LO3: SQL: Intermediate

CS3200 Database design (fa18 s2)

https://northeastern-datalab.github.io/cs3200/

Version 9/13/2018

Announcements!

- It is ok to make mistakes in class. Making mistakes in class is the best thing that can happen to you. You learn and will never make it again.
 - "Create a Culture in Which It Is <u>Okay to Make Mistakes</u> and <u>Unacceptable Not to Learn</u> <u>from Them</u>" ... "Recognize that mistakes are a natural part of the evolutionary process."
 ... "Don't feel bad about your mistakes or those of others. Love them!"
 Ray Dalio. Principles. 2017
 - LIFE & WORK PRINCIPLES RAY DALIO

- Continue bringing your name tags
- Thanks for posting Q&As on Piazza!
- Make use of office hours (OHs)
- Steve Ballmer: another attempt

OHs (Office Hours)

We are here to help. Reach out to us!

OFFICE HOURS

If you don't understand something, please ask questions. We love questions. One of the benefits of attending a university as opposed to reading a book is that you get to interact with faculty, TAs, and your peers. Please reach out to us with your questions on our Plazza site or stop by for office hours.

cs3200fa18s2

	Sun 9/9	Mon 9/10	Tue 9/11	Wed 9/12	Thu 9/13	Fri 9/14	Sat	9/15
Delli I								
9am		······································						
10am								
11am								
12pm		11:45 - 1:25p Class			11:45 - 1:25p Class			
10m		Richards Hall 235			Richards Hall 236			

2pm		2p - 3p OH Wolfgang						
3pm								
4pm								
5pm								
8pm		6p - 8p OH Niklas			6p - 7p OH Nikias			
7pm		Smedemark- Margulies			462 WVH			
		462 WVH						
8pm								
9pm								

Table Alias (Tuple Variables)



Person (<u>pName</u>, address, works_for) University (<u>uName</u>, address)

SELECTDISTINCT pName, addressFROMPerson, UniversityWHEREworks_for = uName

Table Alias (Tuple Variables)





Column Alias (rename attributes)



Person (<u>pName</u>, address, works_for) University (<u>uName</u>, address)



Product (<u>pName</u>, price, category, manufacturer) Company (<u>cName</u>, stockPrice, country)



Product

PName	Price	Category	Manufacturer
Gizmo	\$19.99	Gadgets	GizmoWorks
Powergizmo	\$29.99	Gadgets	GizmoWorks
SingleTouch	\$149.99	Photography	Canon
MultiTouch	\$203.99	Household	Hitachi

Company

CName	StockPrice	Country
GizmoWorks	25	USA
Canon	65	Japan
Hitachi	15	Japan

Q: Find all US companies that manufacture products in the 'Gadgets' category!

SELECT cName FROM WHERE

Product (<u>pName</u>, price, category, manufacturer) Company (<u>cName</u>, stockPrice, country)



Product

PName	Price	Category	Manufacturer	
Gizmo	\$19.99	Gadgets	GizmoWorks	
Powergizmo	\$29.99	Gadgets	GizmoWorks	
SingleTouch	\$149.99	Photography	Canon	
MultiTouch	\$203.99	Household	Hitachi	

Company

CName	StockPrice	Country
GizmoWorks	25	USA
Canon	65	Japan
Hitachi	15	Japan

Q: Find all US companies that manufacture products in the 'Gadgets' category!

SELECT	cName
FROM	Product P, Company
WHERE	country = 'USA'
and	P.category = 'Gadgets'
and	P.manufacturer = cName



Cname
GizmoWorks
GizmoWorks

Product (<u>pName</u>, price, category, manufacturer) Company (<u>cName</u>, stockPrice, country)



Product

PName	Price	Category	Manufacturer
Gizmo	\$19.99	Gadgets	GizmoWorks
Powergizmo	\$29.99	Gadgets	GizmoWorks
SingleTouch	\$149.99	Photography	Canon
MultiTouch	\$203.99	Household	Hitachi

Company

CName	StockPrice	Country
GizmoWorks	25	USA
Canon	65	Japan
Hitachi	15	Japan

Q: Find all US companies that manufacture products in the 'Gadgets' category!

SELECT	DISTINCT cName
FROM	Product P, Company
WHERE	country = 'USA'
and	P.category = 'Gadgets'
and	P.manufacturer = cName





Meaning (Semantics) of conjunctive SQL Queries

Conceptual evaluation strategy (nested for loops):

```
Answer = {}
for x_1 in R_1 do
for x_2 in R_2 do
.....
for x_n in R_n do
if Conditions
then Answer = Answer \cup \{(a_1,...,a_k)\}
return Answer
```

Meaning (Semantics) of conjunctive SQL Queries





Conceptual Evaluation Strategy

- Semantics of an SQL query defined in terms of the following conceptual evaluation strategy:
 - FROM: Compute the cross-product of relation-list.
 - WHERE: Discard resulting tuples if they fail qualifications.
 - SELECT: Delete attributes that are not in target-list.
 - If DISTINCT is specified, eliminate duplicate rows.
- This strategy is probably the least efficient way to compute a query! An optimizer will find more efficient strategies to compute the same answers.



Product (<u>pName</u>, price, category, manufacturer) Company (<u>cName</u>, stockPrice, country)

Q: Find all US companies that manufacture both a product below \$20 and a product above \$25.

SELECT	DISTINCT cName
FROM	
WHERE	



Product (<u>pName</u>, price, category, manufacturer) Company (<u>cName</u>, stockPrice, country)

Q: Find all US companies that manufacture both a product below \$20 and a product above \$25.

SELECTDISTINCT cNameFROMProduct as P, CompanyWHEREcountry = 'USA'andP.price < 20</td>andP.price > 25andP.manufacturer = cName

Wrong! Gives empty result: There is no product with price <20 and >25



Product (<u>pName</u>, price, category, manufacturer) Company (<u>cName</u>, stockPrice, country)

Q: Find all US companies that manufacture both a product below \$20 and a product above \$25.





Product (<u>pName</u>, price, category, manufacturer) Company (<u>cName</u>, stockPrice, country)

Q: Find all US companies that manufacture both a product below \$20 and a product above \$25. Returns companies





Product (<u>pName</u>, price, category, manufacturer) Company (<u>cName</u>, stockPrice, country)

Q: Find all US companies that manufacture both a product below \$20 and a product above \$25.

SELECT	DISTINCT cName
FROM	Product as P1, Product as P2, Company
WHERE	country = 'USA'
and	P1.price < 20
and	P2.price > 25
and	P1.manufacturer = cName
and	P2.manufacturer = cName



P1

F
C

Name	Price	Category	Manufacturer
Bizmo	\$19.99	Gadgets	GizmoWorks
-			

P2

PName	Price	Category	Manufacturer
Powergizmo	\$29.99	Gadgets	GizmoWorks

Company

CName	StockPrice	Country
GizmoWorks	25	USA

SELECT DISTINCT cName FROM Product as P1, Product as P2, Company WHERE country = 'USA' and P1.price < 20 And P2.price > 25 And P1.manufacturer = cName

and P2.manufacturer = cName



Inner Joins



Employee		_	Department		
LastName	DepartmentID		DepartmentID	DepartmentName	
 Rafferty <	31		31	Sales	
Jones	33		33	Engineering	
Steinberg	33		34	Clerical	
Robinson	34	/(35	Marketing	
Smith	34				

SELECT *FROMEmployee E, Department DWHEREE.DepartmentID = D. DepartmentID

	E.LastName	E.DepartmentID	D.DepartmentID	D.DepartmentName
N	Robinson	34	34	Clerical
	Jones	33	33	Engineering
\neg	Smith	34	34	Clerical
	Steinberg	33	33	Engineering
	Rafferty	31	31	Sales

Cross Joins: usually not what you want 🟵

34

34

Smith

Robinson



	Employee		Department	
	LastName	DepartmentID	Departmentl	D DepartmentName
	Rafferty	31	31	Sales
	Jones	33	33	Engineering
	Steinberg	33	34	Clerical
	Robinson	34	35	Marketing
	Smith	34		
-	FROM E	Employee E, De E.DepartmentIE	epartment D) = D. Department	
	E.LastName	E.Department	ID D.DepartmentID	D.DepartmentName
$\overline{}$	Rafferty Jones Steinberg Smith Robinson	31 33 33 34 34 34	31 31 31 31 31	Sales Sales Sales Sales Sales
	Jones Steinberg Smith Robinson	33 33 34 34 34	33 33 33 33 33 33	Engineering Engineering Engineering Engineering Engineering
V	Rafferty Jones Steinberg Smith Robinson	31 33 33 34 34 34	34 34 34 34 34 34	Clerical Clerical Clerical Clerical Clerical Clerical Clerical

35

35

Marketing

Marketing

Definitions (for job interviews?)

- An <u>equi-join</u> is a join in which the joining condition is based on equality between values in the common columns; common columns appear redundantly in the result table
- A <u>natural join</u> is an equi-join in which one of the duplicate columns is eliminated in the result table
- A <u>cross join</u> returns the Cartesian product of rows from tables in the join
 - (i.e. it will produce rows which combine each row from the first table with each row from the second table, that's usually *not* what you want)

Definitions (for job interviews?)

Equi-join			
E.LastName	E.DepartmentID	D.Department/D	D.DepartmentName
Robinson	34	34	Clerical
Jones	33	33	Engineering
Smith	34	34	Clerical
Steinberg	33	33	Engineering
Rafferty	31	31	Sales

Natural join

E.LastName	DepartmentID	D.DepartmentName
Robinson	34	Clerical
Jones	33	Engineering
Smith	34	Clerical
Steinberg	33	Engineering
Rafferty	31	Sales

Cross join

E.LastName	E.DepartmentID	D.DepartmentID	D.DepartmentName
Rafferty	31	31	Sales
Jones	33	31	Sales
Steinberg	33	31	Sales
Smith	34	31	Sales
Robinson	34	31	Sales
Rafferty	31	33	Engineering

Alternative JOIN Syntax

Employee

LastName	DepartmentID	
Rafferty	31	
Jones	33	
Steinberg	33	
Robinson	34	
Smith	34	

Department

DepartmentID	DepartmentName
31	Sales
33	Engineering
34	Clerical
35	Marketing

SELECT	*
FROM	Employee E, Department D
WHERE	E.DepartmentID = D. DepartmentID
AND	E.DepartmentID = 34

SELECT	*
FROM	Employee E JOIN Department D
ON	E.DepartmentID = D. DepartmentID
WHERE	E.DepartmentID = 34

E.LastName	E.DepartmentID	D.DepartmentID	D.DepartmentName
Robinson	34	34	Clerical
Smith	34	34	Clerical



NATURAL JOIN Syntax

Employee

LastName	DepartmentID	
Rafferty	31	
Jones	33	
Steinberg	33	
Robinson	34	
Smith	34	

Department

DepartmentID	DepartmentName
31	Sales
33	Engineering
34	Clerical
35	Marketing

SELECT	*
FROM	Employee E, Department D
WHERE	E.DepartmentID = D. DepartmentID
AND	E.DepartmentID = 34

SELECT *

FROM Employee E NATURAL JOIN Department D

WHERE E.DepartmentID = 34

LastName	DepartmentID	DepartmentName
Robinson	34	Clerical
Smith	34	Clerical

Syntax is not supported by all DBMS's and not recommended

(someone reading the query does not see the join keys. That makes is harder to debug queries or to modify & enhance them)



Using the Formal Semantics





What do these queries compute?











Returns $R \cap (S \cup T)$ if $S \neq \emptyset$ and $T \neq \emptyset$

Using the Formal Semantics





What do these queries compute?









Can seem counterintuitive! But remember conceptual evaluation strategy: Nested loops. If one table is empty -> no looping

Illustration with Python





"Premature optimization is the root of all evil." Donald Knuth (1974)

CREATE Tables & INSERT Values

(Relational Database) Schema





"Schema": describes the structure of data in terms of the relational data model.

A schema includes tables, columns, PKs, FKs, and other constraints

Product(<u>pname</u>, price, category, manufacturer) Company(<u>cname</u>, stockprice, country)

Product.manufacturer is FK to Company



create table Company (

);

create table Product (



create table Company (
 CName char(20)
 StockPrice int,
 Country char(20));

,

create table Product (
 PName char(20),
 Price decimal(9, 2),
 Category char(20),
 Manufacturer char(20),



create table Company (
 CName char(20) PRIMARY KEY,
 StockPrice int,
 Country char(20));

create table Product (
 PName char(20),
 Price decimal(9, 2),
 Category char(20),
 Manufacturer char(20),



create table Company (
 CName char(20) PRIMARY KEY,
 StockPrice int,
 Country char(20));

create table Product (
 PName char(20),
 Price decimal(9, 2),
 Category char(20),
 Manufacturer char(20),
 PRIMARY KEY (PName)



create table Company (
 CName char(20) PRIMARY KEY,
 StockPrice int,
 Country char(20));

```
create table Product (
    PName char(20),
    Price decimal(9, 2),
    Category char(20),
    Manufacturer char(20),
    PRIMARY KEY (PName),
FOREIGN KEY (Manufacturer) REFERENCES Company(CName) );
```



— Drop tables if they already exist

create table Company **DROP TABLE IF EXISTS Product;** IMARY KEY, DROP TABLE IF EXISTS Company; CName ch StockPrice Int Country char create table PName char(Price decimal(9, 2), Category char(20), Manufacturer char(20), PRIMARY KEY (PName), FOREIGN KEY (Manufacturer) REFERENCES Company(CName));



– Populate the tables

```
insert into Company values ('GizmoWorks', 25, 'USA');
insert into Company values ('Canon', 65, 'Japan');
insert into Company values ('Hitachi', 15, 'Japan');
```

insert into Product values ('Gizmo', 19.99, 'Gadgets', 'GizmoWorks'); insert into Product values ('PowerGizmo', 29.99, 'Gadgets', 'GizmoWorks'); insert into Product values ('SingleTouch', 149.99, 'Photography', 'Canon'); insert into Product values ('MultiTouch', 203.99, 'Household', 'Hitachi'); Aggregates
 Groupings
 Having

Aggregation

Car (<u>name</u>, price, maker)





SELECT	count(*)
FROM	Car
WHERE	price > 100

SQL supports several aggregation operations:

sum, count, min, max, avg

Except count, all aggregations apply to a single attribute

Aggregation



Car

SELECT	avg(price)
FROM	Car
WHERE	maker='Toyota'

	<u>Name</u>	Price	Maker	
_	M3	120	BMW	
_	M5	150	BMW	
	Prius	50	Toyota	
	Lexus1	75	Toyota	
	Lexus2	100	Toyota	

Database creates new attribute name (for SQLserver)

(<u>No column na</u>me) 75

Aggregation with rename





Car

Name	Price	Maker
M3	120	BMW
M5	150	BMW
Prius	50	Toyota
Lexus1	75	Toyota
Lexus2	100	Toyota

'n

2

Database creates *our*

new attribute name

Aggregation: Count Distinct





Car

Same as count(*)

We probably want to ignore duplicates: 20 SELECT count(DISTINCT maker) (No column name) FROM Car WHERE price > 100

Simple Aggregation 1/3



Purchase (product, price, quantity)

SELECTsum(price * quantity)FROMPurchase

SELECTsum(price * quantity)FROMPurchaseWHEREproduct = 'Bagel'

What do these queries mean?

Simple Aggregation 2/3

Purchase

	Product	Price	Quantity
	Bagel	3	20
	Bagel	2	20
-	Banana	1	50
_	Banana	2	10
_	Banana	4	10

$$3 * 20 = 60$$

 $2 * 20 = 40$

sum: 100

Database creates new attribute name

SELECTsum(price * quantity)FROMPurchaseWHEREproduct = 'Bagel'

Simple Aggregation 3/3



Purchase

	Product	Price	Quantity			
	Bagel	3	20	3	20	
	Bagel	2	20	2	20	
	Banana	1	50	- sum: 5	* sum: 40	= 200
	Danana	0	10			
	Danana	2	10			
_	Banana	4	10			
	Banana					

SELECT	<pre>sum(price) * sum(quantity)</pre>
FROM	Purchase
WHERE	product = 'Bagel'

(No column name	e)
200	

Grouping and Aggregation

308

Purchase

Product	Price	Quantity	
Bagel	3	20	Г
Bagel	2	20	
Banana	1	50	_
Ranana	2	10	
Danana			
Banana	4	10	

Product	TotalSales
Bagel	40
Banana	20

Notice: we use "sales" for total number of products sold

Find total quantities for all purchases with price over \$1 grouped by product.

From \rightarrow Where \rightarrow Group By \rightarrow Select



Purchase

	Product	Price	Quantity		Product	TotalSales
	Bagel	3	20		Bagel	40
	Bagel	2	20		Banana	20
	Banana	1	50			
	Banana	2	10			
	Banana	4	10		Select cont	ains
	 grouped attributes and aggregates 					
	4 SELECT product, sum(quantity) as TotalSales				Sales	
1	FROM Purchase					
2	2 WHERE price > 1					
	3 GROU	PBY pr	oduct			

Let's confuse the database engine



Purchase

Product	Price	Quantity
Bagel	3	20
Bagel	2	20
Banana	1	50
Banana	2	10
Banana	4	10

SELECT	product, quantity
FROM	Purchase
GROUP BY	product



The DB engine is confused, there is no single quantity for banana (it's an ill-defined query). It should thus return an error (only SQLite misbehaves and returns something, but which makes no sense). Please think this through carefully!