

Introduction to Data Management Review

Paul G. Allen School of Computer Science and Engineering University of Washington, Seattle

December 6, 2024

Review

- Please fill out the course evals: <u>https://uw.iasystem.org/survey/297000</u>
- HW 7 is due today. NO LATE DAYS
- Monday, Dec. 9, Final Review session
 - CSE2 G10
 - 10:30 12:20 (we may finish earlier)



Review of this course

Relational Data Model and SQL

Data is stored in simple, flat relations



No prescription for the physical representation

First Normal Form

1NF

Physical Data Independence

- User writes SQL query:
 - Says what they want

- System responsible for optimizing SQL query
 - How to do it

Physical Data Independence is the main reason why relational model is the most widely used

Nested Loop Semantics



> May group by attributes, e.g. YEAR or expressions, e.g. YEAR/10

Only attributes or exrepssions mentioned in GROUPY may be used here...

SELECT FROM WHERE GROUP BY HAVING ORDER BY





SQL Aggregates



SQL Aggregates



Payroll

UserID	Name	Job	Salary
123	Jack	ТА	50000
345	Allison	ТА	60000
567	Magda	Prof	90000
789	Dan	Prof	100000

SQL Aggregates



Payroll

UserID	Name	Job	Salary
123	Jack	TA	50000
345	Allison	TA	60000
567	Magda	Prof	90000
789	Dan	Prof	100000

Magda

Allison

SQL: NULLs

Three-valued logic:

false = 0; unknown = 0.5; true = 1

$$x AND y = min(x,y);$$

x OR y = max(x,y);
not x = 1-x

Examples:

- true AND unknown =
- true OR unknown =
- unknown AND false =

unknown true false

SELECT FROM Table1 LEFT OUTER JOIN Table2 ON ...

INNER JOIN LEFT OUTER JOIN RIGHT OUTER JOIN FULL OUTER JOIN

Very useful for GROUP BY queries when we need aggregates on empty groups, e.g. count(*)=0

SQL: Witness or Argmin/Argmax

- SQL has the aggregates min(...) and max(...)
- SQL does not have argmin(...) or argmax(...)

SQL: Witness or Argmin/Argmax

- SQL has the aggregates min(...) and max(...)
- SQL does not have argmin(...) or argmax(...)
- Solution 1 using WITH:
 - Compute min or max in temporary table
 - Join main table with temp table to find argmin/argmax

SQL: Witness or Argmin/Argmax

- SQL has the aggregates min(...) and max(...)
- SQL does not have argmin(...) or argmax(...)
- Solution 1 using WITH:
 - Compute min or max in temporary table
 - Join main table with temp table to find argmin/argmax
- Solution 2 using self-joins:
 - Compute min/max from one copy of the table
 - Join with the other table in the HAVING clause

In the FROM clause: better use WITH

- In the FROM clause: better use WITH
- In the SELECT clause: must return single value

- In the FROM clause: better use WITH
- In the SELECT clause: must return single value
- In the WHERE clause:
 - EXISTS or NOT EXISTS
 - IN or NOT IN
 - ALL or ANY
 - They express mathematical quantifiers: ∀,∃

- In the FROM clause: better use WITH
- In the SELECT clause: must return single value
- In the WHERE clause:
 - EXISTS or NOT EXISTS
 - IN or NOT IN
 - ALL or ANY
 - They express mathematical quantifiers: ∀,∃
- Non-monotone* queries require subqueries
- * A query with an aggregate is non-monotone

The 5 basic operations:

- 1. Selection $\sigma_{condition}(S)$
- 2. Projection $\Pi_{attrs}(S)$
- 3. Join $R \bowtie_{\theta} S = \sigma_{\theta}(R \times S)$
- 4. Union ∪
- 5. Set difference –

Add renaming ρ , but we use variables instead

The 5 basic operations:

- 1. Selection $\sigma_{\text{condition}}(S)$
- 2. Projection $\Pi_{attrs}(S)$
- 3. Join $R \bowtie_{\theta} S = \sigma_{\theta}(R \times S)$
- 4. Union U

Monotone

5. Set difference – Non-monotone

Add renaming ρ , but we use variables instead

Two extended operator

 $\hfill \hfill \hfill$

• Group-by aggregate γ_{attr1,attr2,...,agg1,...}



When the query has subqueries then we need to unnest first

```
SELECT P.UserID
FROM Payroll P
WHERE not exists
   (SELECT *
    FROM Regist R
    WHERE P.UserID = R.UserID);
```

When the query has subqueries then we need to unnest first

```
SELECT P.UserID
FROM Payroll P
WHERE not exists
   (SELECT *
    FROM Regist R
    WHERE P.UserID = R.UserID);
```



When the query has subqueries then we need to unnest first

```
SELECT P.UserID
FROM Payroll P
WHERE not exists
   (SELECT *
    FROM Regist R
    WHERE P.UserID = R.UserID);
```





SELECT P.UserID
FROM Payroll P
EXCEPT
SELECT R.UserID
FROM Regist R;

When the query has subqueries then we need to unnest first



Review

Design Theory
The Database Design Process



ER Diagrams



ER Diagrams to Tables

• Each entity set \rightarrow a table

- Each relationship \rightarrow a table with two FKs
 - Except for many-one (or one-one): then add FK

• Each IS_A \rightarrow a FK



















Anomalies

The three types of anomalies

- Redundancy anomaly
- Update anomaly
- Deletion anomaly

All happen because $A \rightarrow B$, where A is no super-key

Definition:

• $A \rightarrow B$ holds if:

any 2 tuples that have same values of A attributes,

also have the same values in the B attribute

Always think about this definition!

Payroll

UserID	Name	Job	Salary
123	Jack	TA	50000
345	Allison	TA	60000
567	Magda	Prof	90000
789	Dan	Prof	100000

FD:

UserID \rightarrow Name, Job, Salary

Payroll

UserID	Name	Job	Salary
123	Jack	TA	50000
345	Allison	TA	60000
567	Magda	Prof	90000
789	Dan	Prof	100000

FD:

UserID \rightarrow Name, Job, Salary

SELECT *

FROM Payroll
WHERE Job = 'TA'

Payroll

UserID	Name	Job	Salary
123	Jack	TA	50000
345	Allison	TA	60000
567	Magda	Prof	90000
789	Dan	Prof	100000

FD:

UserID \rightarrow Name, Job, Salary

SELECT * FROM Payroll WHERE Job = 'TA'

UserID	Name	Job	Salary
123	Jack	TA	50000
345	Allison	TA	60000

Payroll

UserID	Name	Job	Salary
123	Jack	TA	50000
345	Allison	TA	60000
567	Magda	Prof	90000
789	Dan	Prof	100000

FD:

UserID \rightarrow Name, Job, Salary

SELECT *
FROM Payroll
WHERE Job = 'TA'

UserID	Name	Job	Salary
123	Jack	TA	50000
345	Allison	TA	60000

FDs:

UserID \rightarrow Name, Job, Salary Name \rightarrow Job Salary \rightarrow Job

- Goal: remove anomalies
- How:
 - Find set of attributes X such that $X \subsetneq X^+ \subsetneq$ [all-attrs]
 - Split relation into two relation
 - Remember to continue with both relations!

UID	Name	Phone	City
234	Fred	206-555-9999	Seattle
234	Fred	206-555-8888	Seattle
987	Joe	415-555-7777	SF

UID → Name, City

\sim

UID	Name	City	UID	Phone
234	Fred	Seattle	234	206-555-9999
987	Joe	SF	234	206-555-8888
			987	415-555-7777

- BCNF removes all anomalies, may lose some FDs
- 3NF keeps all FDs, may still have some anomalies

- BCNF removes all anomalies, may lose some FDs
- 3NF keeps all FDs, may still have some anomalies

R(City, State, Zip)

City,State \rightarrow Zip Zip \rightarrow State

- BCNF removes all anomalies, may lose some FDs
- 3NF keeps all FDs, may still have some anomalies

```
R(City, State, Zip)
```

City,State
$$\rightarrow$$
 Zip
Zip \rightarrow State

BCNF:

R1(Zip,State)We lost:R2(Zip,City)City,State
$$\rightarrow$$
 Zip

- BCNF removes all anomalies, may lose some FDs
- 3NF keeps all FDs, may still have some anomalies

R(City, State, Zip)

City,State
$$\rightarrow$$
 Zip
Zip \rightarrow State

3NF:

BCNF:

R1(Zip,State)We lost:R2(Zip,City)City,State \rightarrow Zip

R(City, State, Zip)

We have anomalies

Transactions

Two types of SQL workloads:

- Online Analytical Processing (OLAP)
 - Lots of joins, aggregates
 - Rarely any updates
 - Great for data analysis, decision support
- Online Transaction Processing (OLTP)
 - Lots of updates
 - Usually few joins or aggregates
 - Great for data-intensive applications (banking, ...)

Problem: concurrent updates may corrupt the DB

- Transactions: DB remains consistent
- ACID: A and I matter most. C is a consequence.
- Transactions slow down DBMS

Transactions

Using TXNs is easy:

BEGIN TRANSACTION ... COMMIT. (or ROLLBACK)

Implementing TXNs: must have ACID properties

Static database:

- A fixed set of elements: A₁, A₂, ...
- ATXN is a sequence of Read/Write operations

Dynamic database:

The set of elements may increase or decrease





Maybe this is easier to read:



Things to know:

- Serial Schedule
- Serializable Schedule
- Conflict Serializable Schedule
- What happens in a static v.s. a dynamic database



Locks, 2PL, Strict 2PL



Locks, 2PL, Strict 2PL


Locks, 2PL, Strict 2PL



Locks, 2PL, Strict 2PL

- Strict 2PL ensures conflict serializability
 - In particular, it ensures serializability
- But only in a static database

 In a dynamic database need to handle phantoms in order to ensure serializability

- Shared Locks, or Read Locks:
 - Many TXNs can hold a Read Lock

- Exclusive Locks, or Write Locks:
 - Only one TXN can hold a Write Lock
 No other TXN can hold either a Read or Write Lock

L(A) replaced by either S(A) or X(A)

Weaker Isolation Levels



(Review in class what they are)

The Query Engine

The Query Engine



Physical Operators

Join ⋈

- Nested loop join
- Hash-join
- Merge-join

Group-by γ

- Nested loop group-by
- Hash-based group-by
- Sort-based group-by

Selection, projection σ , Π

On-the-fly

Query rewrite rules:

- Selection pushdown: $\sigma_C(R \bowtie S) = \sigma_C(R) \bowtie S$ when *C* refers only to *R*
- Join associativity: $R \bowtie (S \bowtie T) = (R \bowtie S) \bowtie T$
- Many, many others

Basic statistics

- T(R) = number of tuples
- V(R,A) = number of distinct values in R.A
- Histograms

Basic estimation formulas

$$EST[\sigma_{A=const}(R)] = \theta_{A=const} \cdot T(R) = \frac{T(R)}{V(R,A)}$$

$$\theta_{C_1 \text{ and } C_2} = \theta_{C_1} \cdot \theta_{C_2}$$

$$EST[R \bowtie_{A=B} S] = \frac{T(R) \cdot T(S)}{\max(V(R,A),V(S,B))}$$

Assumptions:

- Uniformity
- Independence
- Containment of values
- Preservation of values

```
SELECT *
FROM Payroll x, Regist y, Brand z
WHERE x.UserID = y.UserID
and y.car = z.car
and x.Job = `TA';
```

Est(Q) =

 $T(Payroll) \cdot T(Regist) \cdot T(Brand)$

 $\max(V(Payroll, UserID), V(Regist, UserID)) \cdot \max(V(Regist, car), V(Brand, car)) \cdot V(Payroll, Job)$

(In class: discuss preservation of values)

Memory Hierarchy



Credit: https://15445.courses.cs.cmu.edu/fall2023/

December 6, 2024

- The unit of disk read or write is a block
- Once in main memory, we call it a page
- Block size is fixed. Typically, 4k or 8k or 16k
- Sequential access much faster than random access

- Index = an auxiliary file that facilitates faster access to the data
- Usually a B+ tree, but can also be a hash-table

What do these commands do?

- CREATE INDEX Idx1 ON Payroll(Name)
- CREATE INDEX Idx2 ON Payroll(Salary, Job)
- CREATE INDEX Idx3 ON Payroll(Job, Salary)

Multi-attribute Index

CREATE INDEX Idx1 on Payroll(job, salary);

SELECT *

FROM Payroll

WHERE job='TA';

SELECT *

FROM Payroll

WHERE salary='50000';

(Discussed in class)

SELECT *
FROM Payroll
WHERE job=`TA'
and salary=`50000';

SQL++

(we just finished it)

We covered a lot of material this quarter!

- Details: show your mastery on the final
- High-level concepts: remember them throughout your career

Thanks for a great quarter! (Come on Monday to the final review!)