

Final Review Session

Relational Algebra

- projection (Π) ex. $\Pi_{id, citation}(\text{stops})$
- selection (σ) ex. $\Pi_{id, citation}(\sigma_{age < 50}(\text{stops}))$
- renaming (ρ) ex. $\rho_{\text{stops_data}(id \rightarrow \text{person_id})}(\text{stops})$
- cartesian / cross product (\times) ex. $\sigma_{\text{stops.location} = \text{zips.location}}(\text{stops} \times \text{zips})$

SQL Basics	SELECT [DISTINCT] (5)
FROM	(1)
WHERE	(2)
GROUP BY	(3)
HAVING	(4)
ORDER BY	(6)
LIMIT	(7)

Joins

- theta join (\bowtie_θ): cartesian \rightarrow selection condition
 - ex. $A \bowtie_{A.\text{column 2} > B.\text{column 2}} B$
- equi join: cartesian \rightarrow equality condition
 - ex. $A \bowtie_{A.\text{column 2} = B.\text{column 2}} B$
- natural join: when two relations share same attribute; dupes removed
 - ex. $A \bowtie B$ (will not have two location columns)

ex.

equi join

equi join

rename
student \rightarrow a
teachers \rightarrow b

$\Pi_{name, instructor_for}(\sigma_{class = instructor_for}$
 $(student \times \sigma_{b.area = "datascience"}(teachers)))$

Subqueries

Def: Parenthesized SQL statement used within an outer SQL statement

- As a scalar (singular value)
- To check if a query is not empty via the EXISTS clause

Window Functions

Def: fxn that uses values from more than one row to produce a value (potentially diff) for each row

\hookrightarrow agg functions similar but produce single value for all rows

uses: cumulative sum, determine rank within group, normalize a row using a group avg

3 components of window function:

- 1 PARTITION BY: groups the rows so that functions will be calculated within these groups instead of the entire set of rows
- 2 ORDER BY: sort rows in window frame if the order of rows is important (e.g. running totals)
- 3 RANGE: specify the window frame's relation to current row

ex.)

```
SELECT student_name, order_id, sum(price)
OVER (PARTITION BY student_id) AS total_spent
FROM Orders
ORDER BY order_id;
```

ORDER BY order_id;

Summary

Views

↳ virtual: output not stored, computed on demand

CTEs: used within larger queries

materialized: stored on disk, periodically refreshed

Strings

Regex

```
REGEXP_REPLACE (source, pattern, replacement [, flags])
```

↑ optional

"g": all matches

Substrings

```
SUBSTRING(string, start, length)
```

```
SUBSTR(string, n, m)
```

Other

```
STRPOS(string, substring)
```

↳ ex. STRPOS('smickerdoodle', 'doodle')

Sampling

```
ORDER BY RANDOM LIMIT n;
```

```
TABLESAMPLE BERNoulli(p);
```

```
TABLESAMPLE system(p);
```

↳ by memory page NOT row

DML (Data manipulation Language)

Insertion

INSERT INTO Relation VALUES (<list>)

↳ ex. INSERT INTO stops VALUES (5000, 'Asian', 24, 'West Oakland')

ex. INSERT INTO stops (id, race, age, location) VALUES (5000, 'Asian', 24, 'West Oakland')

INSERT INTO Relation (<Subquery>)

Updating

UPDATE Relation

SET < list of new column value assignments >

WHERE < condition >

↳ ex.) UPDATE stops

SET age = 18

WHERE age IS NULL ;

Deleting

DELETE FROM < relationname >

DELETE FROM < relationname > WHERE < condition >

↳ ex. age < 18

DDL (Data Definition Language)

CREATE: defines schema

DROP: deletes both schema and instance

ALTER: redefines schema

ex.) CREATE TABLE zips(

location VARCHAR(20)

zipcode INTEGER

);

Constraints

PRIMARY KEY vs. UNIQUE

- there can only be one primary key
- there can be many UNIQUE constraints
- NO attribute of PK can be null
- Attributes within UNIQUE can be null, can be tuples w/ null values for UNIQUE attributes

Foreign Keys

a field in one table, that refers to the PRIMARY KEY in another table

↳ ex.) CREATE TABLE Cast_info

(person_id INT FOREIGN KEY REFERENCES Actor(actor_id))

↳ ex.) location VARCHAR(20) NOT NULL,

arrest BOOLEAN DEFAULT False,

PRIMARY KEY (stopID)

UNIQUE (personID, stopTime)

Performance Tuning

Indexing: speeds up reads on a specific key so queries are faster, reduce scans

CREATE INDEX nameIdIndex ON Actor(name,id);

DROP INDEX idNameIndex;

high cardinality = high number of distinct values

↳ low cardinality good when you know data is clustered on the indexed attribute

↳ **CLUSTER** stores ON ageIndex

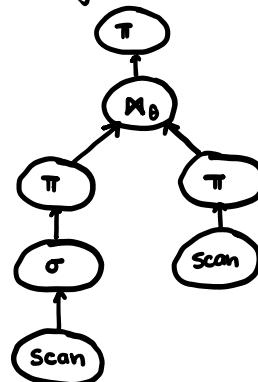
Query Optimization

useful for reducing the size of relational instances and avoiding blow-up of joins

predicate pushdown: pre-select rows before join

projection pushdown: pre-select columns join

ex.)
SELECT c.name, s.name, s.prevMargin
FROM Candidate AS C
INNER JOIN States AS S
ON c.homeStateID != s.sid
WHERE s.population > 5000000 AND



Performance Joins

nested loop join

↳ for loops, for every tuple of R and S check if it matches

sort-merge join

↳ walk down the two runs in their sorted order, merge matching tuples together

hash join

↳ tuples hashed into each bucket should be read from both R and S at same time

Data Preparation

gaussian outliers : removes values more than K z-scores away

windorizing : replace outlier values with those just before those removed on a specified boundary

median absolute deviation (MAD) : $MAD(x) = \text{median}(|x_i - \tilde{x}|)$

$$1\text{ SD} = 1.4826 \text{ MADs}$$

Hampel x84 : filtering outliers $K \cdot 1.4826 \text{ MADs}$ away

levenshtein distance : count # of insertions, deletions, mutations (replacements)
btw strings

blocking + matching :

↳ blocking : q-gram blocking where we block for substrings of length q
large values of q leads to bad recall
small values of q leads to decreased precision

↳ matching : clustering entities into equivalent classes, compute distance, and
clustering class values by minimizing distance

imputation : method to fill in NULLS

↳ interpolation

★ run
run start
run end
run size
run rank

Multivitamin #2 Q5

`int` stricter than `float`

Region	Route	Avg_pokemon_HP	Number_caught	Legendary_caught?
Kalos	21	75	593	0
Hoenn	Route 121	62	3960	0
Alola	11	81	2363.0	1
NULL	115	60	NULL	0
Johto	32	94	4735	0
Hoenn	Route 108	121	1294	1
NULL	Route 34	55	2735.0	0
Kanto	8	78	NULL	0
Sinnoh	205	83	NULL	1
Kalos	17	70	9948	0

Q5.2

3 Points

We decide to use MDL to determine the optimal type for column 4, `Number_caught`. Assume floats are 32 bits, ints are 16 bits, booleans are 1 bit, and a character is 8 bits.

Q5.2.1

Suppose the initial type chosen is integer. What would be the result of the sum in the MDL equation?

144

Q5.2.2

Now we choose the float type. What would be the result of the sum in the MDL equation?

224

5.2.1

$$1b + 1b + 32 + 0 + 1b + 1b + 32 + 1b$$

$$= 144$$

5.2.2

$$32 + 32 + 32 + 0 + 32 + 32 + 32 + 0 + 0 + 32$$

$$= 224$$

choose integer since MDL wants to minimize calculated sum

Date	Week	Client	Project	Time_spent (Hrs)
Mon 3/8/21	7	Cal	Website	3
Mon 3/8/21	7	UCLA	Newsletter	2
Tues 3/9/21	7	Cal	Website	2
Tues 3/9/21	7	UCLA	Newsletter	4
Wed 3/10/21	7	Cal	Hiring	15
Wed 3/10/21	7	UCLA	Website	2
Thurs 3/11/21	7	Cal	Merch Launch	4
Thurs 3/11/21	7	UCLA	Hiring	10
Fri 3/12/21	7	Cal	Newsletter	2
Fri 3/12/21	7	UCLA	Scheduling	24
Mon 3/15/21	8	Cal	Scheduling	37
Mon 3/15/21	8	UCLA	Merch Launch	2

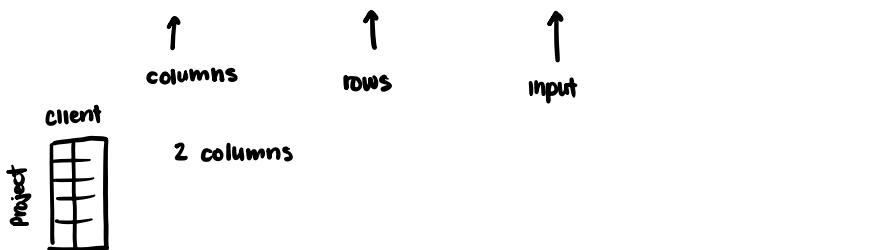
Q4.3.1

We execute the following pivot operation:

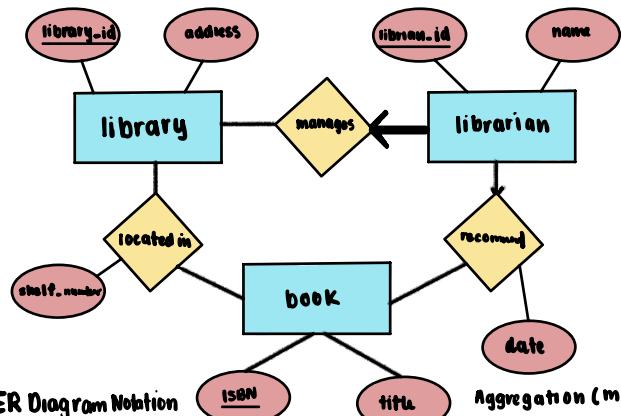
`Pivot(Client, Project, Time_spent(Hrs), sum)`

How many columns will be in the resulting pivot table? (NOTE: Do not treat row labels as a column)

`Pivot(<Column 1>, <Column 2>, <column3>, <Agg function>)`



Entity Relationship Diagram



can = "optional"

Functional Dependency + Normalization

$X \rightarrow Y$ iff: $t.X = u.X \rightarrow t.Y = u.Y$

$X \rightarrow Z$ iff: $t.X = u.X + t.Y = u.Y \rightarrow t.Z = u.Z$

Normalization splitting/decomposing relations minimize redundancy

Functional dependencies constraints b/w two sets of attributes

More mongo

db.tabu.aggregate([

```

    { $group: { _id: "field", newField: { $avgFunc: "$field" } },
    { $match: { column: { $gte: 150000 } },
    { $sort: { field: -1 } },
    { $project: { field: 0 } },
    { $limit: 10 }
  ]
  
```

Aggregation (ML): db.paintings.aggregate([

```

    { $group: { _id: "artistId", totalPrice: { $sum: "$purchasePrice" } },
    { $match: { "totalPrice": { $gte: 10000000 } },
    { $sort: { "totalPrice": -1 } },
  ]
  
```

DBMS	MDB
Database Relation Row Column	Database collection Document Field

Updates (ML): db.paintings.updateMany([

```

    { "purchasePrice": { $gt: 5000000 } },
    { $inc: { "purchasePrice": 1000000 },
    { $set: { "status": "appreciated" } }
  ]
  
```

ER Diagram Notation

Edge / Arrow	Bounds	Use Case
—	[0-many]	"not none, up to many"
—	[1-many]	"at least one"
—→	[0-1]	"at most one"
→	[1-1]	"exactly one"

Olap in PostgreSQL

year	month	sales
2022	January	100
2022	February	150
2022	March	200
2023	January	180
2023	February	120

SELECT year, month, sum(sales)

FROM sales-data
GROUP BY ROLLUP(year, month);

SELECT year, month, sum(sales)

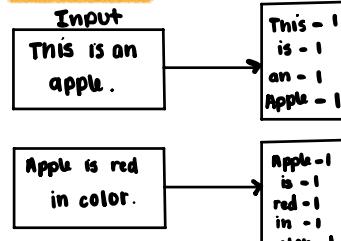
FROM sales-data
GROUP BY CUBE (year, month);

ROLLUP - generates a result that shows aggregates for a hierarchy of values

CUBE - generates a result set that shows all possible combos of aggs for the specified group by elements

year	month	sum
2022	January	100
2022	February	150
2022	March	200
2022	NULL	450
2023	January	180
2023	February	120
2023	NULL	300
NULL	NULL	750

MapReduce



shuffle

This - 1
is - 1
an - 1
Apple - 1
Apple - 1
red - 1
in - 1
color - 1

Reducer

This - 1
is - 2
an - 1
Apple - 2
red - 1
in - 1
color - 1

Output

This - 1
is - 2
an - 1
Apple - 2
red - 1
in - 1
color - 1

MongoDB

Q: list of mems for movieId 192

db.movies.agg([{\$unwind: "\$cast"},

{ \$match: { movieId: { \$eq: 192 } } }])

Q: movies w/ rating higher than 7

db.ratings.find({ "ratingInRating": { \$gt: 7 } })

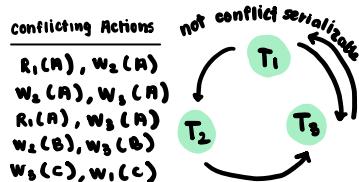
Q: Leo DiCaprio movies

db.movies.find({ cast: "Leo DiCaprio" })

Transactions

Time	1	2	3	4	5	6	7	
T1	R(A)							w(C)
T2		w(A)	w(B)					
T3			w(A)	w(B)	w(C)			

conflicting Actions
R₁(A), W₂(A)
W₂(A), W₃(A)
R₁(A), W₃(A)
W₂(B), W₃(B)
W₃(C), W₁(C)



Vocab

Shared(S) Lock: allow mul trans to read; prevents writing
Exclusive(X) Lock: exc. rights to read + write

Strict 2-Phase Locking: growing (acquire) + shrinking (release)

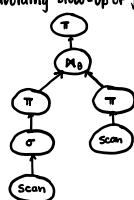
Deadlock: when two transactions wait for held resources

Pros (+)	Cons (-)
• remote redund. • min update / delete anomalies	• joins are costly

Query Optimization

useful for reducing the size of relational instances and avoiding blow-up of joins
predicate pushdown: pre-select rows before join
projection pushdown: pre-select columns join

ex.)
`SELECT c.name, s.name, s.population
 FROM candidate AS c
 INNER JOIN states AS s
 ON c.homestateID != s.sid
 WHERE s.population > 800000 AND`



3 components of window function:

- 1 **PARTITION BY:** groups the rows so that functions will be calculated within those groups instead of the entire set of rows
 - 2 **ORDER BY:** sort rows in window frame if the order of rows is important (e.g. running totals)
 - 3 **RANGE:** specify the window frame's relation to current row
- ex.)
`SELECT student_name, order_id, sum(price)
 OVER (PARTITION BY student_id) AS total_spent
 FROM orders
 ORDER BY order_id;`
- RANK() OVER(PARTITION BY ORDER BY DESC) AS rank

Summary

Views

- ↳ **virtual:** output not stored, computed on demand
- ↳ **CTEs:** real within larger queries
- ↳ **materialized:** stored on disk, periodically refreshed

ACID Principles: properties of database transactions

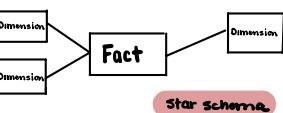
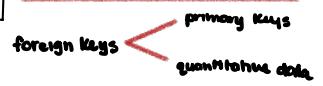
Atomicity: "all or nothing" nature of transactions

Consistency: ensures valid states (i.e. integrity constraints like uniqueness, PKs)

Isolation: ensures that the concurrent execution of diff transactions leaves database in same state as if transactions executed sequentially

Durability: ensures that once transaction is committed, effects → permanent even in the case of system failure

Physical models of Data Cubes



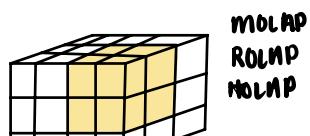
- still one fact table
- many dimension tables
- dimension table connected to others

Data Cube OLAP Queries

OL (Analytical) P: reads / sums large volumes of multidimensional data

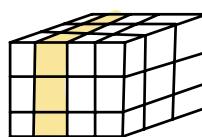
OL (Transactional) P: read, insert, update, delete data; smaller storage

↳ each element in the cube map the unique intersection of all combo of values along dim.

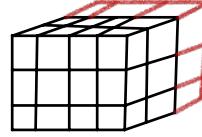


Dice: get range partition on one or more dimensions (subcube)

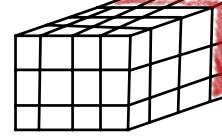
OLAP systems fundamentally provide comprehensive, summarized views of data in data warehouses. Sum will happen across multiple dimensions of data...



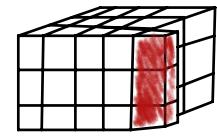
Slice: to get equality condition on a dimension



roll-up: decrease granularity to sum at higher level of dim hierarchy
Summarize



drill-down: increase granularity to sum at low level of a dim hierarchy
lower level



CROSS-TAB: pivot table, dist of 1 var relative to another

Relational Algebra

Union : R ∪ S

Natural Join : R ⋈ S

S

- **projection (π)** ex. $\pi_{id, citation}(\text{stops})$ Push ↓
 • **selection (σ)** ex. $\sigma_{id, citation \mid \text{age} < 50}(\text{stops})$
- **renaming (ρ)** ex. $\rho_{\text{stops} \rightarrow \text{data}(id \rightarrow \text{person.id})}(\text{stops})$ (stop=)
- **cartesian / cross product (x)** ex. $\sigma_{\text{stops}, \text{location}} = \text{zips}, \text{location}$ (stops × zips)

DML (Data manipulation Language)

Insertion

`INSERT INTO relation VALUES (<list>)`

↳ ex. `INSERT INTO stops VALUES (5000, 'Asian', 24, 'West Oakland')`

Data Models

- **data frame -** row / col labels
- **matrix -** no labels
- **relations -** column labels

Tensors - generalization of matrix

- Rank 0, scalar
- Rank 1, vector
- Rank 2, matrix
- Rank 3, tensor

Updating

`UPDATE Relation`

`SET < list of new column value assignments >`

`WHERE < condition >`

↳ ex. `UPDATE stops`

`SET age = 10`

`WHERE age IS NULL ;`

Deleting

`DELETE FROM < relationname >`

`DELETE FROM < columnname > WHERE < condition >`

Spreadsheets conceptual model

↳ ex. `age < 18`

A spreadsheet workspace comprises many sheets

- each sheet has cells
- a spreadsheet is structured around cells
- cells contain one or more of:

↳ value

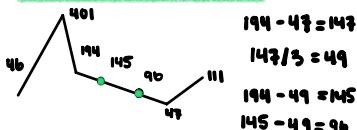
↳ formula (arithmetic / special functions)

Access control

`GRANT [privileges] ON object TO users [WITH GRANT OPTION]`

`REVOKE privileges ON object TO users [RESTRICT] CASCADE`

General Interpolation



Functional dependencies

X	Y	→	Z
Address	SSN		phone number
10 Green	123-456-789		(201) 233-1456
10 Green	123-456-789		(201) 123-3439
431 Purple	987-654-321		(145) 241-2131
431 Purple	987-654-321		(312) 123-1287
...

↓ ↓

Address	SSN
10 Green	123-456-789
431 Purple	987-654-321
...	...

SSN	phone number
123-456-789	(201) 233-1456
123-456-789	(201) 123-3439
987-654-321	(145) 241-2131
987-654-321	(312) 123-1287
...	...

Update anomalies:

If we update address of person with phone number (201) 233-1456

... there will be two addresses of that person

Deletion anomalies:

If we delete the phone number for a person

... we need to check another phone number exists for said person or we lose address info (dangling points)

K 10^3 bytes clust. index ↑ sort time, order, min/max ↓ writing
 M 10^6 views - memory - every reference
 G 10^9 mat. views - disk - query must make (low recom)
 T 10^{12} (te - memory - same as query (visual))
 P 10^{15} table - disk - permanent data
 E 10^{18}
 page = 8kb

hierarchy - granularity lvl (t:me = mins/hrs, days)
 numer. gran. - sign. digits, roll up → coarser (up in hier)
 ↳ less specific
 ✓ delete from child
 X delete from parent (unless cascade)

Database = records + record - size
 Pages needed = $\frac{\text{usable space per page}}{\text{bytes/records}}$ = $\lambda \rightarrow \frac{\text{records}}{\lambda}$
 df.melt = unpivot df.pivot (x, y, val)

Sparse = many NaNs

Epoch - 4 byte int
 1 int + 16 = 2 bytes
 32 = 4 bytes
 64 = 8 bytes
 String - timestamp - 20 bytes double 54 = 8 bytes
 (1 byte for each spot) page = 8 kb
 SQL timestamp = 8 bytes
 Extract (Year from timestamp) = 1 letter
 = 1 byte
 TZ - timestamp (epoch) at Time Zone 'UTC' as timestamp_utc
 String :: int → convert to int
 int 32 = $2^{32} = 4 \cdot 3 \times 10^9$ bytes = 28 bytes
 16 = 65K 64 = 10^{19} 8 bits

Nested - small table, ind on join col X big X with index (n^2)
 NOT RECOMMENDED (!=) small table
 Sort-merge - medium+ table, data sorted/index both on join
 ↳ join needs to sort
 Hash - large tables, no index, join join X range <, >

order by random() expensive (full table)
 tablesample Bernoulli(p) factor (p for each row)
 tablesample System(p) factor (selects whole page)

DF - rows, cols have labels
 Mat = no labels
 relational col. labels
 sum c - scalar (univ)
 1 - vector
 2 - matrix X
 3 - tensor (its all)

MAD = median of deviations
 Hampel = $1.4826 \cdot \text{MAD}$
 2 Hampel = 2.0
 MDL - takes default if possible
 percen - inc (0.5) within group (order by devicetime)

Mongo
 .find({ "key": "value" }, ...)
 \$match: {
 "lat": { "\$gte": 37, "\$lte": 21 } }
 \$group ... arg_star: { "\$avg": "\$stars" }
 .updateMany() → '\$text'
 '\$search': "a b c"
 '\$set': { "field": "val" }
 \$all, \$in, \$exists, \$unwind, \$lookup
 \$insertMany \$updateMany w/ \$set
 \$lookup('value', where, return)

\$group: {
 _id: '\$device',
 count: { \$sum: 1 } }
 Conflict = order is not valid (T₂ between T₁, functions)
 T₁ [RJ WP] S T₁ [] [RJ WP] S T₂ [RP WP]
 T₂ [] [RP WP] S if conflict
 cycles = unserializable

Serial
 T₁ [] [RJ WP] S T₁ [] [RJ WP] S Z
 T₂ [RP WP] T₂ [RP WP]

T₁ [] [RJ WP] T₁ [] [RJ WP] X
 T₂ [RP WP] T₂ [RD RJS WJS WP]

but serializable

Conflict = diff trans, same object, one is write
 i.e.: T₁ → R1(x) T₂ → W2(x) T₁ → W1(x)
 DBT from { \$ref: { 'name': ' } }
 roll up = more general day → month
 \$ = anchor (lexical)
 mult X by 8 → X = X < < 3
 2³ splits Shuffled
 cont divide

rollup(yrs, mon, val) = 2023 Jun 100
 n+1 2023 Feb 200
 2023 null 300

cube(yrs, mon, val) = 2023 Jun 100
 2ⁿ 2023 Feb 200
 2023 null 300
 2024 Jan 1000
 null Jan 1200

cube math = $(\# \text{unique} + 1) \cdot (\# \text{unique} + 1)$

serial = ordered = can do all T₁ x y z
 and T₂ x y z

same type = M, T

col. label = R, df

row reference by address: DF, M, Γ, S

Res. sample = SRS n big O, n rows
total not known (streams)

Ex. res =