Updated 9/19/2022

## Topic 1: SQL L04: SQL intermediate

Wolfgang Gatterbauer

CS3200 Database design (fa22)

https://northeastern-datalab.github.io/cs3200/fa22s3/ 9/19/2022

#### Class warm-up

- Last class summary
- HW 1 groups: contacting best practice
- How to think about resources in the schedule?
- Thanks for posting and answering on Piazza!
- SQL today: aggregates, grouping, having clause
- SQL next time: nested queries

Setup PostgreSQL, L01-Introduction, L01-SQL

Setup Gradiance, L02-SQL, SAMS Ch 1-3, SDK 3.1-3.3, 3.4.4, 3.4.5, 3.9, 4.1, 4.4.5, 4.5.1

L03-SQL, SAMS 4 & 12

SAMS Ch 5-9, SDK 3.7

SAMS Ch 10-17, SDK 3.8

• SQL later: Nulls, outer joins, "witnesses" (traditionally students find this topic the conceptually most difficult)

## Big IMDB schema (Postgres)



## Small IMDB schema This is a far smaller movie database in ---- 300 our SQL folder



## Character types

- character varying = varchar(n)
  - variable-length character string (limited length)
- character(n) = char(n)
  - fixed length character string, blank padded (limited length)
  - in some databases faster than varchar for index look-ups (not postgres):
    - varchar may require more string manipulations; must perform some form of length checking

#### • text

- variable-length character string (unlimited length)
- Not part of SQL standard (e.g. Oracle does not know "text")

Source for tip: <u>https://www.postgresql.org/docs/12/datatype-character.html</u> Wolfgang Gatterbauer. Database design: <u>https://northeastern-datalab.github.io/cs3200/</u>

Name	Description
character varying( <b>n</b> ), varcha	ar ( <b>n</b> ) variable-length with limit
character( <b>n</b> ), char( <b>n</b> )	fixed-length, blank padded
text	variable unlimited length

#### Tip

There is no performance difference among these three types, apart from increased storage space when using the blank-padded type, and a few extra CPU cycles to check the length when storing into a length-constrained column. While character(n) has performance advantages in some other database systems, there is no such advantage in PostgreSQL; in fact character(n) is usually the slowest of the three because of its additional storage costs. In most situations text or character varying should be used instead.

Aggregates
 Groupings
 Having

## Simple Aggregation 1/3

Purchase (product, price, quantity)



#### Purchase

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

SELECT	<pre>sum(price * quantity)</pre>
FROM	Purchase

SELECTsum(price \* quantity)FROMPurchaseWHEREproduct = 'Bagel'

What do these queries mean?

(next slides)

## Simple Aggregation 2/3



#### **Purchase**

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

SELECTsum(price \* quantity)FROMPurchaseWHEREproduct = 'Bagel'



## Simple Aggregation 2/3

#### **Purchase**

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

Database creates new attribute name

SELECTsum(price \* quantity)FROMPurchaseWHEREproduct = 'Bagel'

SUM 100





## Simple Aggregation 3/3



#### **Purchase**

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

SELECTsum(price) \* sum(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
Banana	2	10	]		
Banana	4	10	-		

SELECTsum(price) \* sum(quantity)FROMPurchaseWHEREproduct = 'Bagel'

## Grouping and Aggregation



#### **Purchase**

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

# Q: Find total quantities for all purchases with price above \$1 grouped by product.

## Grouping and Aggregation



#### Purchase

			_		
Product	Price	Quantity		Product	TotalQuantities
Bagel	3	20		Bagel	?
Bagel	2	20		Banana	?
Banana	1	50			
Banana	2	10			
Banana	4	10	]		

## Q: Find total quantities for all purchases with price above \$1 grouped by product.

## Grouping and Aggregation



#### Purchase

Product	Price	Quantity		Product	TotalQuantities
Bagel	3	20		Bagel	40
Bagel	2	20		Banana	20
 Banana	1	50	_		
Banana	•				
Banana	2	10			
Banana	4	10			

## Q: Find total quantities for all purchases with price above \$1 grouped by product.

#### From $\rightarrow$ Where $\rightarrow$ Group By $\rightarrow$ Select



#### **Purchase**

Ρ	roduct	Price	Quantity		Product	TotalQuantities
B	agel	3	20		Bagel	40
B	agel	2	20		Banana	20
B	anana	1	50			
B	anana	2	10			
B	anana	4	10		Select cor	ntains
	<ul> <li>grouped attributes</li> <li>and aggregates</li> </ul>					
4 1 2 3	SELEC FROM WHER GROU	CT p F E p P BY p	oroduct, <mark>sum</mark> Purchase orice > 1 oroduct	n(quantity	/) <mark>as</mark> Total	Quantities

## 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



database return for Banana?

SELECTproduct, quantityFROMPurchaseGROUP BYproduct



## 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
<b>GROUP BY</b>	product

The DB engine is confused, and rightly so: <u>there is no single quantity</u> 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 through this example carefully!







Our colorful hands represent "team exercises"

SELECT color, avg(numc) anc FROM Shapes GROUP BY color SELECT numc FROM Shapes GROUP BY numc

ightarrow?







SELECT color, avg(numc) anc FROM Shapes GROUP BY color



color	anc
blue	4
orange	5

SELECT numc FROM Shapes GROUP BY numc







SELECT color, avg(numc) anc FROM Shapes GROUP BY color



color	anc
blue	4
orange	5

SELECT numc FROM Shapes GROUP BY numc

numc

3

4

5

6









SELECT color, avg(numc) anc FROM Shapes GROUP BY color



color	anc
blue	4
orange	5



SELECT numc FROM Shapes GROUP BY numc

6

Same as:

SELECT DISTINCT numc FROM Shapes



### Another Example



#### **Purchase**

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

SELECT product, sum(quantity) as SumQ, max(price) as MaxP FROM Purchase GROUP BY product

#### Another Example



#### **Purchase**

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

Product	SumQ	MaxP
Bagel	40	3
Banana	70	4

SELECT product, sum(quantity) as SumQ, max(price) as MaxP FROM Purchase GROUP BY product

Next, focus only on products with at least 50 total quantity

### Group qualification



Q: Similar to before, but only products with at least 50 sales.

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

Product	SumQ	MaxP
Banana	70	4

SELECT product, sum(quantity) as SumQ, max(price) as MaxP FROM Purchase GROUP BY product HAVING sum(quantity) > 50 Quiz



What does this query return over the given database?

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

SELECT product, sum(quantity) as SumQ
FROM Purchase
WHERE quantity > 15
GROUP BY product
HAVING sum(quantity) > 40

#### Quiz



What does this query return over the given database?



## General form of SQL Query

The clauses of a SQL query always appear in this order; you can't reorder them



S: "projection": may contain attributes  $a_1, \ldots, a_k$ and/or any aggregates but no other attributes

C1: is any condition ("comparison predicate" or "join condition") on the attributes in R<sub>1</sub>,...,R<sub>n</sub>

C2: is any condition on aggregates and on attributes  $a_1, \ldots, a_k$  (except if PK)

#### Evaluation ("logical order")

- 1. Evaluate FROM
- 2. WHERE, apply condition C1
- 3. GROUP BY the attributes  $a_1, ..., a_k$
- 4. Apply condition C2 to each group (may have aggregates)
- 5. Compute aggregates in S and return the result
- 6. Sort rows by ORDER BY clause

The logical order is useful for understanding, but not always correct. The ANSI SQL standard does not require a specific processing order and leaves that to the implementation. Recall our intro example with SELECT DISTINCT and order by! Notice that that example can't be explained with the order shown here (you need to assume that the ORDER clause still has access to all attributes)

## Conceptual Evaluation Strategy with GROUP BY

- The cross-product of relation-list is computed (FROM), tuples that fail qualification are discarded (WHERE), then:
- GROUP BY: the remaining tuples are partitioned into groups by the value of attributes in grouping-list.
- HAVING: the group-qualification is applied to eliminate some groups
  - Expressions in group-qualification must have <u>a single value per group</u>!
  - In effect, an attribute in group-qualification that is not an argument of an aggregate op must also appear in grouping-list. (exception: Some DBMSs can exploit primary key semantics!)
- One answer tuple is generated per qualifying group.

#### Don't use new Alias in HAVING clause



What does this query return over the given database?

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

SELECT product, sum(quantity) as SumQ
FROM Purchase
WHERE quantity > 15
GROUP BY product
HAVING SumQ > 35

### Don't use new Alias in HAVING clause



What does this query return over the given database?

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



Error in PostgreSQL and (SQL server)! Reason: HAVING is evaluated before SELECT! (However, SQLite works: different implementation)

SELECTproduct, sum(quantity) as SumQFROMPurchaseWHEREquantity > 15GROUP BY productHAVINGSumQ > 35

### Don't use new Alias in HAVING clause



What does this query return over the given database?

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

ProductSumQBanana50Bagel40

Works! Notice the new sorting

SELECT product, sum(quantity) as SumQ
FROM Purchase
WHERE quantity > 15
GROUP BY product
HAVING sum(quantity) > 35
ORDER BY sumQ desc

Side-topic

# Some more history

Wolfgang Gatterbauer. Database design: https://northeastern-datalab.github.io/cs3200/

214

### SQL: Declarative Programming

SQL

select (e.salary / (e.age - 18)) as comp
from employee as e
where e.name = "Jones"

<u>Declarative Language</u>: you say what you want without having to say how to do it.

<u>Procedural Language</u>: you have to specify exact steps to get the result.

#### SQL: was not the only Attempt

```
SQL select (e.salary / (e.age - 18)) as comp
where e.name = "Jones"
```

```
QUEL retrieve (comp = e.salary / (e.age - 18))
where e.name = "Jones"
```

Commercially not used anymore since ~1980

# Arithmetic & Date functions

### Arithmetic expressions





### Arithmetic expressions







### Arithmetic expressions







## Date functions are database-specific



#### Worker

Name	Birthdate
Max	1980-01-01
Fred	1979-02-01
Susan	1990-01-31
Tilda	1988-01-01

SELECT date('now')-date(birthdate) as age FROM Worker



This is here SQLite semantics. Date functions are different between different databases. You do not need to remember those. In real life, you may need to look up how your DB handles date functions, e.g., <u>http://www.sqlite.org/lang\_datefunc.html</u>, <u>https://www.postgresql.org/docs/current/functions-datetime.html postgres</u>

SELECT extract(year from age(birthdate)) as age FROM Worker SELECT date\_part('year', age(birthdate)) as age FROM Worker

# Practicing more Joins



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

Q: Find all US companies that manufacture at least two different products.





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

Q: Find all US companies that manufacture at least two different products.

SELECTDISTINCT cNameFROMProduct P1, Product P2, CompanyWHEREcountry = 'USA'andP1.manufacturer = cNameandP2.manufacturer = cNameand"the product should be different"





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

Q: Find all US companies that manufacture at least two different products.

SELECT DISTINCT cName
FROM Product P1, Product P2, Company
WHERE country = 'USA'
and P1.manufacturer = cName
and P2.manufacturer = cName
and P1.pName <> P2.pName







PNar	ne	Price	Category	Manufacturer	
Gizm	10	\$19.99	Gadgets	GizmoWorks	
P2	$\diamond$				_
PNa	ne	Price	Category	Manufacturer	
Powe	ergizmo	\$29.99	Gadgets	GizmoWorks	ľ
					-

Company

CName	StockPrice	Country
GizmoWorks	25	USA

**SELECT DISTINCT cName FROM** Product P1, Product P2, Company WHERE country = 'USA' P1.manufacturer = cName and and P2.manufacturer = cName P1.pName <> P2.pName and



## Question: what about 3, 4, 5, etc. different products?



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

Q: Find all US companies that manufacture at least two different products.





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



#### Product

Product				_	Company		
PName	Price	Category	Manufacturer		CName	StockPrice	Country
Gizmo	\$19.99	Gadgets	GizmoWorks		GizmoWorks	25	USA
Powergizmo	\$29.99	Gadgets	GizmoWorks		Canon	65	Japan
SingleTouch	\$149.99	Photography	Canon		Hitachi	15	Japan
MultiTouch	\$203.99	Household	Hitachi				

SELECT	cName
FROM	Product P, Company C
WHERE	manufacturer = cName
	country = 'USA'

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



PName	Price	Category	Manufacturer	CName	StockPrice	Country
Gizmo	\$19.99	Gadgets	GizmoWorks	GizmoWorks	25	USA
Powergizmo	\$29.99	Gadgets	GizmoWorks	GizmoWorks	25	USA
SingleTouch	\$149.99	Photography	Canon	Canon	65	Japan
MultiTouch	\$203.99	Household	Hitachi	Hitachi	15	Japan

SELECT cName FROM Product P, Company C WHERE manufacturer = cName and country = 'USA' GROUP by cName HAVING count(\*) >= 2

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



PName	Price	Category	Manufacturer	CName	StockPrice	Country
Gizmo	\$19.99	Gadgets	GizmoWorks	GizmoWorks	25	USA
Powergizmo	\$29.99	Gadgets	GizmoWorks	GizmoWorks	25	USA
Oliver Le Terrele	¢440.00		0	0	<u>ог</u>	La va ava
	ψ1 <del>4</del> 3.33	၊ ဂ၊၀း၀ၝ၊apriy	Canon		00	Japan
MA. HALT	<b>©</b>		L l'Ac chi	L lite a la l	45	La man
	ψ200.99	riousenoiu	ппасти	ппаст		Japan

SELECT cName FROM Product P, Company C WHERE manufacturer = cName and country = 'USA' GROUP by cName HAVING count(\*) >= 2

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



	M1W			$\hat{\nabla}$	SUM	
PName	Price	Category	Manufacturer	CName	StockPrice	Country
Gizmo	\$19.99	Gadgets	GizmoWorks	GizmoWorks	25	USA
Powergizmo	\$29.99	Gadgets	GizmoWorks	GizmoWorks	25	USA
Single Touch	\$149.99	Photography	Canon	Canon	65	Japan
MultiTouch	\$203.99	Household	Hitachi	Hitachi	15	Japan

SELECT cName
FROM Product P, Company C
WHERE manufacturer = cName
and country = 'USA'
GROUP by cName
HAVING count(\*) >= 2

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



#### 仑

				<b>•</b>		
PName	Price	Category	Manufacturer	CName	StockPrice	Country
Gizmo	\$19.99	Gadgets	GizmoWorks	GizmoWorks	25	USA
Powergizmo	\$29.99	Gadgets	GizmoWorks	GizmoWorks	25	USA
	¢440.00		0	0		
	ψ149.99		Carlon	Carlon	05	Japan
MALLET TALLA	<b>#000 00</b>		Lite elst			
Multi Touch	ψ200.99	riousenoiu	Паст	Тпаст	15	Japan

count(\*) >= 2 ✓

SELECT cName
FROM Product P, Company C
WHERE manufacturer = cName
and country = 'USA'
GROUP by cName
HAVING count(\*) >= 2

