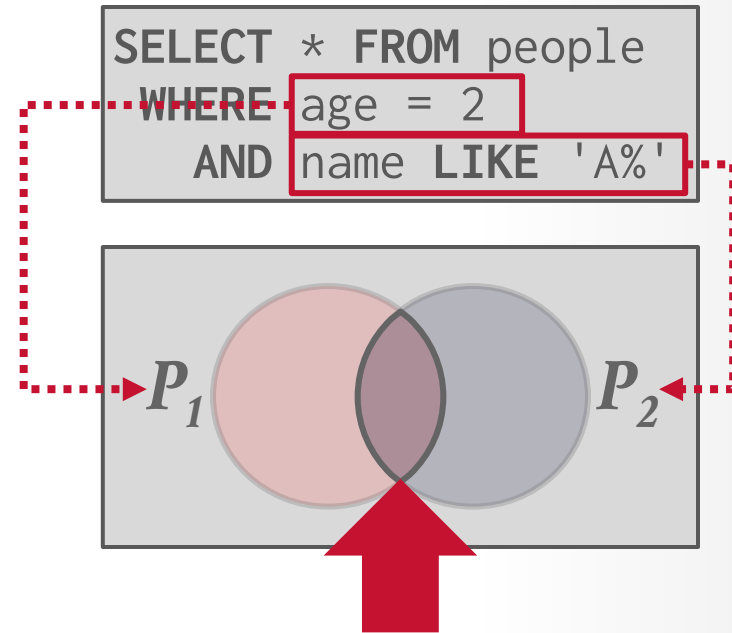
We now can compute selectivities for individual predicates, but what happens if there are multiple predicates in a query?

→ Even though the predicates are on the same table, the attributes may have different distributions.

Example:

→ $sel(age = 2) \approx 0.053$

→ $sel(name LIKE 'A\%') \approx 0.1$



ASSUMPTIONS

Assumption #1: Uniform Data

- The distribution of values (except for the heavy hitters) is the same within a histogram bucket.

Assumption #2: Independent Predicates

- The predicates on attributes are independent. The selectivity of the conjunction of two or more predicates is estimated as the product of their individual selectivities.

MULTIPLE SELECTION CONDITION

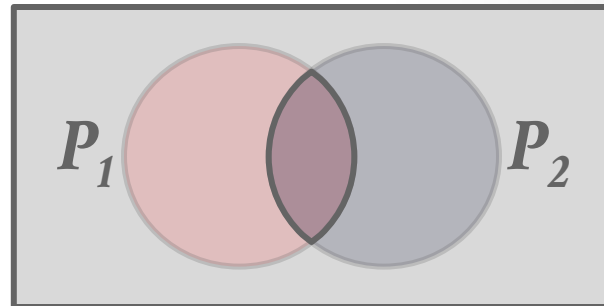
Conjunction:

→ $sel(P1 \wedge P2) = sel(P1) \times sel(P2)$

→ Example: $sel(age=2 \wedge name \text{ LIKE 'A\%' })$
 $\approx sel(age=2) \times sel(name \text{ LIKE 'A\%' })$
 $\approx 0.053 \times 0.1 \approx 0.0053$

This assumes that the predicates are independent.

```
SELECT * FROM people  
WHERE age = 2  
      AND name LIKE 'A%'
```



MULTIPLE SELECTION CONDITION

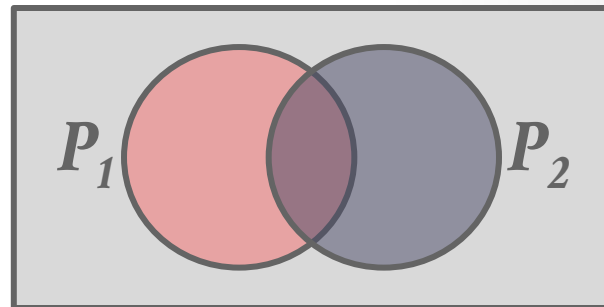
Disjunction:

→ $sel(P1 \vee P2)$

$\approx sel(P1) + sel(P2) - sel(P1 \wedge P2)$

$\approx sel(P1) + sel(P2) - sel(P1) \times sel(P2)$

```
SELECT * FROM people
WHERE age = 2
      OR name LIKE 'A%'
```



MULTIPLE SELECTION CONDITION

Disjunction:

→ $sel(P1 \vee P2)$

$$\approx sel(P1) + sel(P2) - sel(P1 \wedge P2)$$

$$\approx sel(P1) + sel(P2) - sel(P1) \times sel(P2)$$

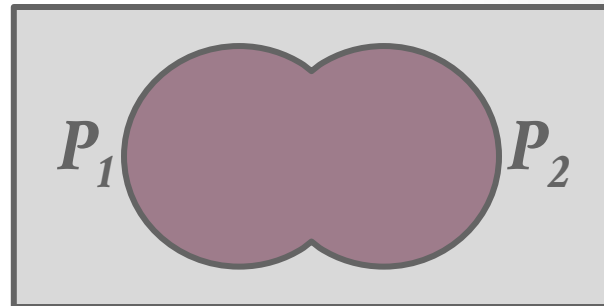
→ Example: $sel(age=2 \vee name \text{ LIKE } 'A\%')$

$$\approx 0.053 + 0.1 - (0.053 \times 0.1)$$

$$\approx 0.1477$$

This again assumes that the selectivities are **independent**.

```
SELECT * FROM people
WHERE age = 2
      OR name LIKE 'A%'
```



CORRELATED ATTRIBUTES

Consider a database of automobiles:

→ # of Makes = 10, # of Models = 100

Then the following query shows up:

→ `(make="Honda" AND model="Accord")`

With the independence and uniformity assumptions, the selectivity is:

→ 1

$/10 \times 1/100 \approx 0.001$

But since only Honda makes Accords the real selectivity is *$1/100 = 0.01$*

JOIN SIZE ESTIMATION

Given a join of R and S , what is the range of possible result sizes in # of tuples?

→ In other words, for a given tuple of R , how many tuples of S will it match?

Assume each key in the inner relation will exist in the outer table.

ASSUMPTIONS

Assumption #1: Uniform Data

- The distribution of values (except for the heavy hitters) is the same within a histogram bucket.

Assumption #2: Independent Predicates

- The predicates on attributes are independent. The selectivity of the conjunction of two or more predicates is estimated as the product of their individual selectivities.

Assumption #3: Containment Principle

- The domain of join keys overlap such that each key in the inner relation will also exist in the outer table.

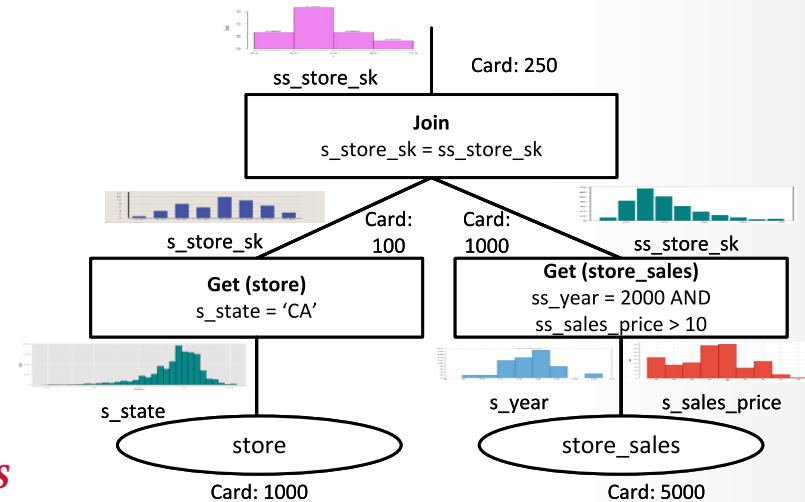
JOIN SIZE ESTIMATION

Align the buckets of the histograms so that their boundaries agree.

Compute a per bucket estimation of join size.

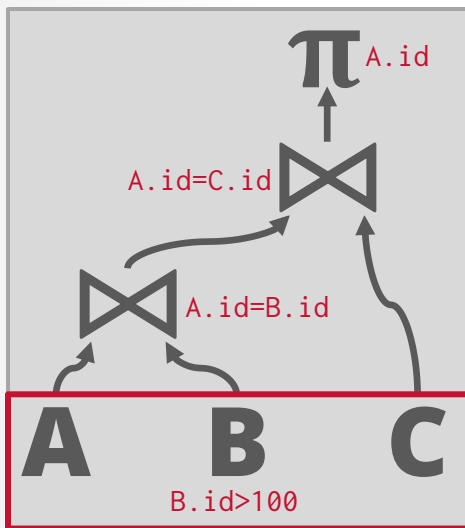
- Containment assumes that for each group g_R of tuples in R , it has a corresponding group g_S in S .
- Each tuple in g_R will match with tuples in g_S

Aggregate partial frequencies from joining each pair of buckets to get cardinality of the whole join.



ERROR PROPAGATION

```
SELECT A.id
FROM A, B, C
WHERE A.id = B.id
      AND A.id = C.id
      AND B.id > 100
```



Compute the cardinality of base tables

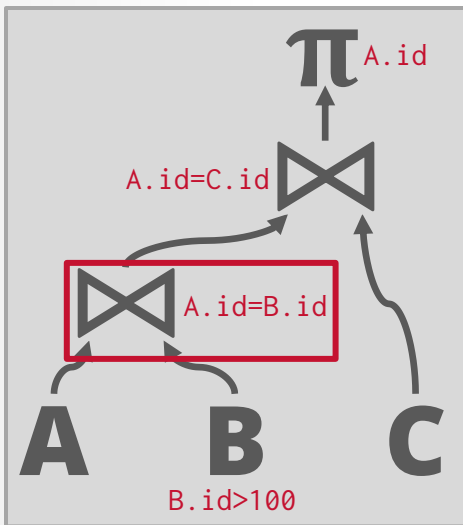
$$A \rightarrow |A|$$

$$B.id > 100 \rightarrow |B| \times \text{sel}(B.id > 100)$$

$$C \rightarrow |C|$$

ERROR PROPAGATION

```
SELECT A.id
FROM A, B, C
WHERE A.id = B.id
      AND A.id = C.id
      AND B.id > 100
```



Compute the cardinality of base tables

$$A \rightarrow |A|$$

$$B.id > 100 \rightarrow |B| \times sel(B.id > 100)$$

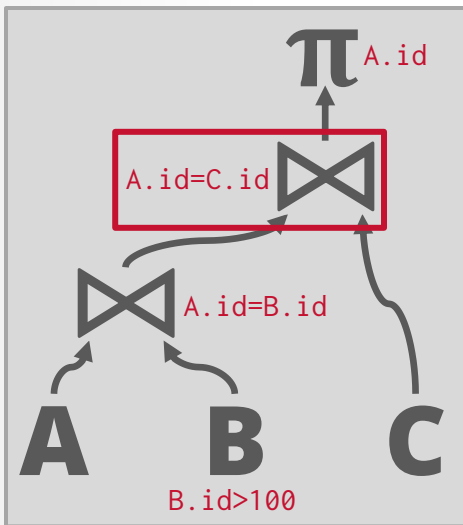
$$C \rightarrow |C|$$

Compute the cardinality of join results

$$A \bowtie B \approx (|A| \times |B|) / \max(sel(A.id = B.id), sel(B.id > 100))$$

ERROR PROPAGATION

```
SELECT A.id
FROM A, B, C
WHERE A.id = B.id
      AND A.id = C.id
      AND B.id > 100
```



Compute the cardinality of base tables

$$A \rightarrow |A|$$

$$B.id > 100 \rightarrow |B| \times \text{sel}(B.id > 100)$$

$$C \rightarrow |C|$$

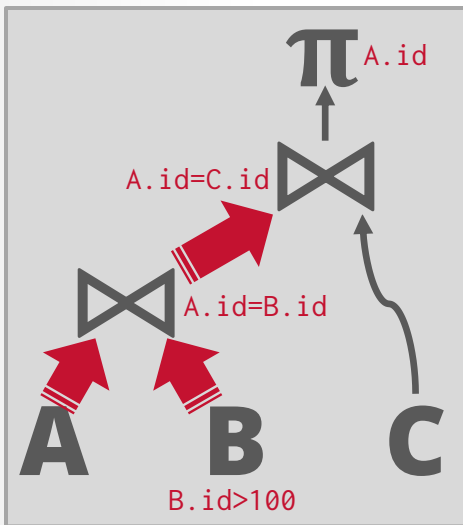
Compute the cardinality of join results

$$A \bowtie B \approx (|A| \times |B|) / \max(\text{sel}(A.id = B.id), \text{sel}(B.id > 100))$$

$$(A \bowtie B) \bowtie C \approx (|A| \times |B| \times |C|) / \max(\text{sel}(A.id = B.id), \text{sel}(B.id > 100), \text{sel}(A.id = C.id))$$

ERROR PROPAGATION

```
SELECT A.id
FROM A, B, C
WHERE A.id = B.id
      AND A.id = C.id
      AND B.id > 100
```



Compute the cardinality of base tables

$$A \rightarrow |A|$$

$$B.id > 100 \rightarrow |B| \times sel(B.id > 100)$$

$$C \rightarrow |C|$$

Compute the cardinality of join results

$$A \bowtie B \approx (|A| \times |B|) / \max(sel(A.id=B.id), sel(B.id > 100))$$

$$(A \bowtie B) \bowtie C \approx (|A| \times |B| \times |C|) / \max(sel(A.id=B.id), sel(B.id > 100), sel(A.id=C.id))$$

ESTIMATOR QUALITY

Evaluate the correctness of cardinality estimates generated by DBMS optimizers as the number of joins increases.

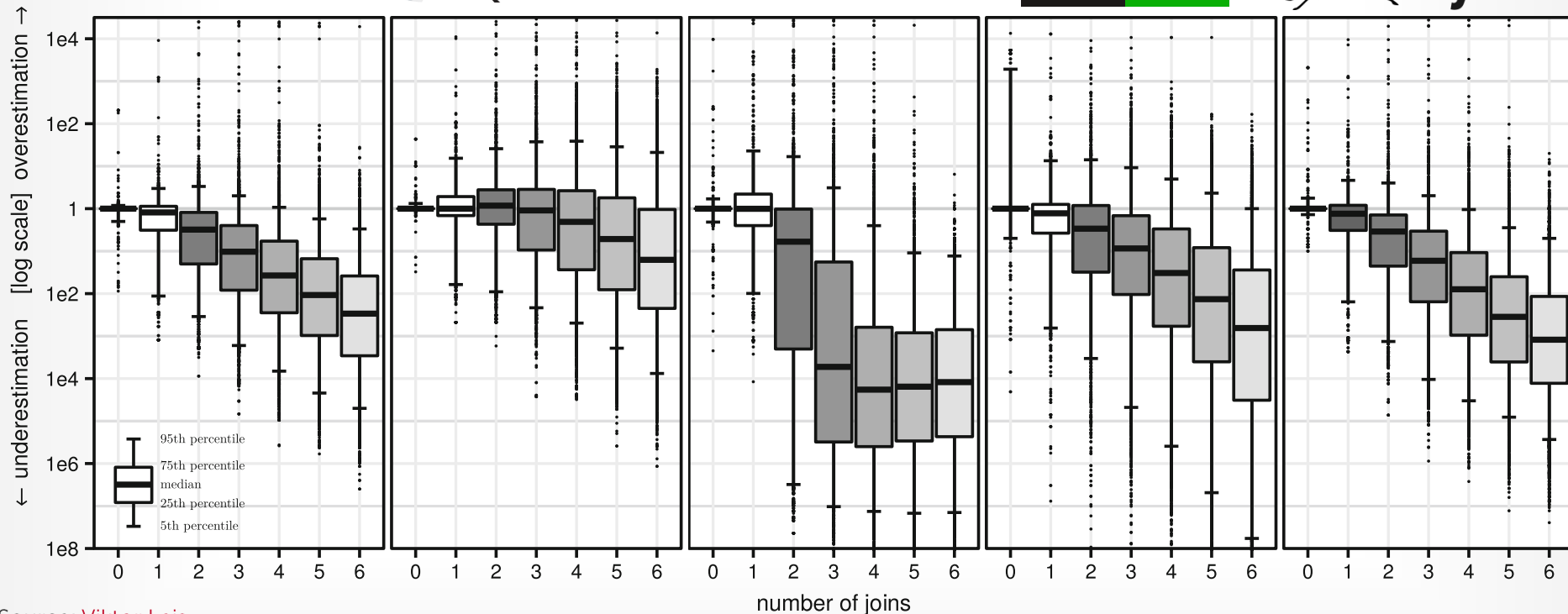
- Let each DBMS perform its stats collection.
- Extract measurements from query plan.

Compared five DBMSs using 100k queries from the JOB workload based on IMDB data.



HOW GOOD ARE QUERY OPTIMIZERS, REALLY?
VLDB 2015

ESTIMATOR QUALITY

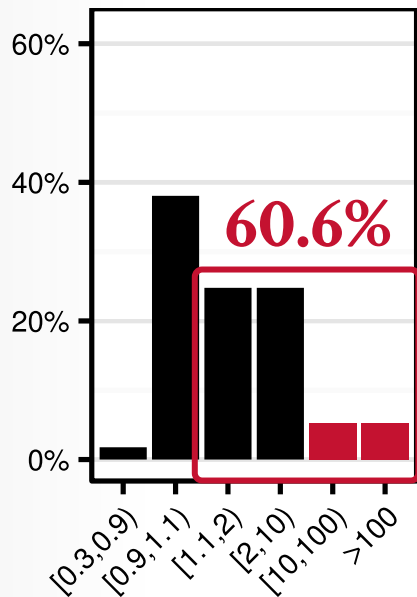


Source: [Viktor Leis](#)

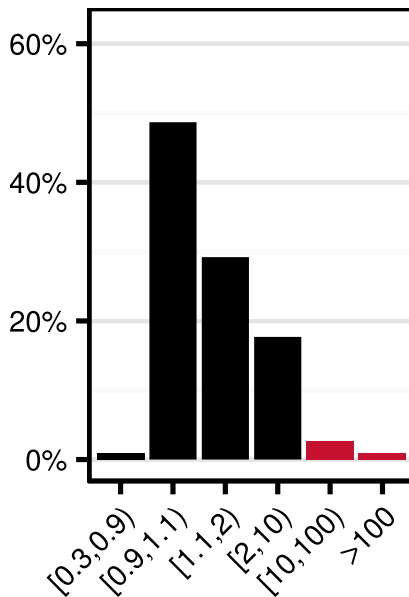
EXECUTION SLOWDOWN

PostgreSQL v9.4 – JOB Workload

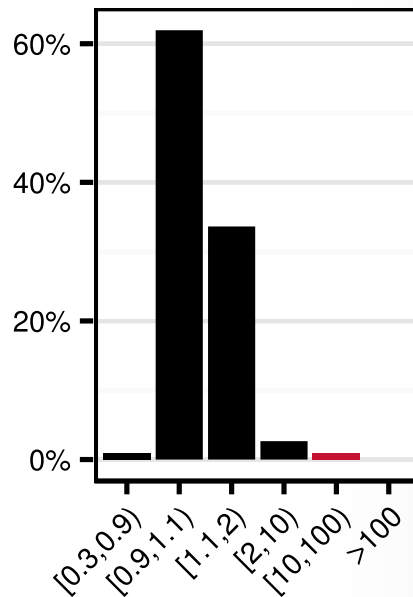
Default Planner



No NL Join



Dynamic Rehashing



Slowdown compared to using true cardinalities

Source: [Viktor Leis](#)

LESSONS FROM THE GERMANS

Query opt is more important than a fast engine

→ Cost-based join ordering is necessary

Cardinality estimates are routinely wrong

→ Try to use operators that do not rely on estimates

Hash joins + seq scans are a robust exec model

→ The more indexes that are available, the more brittle the plans become (but also faster on average)

Working on accurate models is a waste of time

→ Better to improve cardinality estimation instead

IMPLEMENTATIONS

Postgres

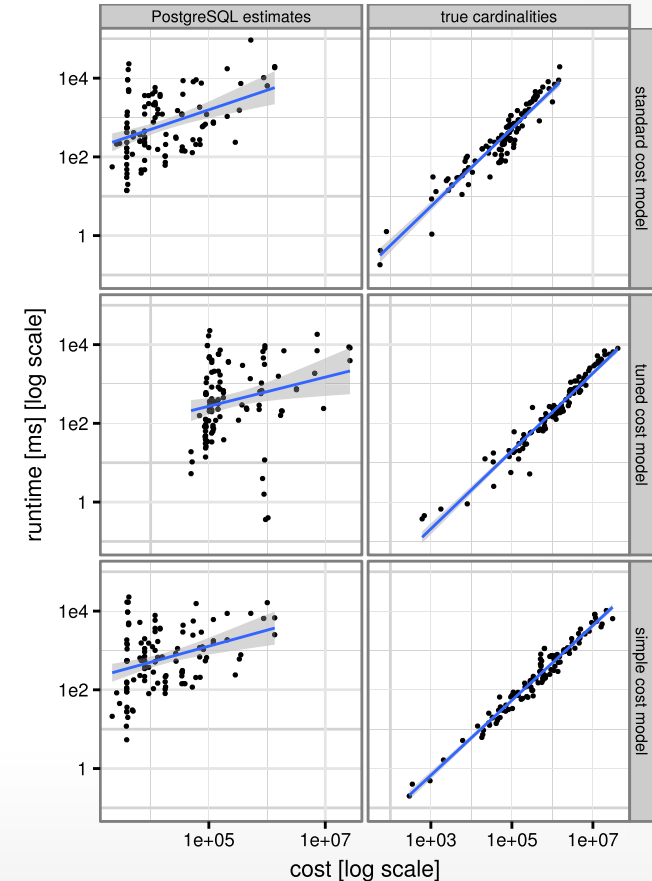
Microsoft SQL Server

POSTGRES COST MODEL

Cost of an operator is a weighted sum of the # of accessed disk pages and the amount of data processed in memory.

- Distinguishes between sequential and random I/O.
- Requires manual tuning the weights based on hardware characteristics.

The Germans replaced Postgres' cost model with simple cost model to see what happens with query plans...



postgres / src / backend / optimizer / path / costsize.c

6605 lines (5909 loc) · 214 KB

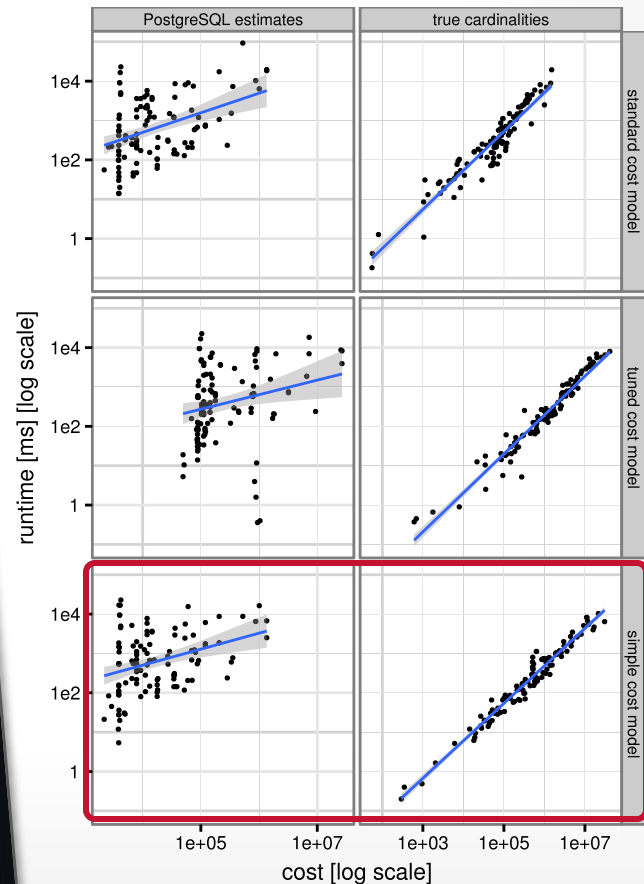
Code Blame

```

1 /*
2  *
3  * costsize.c
4  *   Routines to compute (and set) relation sizes and path costs
5  *
6  * Path costs are measured in arbitrary units established by these basic
7  * parameters:
8  *
9  * seq_page_cost      Cost of a sequential page fetch
10 * random_page_cost   Cost of a non-sequential page fetch
11 * cpu_tuple_cost      Cost of typical CPU time to process a tuple
12 * cpu_index_tuple_cost Cost of typical CPU time to process an index tuple
13 * cpu_operator_cost   Cost of CPU time to execute an operator or function
14 * parallel_tuple_cost Cost of CPU time to pass a tuple from worker to leader backend
15 * parallel_setup_cost Cost of setting up shared memory for parallelism
16 *
17 * We expect that the kernel will typically do some amount of read-ahead
18 * optimization; this in conjunction with seek costs means that seq_page_cost
19 * is normally considerably less than random_page_cost. (However, if the
20 * database is fully cached in RAM, it is reasonable to set them equal.)
21 *
22 * We also use a rough estimate "effective_cache_size" of the number of
23 * disk pages in Postgres + OS-level disk cache. (We can't simply use
24 * NBuffers for this purpose because that would ignore the effects of
25 * the kernel's disk cache.)
26 *
27 * Obviously, taking constants for these values is an oversimplification,
28 * but it's tough enough to get any useful estimates even at this level of
29 * detail. Note that all of these parameters are user-settable, in case
30 * the default values are drastically off for a particular platform.
31 *
32 * seq_page_cost and random_page_cost can also be overridden for an individual
33 * tablespace, in case some data is on a fast disk and other data is on a slow
34 * disk. Per-tablespace overrides never apply to temporary work files such as
35 * an external sort or a materialize node that overflows work_mem.
36 *
37 * We compute two separate costs for each path:
38 *   - the estimated cost to fetch all tuples
39 *   - the estimated cost to fetch all tuples plus the cost of any
40 *     index scans that will be performed on the tuples.
41 */

```

MODEL

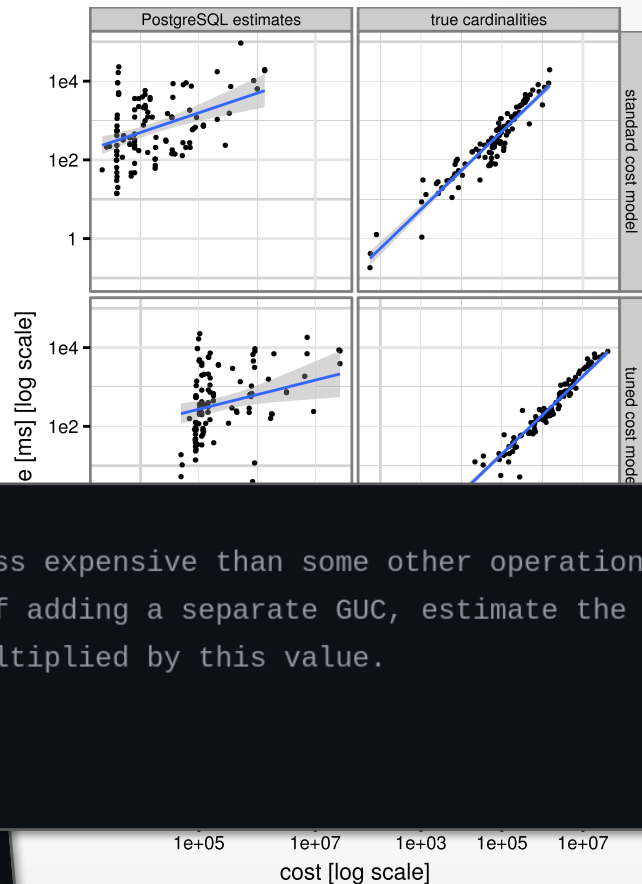


MODEL

```

6605 lines (5909 loc) · 214 KB
Code Blame
1 /*
2  *
3  * costsize.c
4  *   Routines to compute (and set) relation sizes and path costs
5  *
6  * Path costs are measured in arbitrary units established by these basic
7  * parameters:
8  *
9  * seq_page_cost      Cost of a sequential page fetch
10 * random_page_cost   Cost of a non-sequential page fetch
11 * cpu_tuple_cost      Cost of typical CPU time to process a tuple
12 * cpu_index_tuple_cost Cost of typical CPU time to process an index tuple
13 * cpu_operator_cost   Cost of CPU time to execute an operator or function
14 * parallel_tuple_cost Cost of CPU time to pass a tuple from worker to leader backend
15 * parallel_setup_cost Cost of setting up shared memory for parallelism
16 *
17 * We expect that the kernel will typically do some amount of read-ahead
18 * optimization; this in conjunction with seek costs means that seq_page_cost
19 * is normally considerably less than random_page_cost. (However, if the
20 * database is fully sequential, they should be equal.)
21 *
22 * We also use a random_page_cost multiplier to account for the fact that
23 * disk pages in PostgreSQL are not necessarily contiguous.
24 * NBuffers for the buffer pool.
25 * the kernel's disk access cost.
26 *
27 * Obviously, taking disk access into account is a bit of a
28 * but it's tough.
29 * detail. Note that the default value of random_page_cost is 4.
30 *
31 *
32 * seq_page_cost
33 * tablespace, in
34 * disk. Per-tablespace
35 * an external sort or a materialize node that overwrites
36 *
37 * We compute two separate costs for each path:
38 *   - the estimated cost to fetch all tuples
39 *   - the estimated cost to fetch all tuples and execute the query
40 */
41
42 #define APPEND_CPU_COST_MULTIPLIER 0.5
43
44 /*
45 * Append and MergeAppend nodes are less expensive than some other operations
46 * which use cpu_tuple_cost; instead of adding a separate GUC, estimate the
47 * per-tuple cost as cpu_tuple_cost multiplied by this value.
48 */
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100

```



SQL SERVER

Each operator's cost is based on producing all its output and the cost producing the first output.

→ Example: Nested-Loop Join vs. Index Nested-Loop Join

Maintains "row goals" to track the number of tuples are needed above in the query plan.

→ Example: Top-K / LIMIT clause

Combine predicates using an exponential backoff.

Given a table R and predicate selectivities $S_1, S_2, S_3, \dots, S_n$, where S_1 is the most selective and S_n the least:

$$\text{Estimate} = |R| \times S_1 \times \sqrt{(S_2)} \times \sqrt{(\sqrt{(S_3)})} \times \sqrt{(\sqrt{(\sqrt{(S_4))})})} \dots$$

PARTING THOUGHTS

Cardinality estimation is the hardest part of building a query optimizer.

- Lots of simplifying assumptions to make the problem more tractable.
- But developers often don't revisit those assumptions...

Research suggests that accurate cardinality estimation is more important than sophisticated cost models.

All of this seems like it is the perfect usecase for machine learning...

PART TWO THOUGHTS

Cardinality estimation
a query optimizer.

→ Lots of simplifying
tractable.

→ But developers oft

Research suggests
estimation is more
cost models.

All of this seems
machine learning...

Peep this cardinality estimation work



From:

To:

Date:

Spam Status:

Me

2/11/25 2:06 PM

Spamassassin

You're a huge pile shit and I am going to fuck you up when I see you again. I want my money back. Don't think you can roll up in Brooklyn like that again making big assumptions about how do we things in the streets.

Speaking of big assumptions, check out this hot piece of work:

https://github.com/microsoft/documentdb/blob/a0c44348fb473daf75b8950d0118201c0b3e6d19/pg_documentdb_core/src/planner/selectivity.c

Unmatched.

Erik Darling

<https://erikdarling.com/>

I'll make your SQL Server faster in exchange for money.



main

documentdb / pg_documentdb_core / src / planner / selectivity.c



Code

Blame

Raw



```
1  /*-----
2   * Copyright (c) Microsoft Corporation. All rights reserved.
3   *
4   * src/planner/selectivity.c
5   *
6   * Implementation of selectivity functions for BSON operators.
7   *
8   *-----
9   */
10 #include <postgres.h>
11 #include <fmgr.h>
12
13
14 PG_FUNCTION_INFO_V1(bson_operator_selectivity);
15
16
17 /*
18  * bson_operator_selectivity returns the selectivity of a BSON operator
19  * on a relation.
20  */
21 Datum
22 bson_operator_selectivity(PG_FUNCTION_ARGS)
23 {
24     /* dumbest possible implementation: assume 1% of rows are returned */
25     PG_RETURN_FLOAT8(0.01);
26 }
```

TS

y estimation work

25 2:06 PM

assassin

going to fuck you up when I see you again. I
you can roll up in Brooklyn like that again
w do we things in the streets.

ok out this hot piece of work:

mentdb/blob/a0c44348fb473daf75b8950d011
re/src/planner/selectivity.c

change for money.



NEXT CLASS

Project #2 Proposals