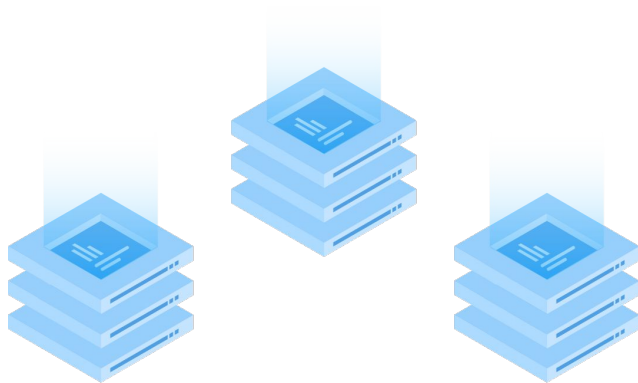


Discovering Statistical Properties for Query Optimization

Presented by Han Fei



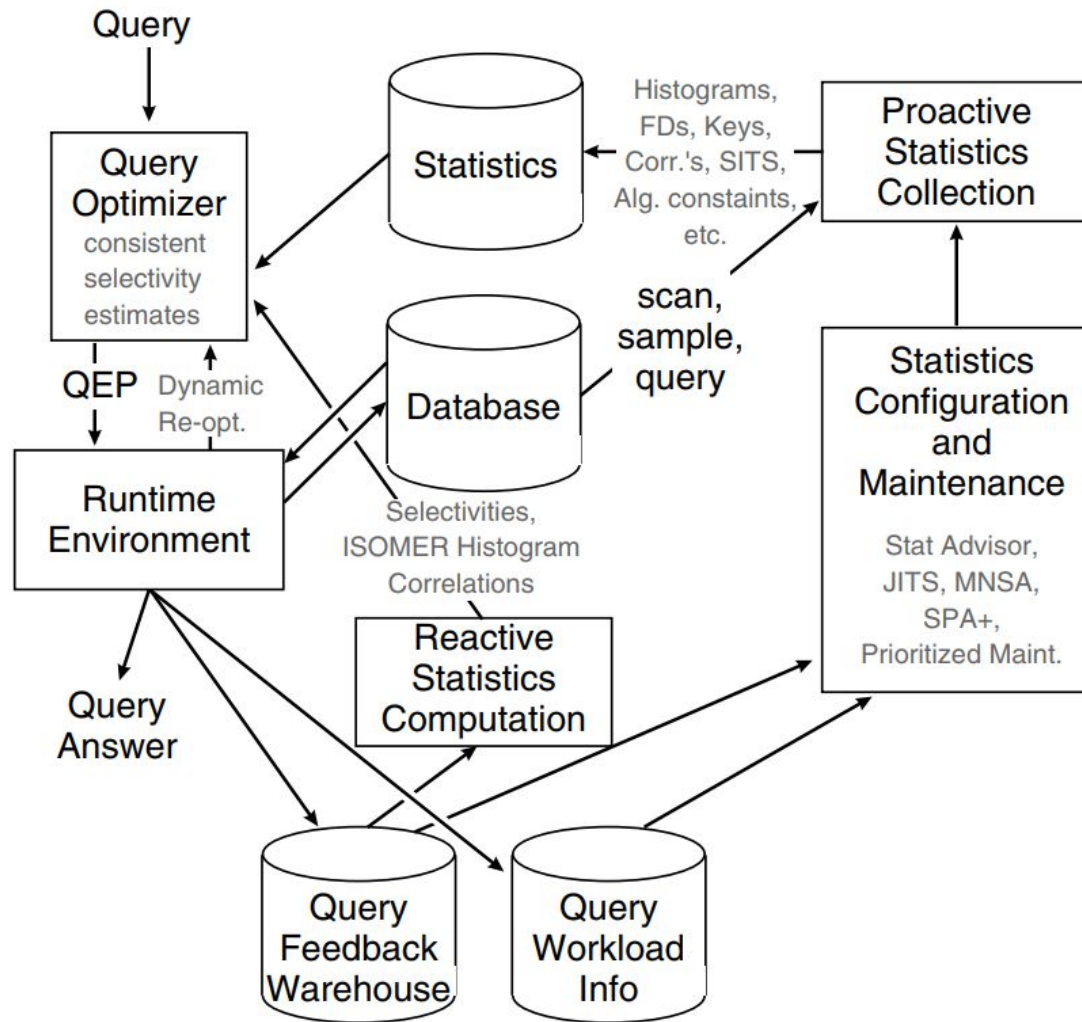
Part I - Introduction to Query Optimization



The Process of Query Optimization

- Enumerate all possible plans by transformation rules (Cascade Optimzier)
- Get the ordinary table/columns cardinality by simple predicate
- Derive cardinality for every subplan
- Estimate the cost of plan and choose the best one





Most Common Statistics : Histogram

- Bucket Scheme
 - Equi-Width
 - Equi-Depth
 - Max Diff
 - V-Optimal
- Estimation Scheme
 - Continuous Spread Assumption
 - Four Level Tree



What do We Need More for Modern Databases?

- Correlation Discovery
- 'soft' Functional Dependency
- Algebraic constraints
- Holes in joins
- Workload-aware Methods
- Incremental Maintenance



Part II - Proactive Methods



How do we estimate multiple predicates?

- attribute value independence assumption
 - the distributions of individual attributes are independent of each other
- join uniformity assumption
 - a tuple from one relation is equally likely to join with any tuple from the second relation



Automated Correlation Discovery

- Consider a table with tree attributes:
 - Education
 - Income
 - Home-holder
- some of the correlations between attributes might be indirect ones, mediate by others.
 - a high-school dropout who owns a successful Internet startup is more likely to own a home than a highly educated beach bum.

E	I	H	$P(E, I, H)$
h	l	f	0.27
h	l	t	0.03
h	m	f	0.105
h	m	t	0.045
h	h	f	0.005
h	h	t	0.045
c	l	f	0.135
c	l	t	0.015
c	m	f	0.063
c	m	t	0.027
c	h	f	0.006
c	h	t	0.054
a	l	f	0.018
a	l	t	0.002
a	m	f	0.042
a	m	t	0.018
a	h	f	0.012
a	h	t	0.108

Conditionally Independent

- $P(H = h \mid E = e, I = i) = P(H = h \mid I = i)$
- We just need hold some marginal distribution
 - $P(E)$
 - $P(I \mid E)$
 - $P(H \mid I)$
- Then $P(H, E, I) = P(E) P(I \mid E) P(H \mid I)$

E	$P(E)$
h	0.5
c	0.3
a	0.2

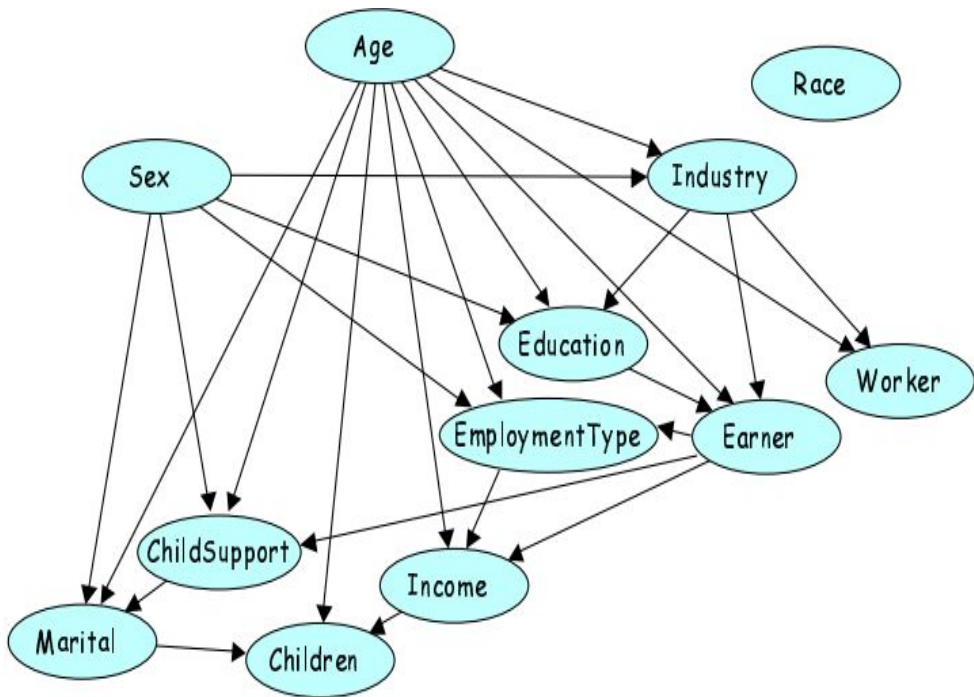
I	E	$P(I \mid E)$
l	h	0.6
m	h	0.3
h	h	0.1
l	c	0.5
m	c	0.3
h	c	0.2
l	a	0.1
m	a	0.3
h	a	0.6

H	I	$P(H \mid I)$
t	l	0.1
f	l	0.9
t	m	0.3
f	m	0.7
t	h	0.9
f	h	0.1

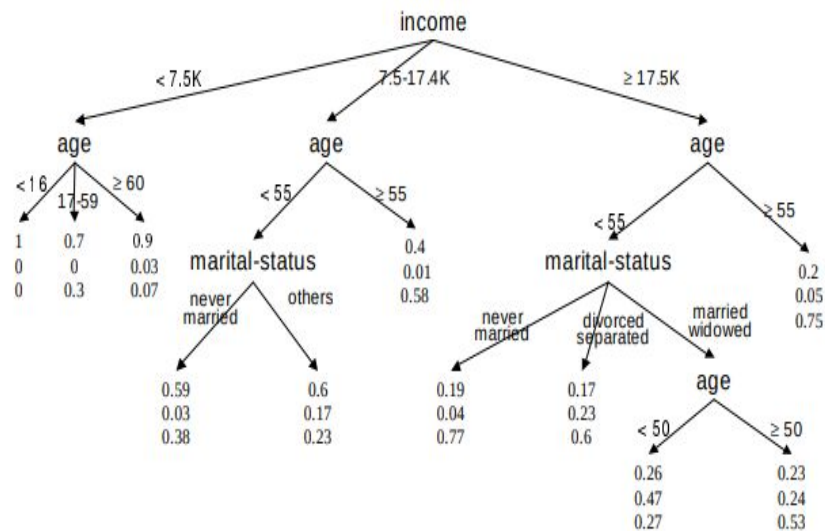
Use Graphical Model to detect correlation

- We can build a Bayesian network which consists of two component
 - a DAG whose nodes correspond to A_1, \dots, A_n , edges denote a direct dependency of A_i on its parents(A_i)
 - conditional probability distribution





(a)



(b)

Part III - Reactive Methods



Feedback based Histogram

- monitor queries on specified column gathering estimated and actual cardinalities
- create/refine maximum-entropy distribution at given condition

Max Entropy Principle

- We define every bucket as $\{l_i, r_i, f_i\}$, of which f_i is relative frequency.
- $H(R) = - \sum (m_i / \ln (m_i / h_i))$
- Use this principle to determine the boundaries

Meeting a Space Budget - Pruning



Figure 6: Set of bins before pruning

$h_1 = 4$	$m_1 = 0.2$
$h_2 = 2$	$m_2 = 0.1$
$h_3 = 3$	$m_3 = 0.1$
$h_4 = 2$	$m_4 = 0.2$
$h_5 = 3$	$m_5 = 0.3$
$h_6 = 4$	$m_6 = 0.05$
$h_7 = 3$	$m_7 = 0.05$
$max = 5$	$k = 7$

Meeting a Space Budget - Pruning

$$err(r_x, r_{x+1}) = h_x \left| \frac{m_x}{h_x} - \left(\frac{m_x + m_{x+1}}{h_x + h_{x+1}} \right) \right| + h_{x+1} \left| \frac{m_{x+1}}{h_{x+1}} - \left(\frac{m_x + m_{x+1}}{h_x + h_{x+1}} \right) \right| \quad (6)$$

The error equals zero if and only if two bins imply uniformity.

In this example merging bins 1,2 and 4,5 minimizes the error because:

$$err(r_1, r_2) = 4 \left(\left| \frac{0.2}{4} - \frac{0.2+0.1}{4+2} \right| \right) + 2 \left(\left| \frac{0.1}{2} - \frac{0.2+0.1}{4+2} \right| \right) = 0$$

$$err(r_4, r_5) = 2 \left(\left| \frac{0.2}{2} - \frac{0.2+0.3}{2+3} \right| \right) + 3 \left(\left| \frac{0.3}{3} - \frac{0.2+0.3}{2+3} \right| \right) = 0$$

Thank You !



Shanghai Meetup No.82 现场群

