

SQL AND SPREADSHEETS

Instructors: Crista Lopes
Copyright © Instructors.

Data-Centric Programming

- Focus on data
- Interactive data: SQL
- Dataflow: Spreadsheets
- Dataflow: Iterators, generators, coroutines



SQL

Standard Query Language

History

- Data banks since 1950s
- Disks (direct access storage) in 1960s
- How to store and retrieve data from disk
 - ▣ Efficiently, cleanly
- Before 1970:
 - ▣ Hierarchical models (trees)
 - ▣ Network models (graphs)
- E. Codd, 1970:
 - ▣ Relational model

Relational model

- Logic deductive system minus deductions ☺
 - ▣ Data independence – isolate applications from data representations
 - ▣ Data inconsistency
- “Relation” as in Mathematics:
 - ▣ Given sets S_1, S_2, \dots, S_n : R is a relation on these sets iff $R = \{\{e_1, e_2, \dots, e_n\}, \dots\}$ where $e_i \in S_i$
 - ▣ $R \subseteq S_1 \times S_2 \times \dots \times S_n$
(R is a subset of the Cartesian product)

Math: Relations vs. Functions

- Relation >> Function
- Function = relation where each element of the domain corresponds to one element of the range

Int	Int
-2	3
5	3
7	3
12	3

Relation?
Function?

Int	Int
-2	3
-2	0
7	3
12	1

Relation?
Function?

Relations in data bases

- Relations define subsets of the domain:

$$R \subseteq S_1 \times S_2 \times \dots \times S_n$$

supply (supplier part project quantity)

1	2	5	17
1	3	5	23
2	3	7	9
2	7	5	4
4	1	1	12

Supply is a relation (subset) from
supplier x part x project x quantity → supplier x part x project x quantity
And
supplier, part, project, quantity all subsets of Int

Relations in data bases

- Relations may include repeated domains

component (part part quantity)

1	5	9
2	5	7
3	5	2
2	6	12
3	6	3

“Part 1 is a subpart of part 5, and there needs to be 9 part 1s to make a part 5”

Component is a relation (subset) from
part x part x quantity → part x part x quantity

And

part, part, project, quantity all subsets of Int

Relationships

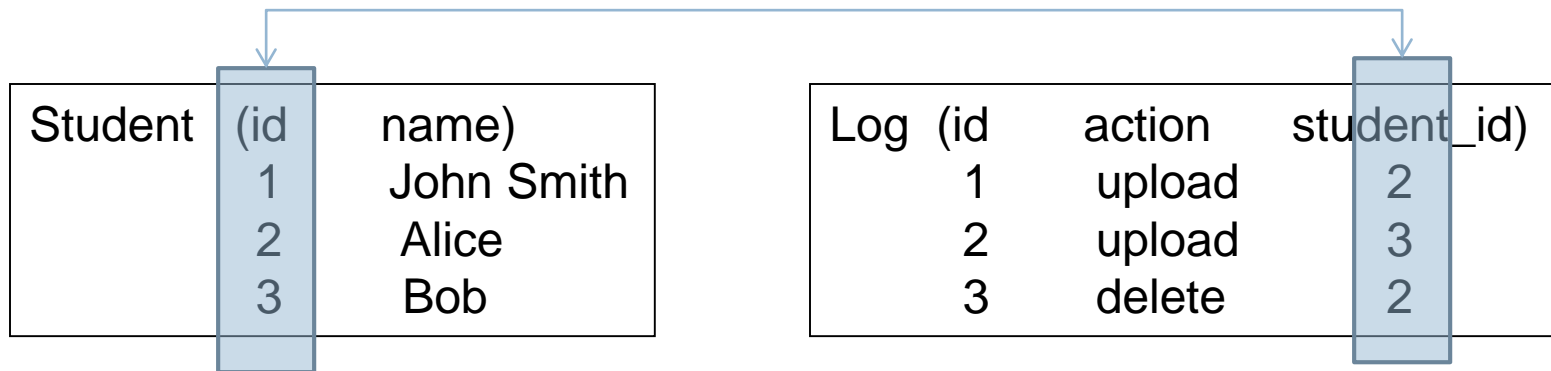
- A relationship is an equivalence class of those relations that are equivalent under permutation of domains
- Order not important
- Same domains are distinguished by role names
- User-facing model

component (**subpart** **part** quantity)

1	5	9
2	5	7
3	5	2
2	6	12
3	6	3

Cross-references

- Elements of a relation can cross-reference elements of the same or another relation
- Done via Keys



Operations on relations

- Permutation
 - ▣ Interchanging columns yields converse relations
- Subsetting
 - ▣ Selecting only a subset of tuples
- Projection
 - ▣ Selection of only a subset of columns
- Join
 - ▣ Merging two or more relations without loss of information

Relational Model → SQL

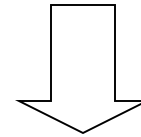
- Data Definition Language (DDL)
 - ▣ Create/alter/delete tables and their attributes
- Data Manipulation Language (DML)
 - ▣ Query one or more tables
 - ▣ Insert/delete/modify tuples in tables

Subsetting

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

```
SELECT *  
FROM Product  
WHERE category='Gadgets'
```



“selection”

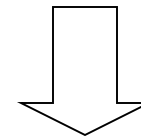
PName	Price	Category	Manufacturer
Gizmo	\$19.99	Gadgets	GizmoWorks
Powergizmo	\$29.99	Gadgets	GizmoWorks

Projection+Subsetting

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

```
SELECT PName, Price, Manufacturer
FROM Product
WHERE Price > 100
```



“selection” and
“projection”

PName	Price	Manufacturer
SingleTouch	\$149.99	Canon
MultiTouch	\$203.99	Hitachi

Joins

Product (pname, price, category, manufacturer)

Company (cname, stockPrice, country)

Find all products under \$200 manufactured in Japan;
return their names and prices.

```
SELECT PName, Price  
FROM Product, Company  
WHERE Manufacturer=CName AND Country='Japan'  
AND Price <= 200
```

Join
between Product
and Company

Joins

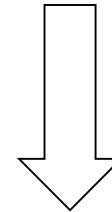
Product

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

```
SELECT PName, Price
FROM Product, Company
WHERE Manufacturer=CName AND Country='Japan'
AND Price <= 200
```



PName	Price
SingleTouch	\$149.99

Full SQL

- Very powerful query language
 - ▣ Ordering, Grouping, aggregation, rich type system, ...
- Declarative
 - ▣ Say *what* you want, not *how* you want it to happen
 - ▣ Nothing related to query processing or internal data representations



Spreadsheets

Spreadsheets

ProductSalesExport-1997 (2).xls [Compatibility Mode] - Microsoft Excel

	A	B	C	D	E	F
	Product	Order	Date	Unit Price	Quantity	Total
2	Alice Mutton	10415	1/15/1997	\$31.20	2	\$62.40
3		10430	1/30/1997	\$31.20	45	\$1,404.00
4		10431	1/30/1997	\$31.20	50	\$1,560.00
5		10444	2/12/1997	\$31.20	10	\$312.00
6		10523	5/1/1997	\$39.00	25	\$975.00
7		10530	5/8/1997	\$39.00	40	\$1,560.00
8		10550	5/28/1997	\$39.00	8	\$312.00
9		10564	6/10/1997	\$39.00	16	\$624.00
10		10573	6/19/1997	\$39.00	18	\$702.00
11		10607	7/22/1997	\$39.00	100	\$3,900.00
12		10686	9/30/1997	\$39.00	30	\$1,170.00
13		10696	10/8/1997	\$39.00	20	\$780.00
14		10698	10/9/1997	\$39.00	8	\$312.00
15		10714	10/22/1997	\$39.00	27	\$1,053.00
16		10727	11/3/1997	\$39.00	20	\$780.00
17		10773	12/11/1997	\$39.00	33	\$1,287.00
18		10795	12/24/1997	\$39.00	35	\$1,365.00
19		10801	12/29/1997	\$39.00	40	\$1,560.00
20	Total:			\$670.80	527	\$19,718.40
21						
22						
23	Aniseed Syrup	10405	1/6/1997	\$8.00	50	\$400.00
24		10485	3/25/1997	\$8.00	20	\$160.00
25		10540	5/19/1997	\$10.00	60	\$600.00
26		10591	7/7/1997	\$10.00	14	\$140.00
27		10702	10/13/1997	\$10.00	6	\$60.00
28		10742	11/14/1997	\$10.00	20	\$200.00
29		10764	12/3/1997	\$10.00	20	\$200.00
30	Total:			\$66.00	190	\$1,760.00
31						
32						

Ready

Spreadsheets

- One of the most successful software genres
- Centuries-old accounting practices...
 - ▣ Some cells contain primitive values
 - ▣ Some cells contain values derived from formulas
- ...with computers
 - ▣ Automatic update of derived values when primitive values change

→ Dataflow programming

```

3
4 #
5 # The columns. Each column is a data element and a formula.
6 # The first 2 columns are the input data, so no formulas.
7 #
8 all_words = [(), None]
9 stop_words = [(), None]
10 non_stop_words = [(), lambda : \
11                     map(lambda w : \
12                         w if w not in stop_words[0] else '', \
13                           all_words[0])]
14 unique_words = [(), lambda :
15                  set([w for w in non_stop_words[0] if w!=''])]
16 counts = [(), lambda :
17            map(lambda w, word_list : word_list.count(w), \
18                unique_words[0], \
19                  itertools.repeat(non_stop_words[0], \
20                                    len(unique_words[0])))]
21 sorted_data = [(), lambda : sorted(zip(list(unique_words[0]), \
22                                         counts[0]), \
23                                     key=operator.itemgetter(1),
24                                     reverse=True))]
25
26 # The entire spreadsheet
27 all_columns = [all_words, stop_words, non_stop_words, \
28                unique_words, counts, sorted_data]
29
30 #
31 # The active procedure over the columns of data.
32 # Call this everytime the input data changes, or periodically.
33 #
34 def update():
35     global all_columns
36     # Apply the formula in each column
37     for c in all_columns:
38         if c[1] != None:
39             c[0] = c[1]()
40
41
42 # Load the fixed data into the first 2 columns
43 all_words[0] = re.findall('[a-z]{2,}', open(sys.argv[1]).read()).
44     lower()
45 stop_words[0] = set(open('../stop_words.txt').read().split(','))
46 # Update the columns with formulas
47 update()
48
49 for (w, c) in sorted_data[0][:25]:
50     print w, '- ', c

```

In Python:

Columns = 2-part lists:
data
formula

All formulas run on updates

In OOP

- Columns = Objects with 2 parts, data and formula
- Formulas = Objects with method “execute”
- “Map” function: applies a given function to one or more list of values
 - ▣ Check for equivalents in C++ (Boost maybe?), C# (Select)
 - ▣ Not hard to do by hand: iterate
- Homework: no ugly code!
 - ▣ Think carefully. Model it nicely. Use the right words. Use ADTs.