

# CSE 562 Database Systems

## Query Processing: Cost Analysis

Some slides are based or modified from originals by  
*Database Systems: The Complete Book,*  
Pearson Prentice Hall 2<sup>nd</sup> Edition  
©2008 Garcia-Molina, Ullman, and Widom

*cse@buffalo*

UB CSE 562 Spring 2009

## Outline – Query Optimization

- Overview
- Relational algebra level
  - Algebraic Transformations
- Detailed query plan level
  - Estimate Costs
    - Estimating size of results
    - Estimating # of IOs
  - Generate and compare plans

UB CSE 562 Spring 2009

2

## Estimating cost of query plan

1. Estimating **size** of results
2. Estimating **#** of IOs

UB CSE 562 Spring 2009

3

## Estimating Result Size

- Keep statistics for relation R
  - $T(R)$  : # tuples in R
  - $S(R)$  : # of bytes in each R tuple
  - $B(R)$  : # of blocks to hold all R tuples
  - $V(R, A)$  : # distinct values in R for attribute A

UB CSE 562 Spring 2009

4

## Example

R	A	B	C	D
cat	1	10	a	
cat	1	20	b	
dog	1	30	a	
dog	1	40	c	
bat	1	50	d	

A: 20 byte string

B: 4 byte integer

C: 8 byte date

D: 5 byte string

$$T(R) = 5 \quad S(R) = 37$$

$$V(R,A) = 3 \quad V(R,C) = 5$$

$$V(R,B) = 1 \quad V(R,D) = 4$$

UB CSE 562 Spring 2009

5

## Size Estimates

For  $W = R1 \times R2$

$$T(W) = T(R1) \times T(R2)$$

$$S(W) = S(R1) + S(R2)$$

UB CSE 562 Spring 2009

6

## Size Estimate

For  $W = \sigma_{z=val}(R)$

$$S(W) = S(R)$$

$$T(W) = ?$$

UB CSE 562 Spring 2009

7

## Example

R	A	B	C	D
cat	1	10	a	
cat	1	20	b	
dog	1	30	a	
dog	1	40	c	
bat	1	50	d	

$V(R,A)=3$

$V(R,B)=1$

$V(R,C)=5$

$V(R,D)=4$

$$W = \sigma_{z=val}(R) \quad T(W) = \frac{T(R)}{V(R,Z)}$$

UB CSE 562 Spring 2009

8

## Assumption

Values in select expression  $Z = \text{val}$  are **uniformly distributed** over possible  $V(R,Z)$  values

## Alternate Assumption

Values in select expression  $Z = \text{val}$  are **uniformly distributed** over domain with  $\text{DOM}(R,Z)$  values

## Example

R	A	B	C	D
cat	1	10	a	
cat	1	20	b	
dog	1	30	a	
dog	1	40	c	
bat	1	50	d	

Alternate assumption

$V(R,A)=3$   $\text{DOM}(R,A)=10$

$V(R,B)=1$   $\text{DOM}(R,B)=10$

$V(R,C)=5$   $\text{DOM}(R,C)=10$

$V(R,D)=4$   $\text{DOM}(R,D)=10$

$W = \sigma_{Z=\text{val}}(R)$   $T(W) = ?$

## Example (cont'd)

$$C=\text{val} \Rightarrow T(W) = (1/10)1 + (1/10)1 + \dots \\ = (5/10) = 0.5$$

$$B=\text{val} \Rightarrow T(W) = (1/10)5 + 0 + 0 = 0.5$$

$$A=\text{val} \Rightarrow T(W) = (1/10)2 + (1/10)2 + (1/10)1 \\ = 0.5$$



## Solution #3: Estimate Values in Range

Equivalently:

$f \times V(R,Z)$  = fraction of distinct values

$$T(W) = [f \times V(Z,R)] \times \frac{T(R)}{V(Z,R)} = f \times T(R)$$

## Size Estimate

For  $W = R1 \bowtie R2$

Let  $X$  = attributes of R1

$Y$  = attributes of R2

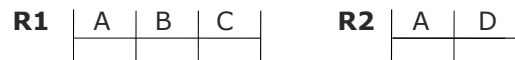
Case 1  $X \cap Y = \emptyset$

Same as  $R1 \times R2$

## Size Estimate

For  $W = R1 \bowtie R2$

Case 2  $X \cap Y = A$



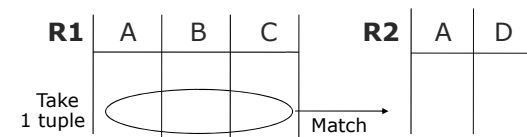
Assumption:

- $V(R1,A) \leq V(R2,A) \Rightarrow$  Every A value in R1 is in R2
- $V(R2,A) \leq V(R1,A) \Rightarrow$  Every A value in R2 is in R1

“containment of value sets”  
(justified by primary key – foreign key relationship)

## Computing $T(W)$

When  $V(R1,A) \leq V(R2,A)$



1 tuple of R1 matches with  $\frac{T(R2)}{V(R2,A)}$  tuples of R2

$$\text{So } T(W) = \frac{T(R2) \times T(R1)}{V(R2, A)}$$

## Size Estimate

- When  $V(R1,A) \leq V(R2,A)$

$$T(W) = \frac{T(R2) T(R1)}{V(R2,A)}$$

- $V(R2,A) \leq V(R1,A)$

$$T(W) = \frac{T(R2) T(R1)}{V(R1,A)}$$

[A is common attribute]

## In General

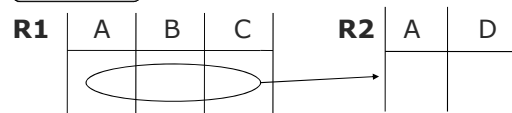
When  $W = R1 \bowtie R2$

$$T(W) = \frac{T(R2) T(R1)}{\max\{ V(R1,A), V(R2,A) \}}$$

## With Alternate Assumption

Values uniformly distributed over domain

Case 2



This tuple matches  $T(R2)/DOM(R2,A)$  so

$$T(W) = \frac{T(R2) T(R1)}{DOM(R2, A)} = \frac{T(R2) T(R1)}{DOM(R1, A)}$$

Assume the same

## In All Cases

- $S(W) = S(R1) + S(R2) - S(A)$  ← size of attribute A

## Size Estimate

For  $W = R \bowtie S$  with Multiple Attributes

$$X \cap Y = \{A, B\}$$

<b>R</b>	A	B	C	<b>S</b>	A	B	D
----------	---	---	---	----------	---	---	---

$$T(W) = \frac{T(R) T(S)}{\max\{V(R,A), V(S,A)\} \max\{V(R,B), V(S,B)\}}$$

## Size Estimate

For  $W = R \bowtie S \bowtie T$  with Multiple Attributes

$$T(W) = ?$$

## Using Similar Ideas...

We can estimate sizes of:

$\pi_{AB}(R)$  ... Section 16.4.2

$\sigma_{A=a \text{ AND } B=b}(R)$  ... Section 16.4.3

Union, intersection, diff ... Section 16.4.7

## Intermediate T,S,V Results

For complex expressions,  
we need intermediate T,S,V results

E.g.  $W = [\underbrace{\sigma_{A=a}(R1)}] \bowtie R2$   
Treat as relation U

$$T(U) = T(R1)/V(R1,A)$$

$$S(U) = S(R1)$$

Also need  $V(U, *)$  !!

## Intermediate T,S,V Results

To estimate **Vs**

E.g.,  $U = \sigma_{A=a}(R1)$

Say R1 has attributes A,B,C,D

$$V(U, A) =$$

$$V(U, B) =$$

$$V(U, C) =$$

$$V(U, D) =$$

## Example

<b>R1</b>	A	B	C	D
cat	1	10	10	
cat	1	20	20	
dog	1	30	10	
dog	1	40	30	
bat	1	50	10	

$$V(R1,A)=3$$

$$V(R1,B)=1$$

$$V(R1,C)=5$$

$$V(R1,D)=3$$

$$U = \sigma_{A=a}(R1)$$

$$V(U,A) = 1 \quad V(U,B) = 1 \quad V(U,C) = \frac{T(R1)}{V(R1,A)}$$

$V(D,U)$  ... somewhere in between

## Possible Guess

$$U = \sigma_{A=a}(R1)$$

$$V(U,A) = 1$$

$$V(U,B) = V(R,B)$$

## Intermediate T,S,V Results

For Joins  $U = R1(A,B) \bowtie R2(A,C)$

$$V(U,A) = \min \{ V(R1, A), V(R2, A) \}$$

$$V(U,B) = V(R1, B)$$

$$V(U,C) = V(R2, C)$$

"preservation of value sets"

## Example

$$Z = R1(A,B) \bowtie R2(B,C) \bowtie R3(C,D)$$

**R1**  $T(R1) = 1000$   $V(R1,A)=50$   $V(R1,B)=100$

**R2**  $T(R2) = 2000$   $V(R2,B)=200$   $V(R2,C)=300$

**R3**  $T(R3) = 3000$   $V(R3,C)=90$   $V(R3,D)=500$

UB CSE 562 Spring 2009

33

## Example

Partial Result:  $U = R \bowtie S$

$$T(U) = \frac{1000 \times 2000}{200}$$

$$V(U,A) = 50$$

$$V(U,B) = 100$$

$$V(U,C) = 300$$

UB CSE 562 Spring 2009

34

## Example

$$Z = U \bowtie R3$$

$$T(Z) = \frac{1000 \times 2000 \times 3000}{200 \times 300}$$

$$V(Z,A) = 50$$

$$V(Z,B) = 100$$

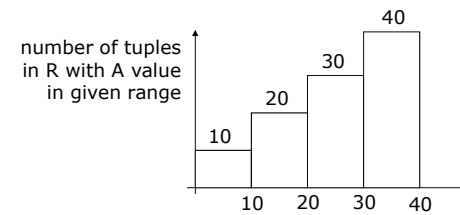
$$V(Z,C) = 90$$

$$V(Z,D) = 500$$

UB CSE 562 Spring 2009

35

## A Note on Histograms



$$\sigma_{A=a}(R1) = ?$$

UB CSE 562 Spring 2009

36

## Summary

- Estimating size of results is an “art”
- Don’t forget:  
Statistics must be kept up to date...  
(cost?)