

Informatics 1: Data & Analysis

Lecture 7: SQL

Ian Stark

School of Informatics
The University of Edinburgh

Tuesday 2 February 2016
Semester 2 Week 4



THE UNIVERSITY
of EDINBURGH



Careers in IT

Job Fair

Wednesday 3 February 2016

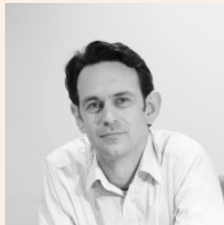
Informatics Forum
1300–1600

<http://www.ed.ac.uk/schools-departments/careers>

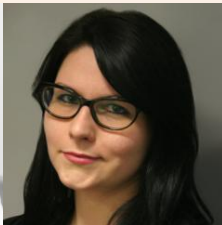
Careers advice and stalls from 40 local, national
and international employers

Careers in Computer Games

Panel and Networking



Michael Boniface
Reloaded



Alice Rendell
Kobojo



Luke Brown
Outplay



Timea Tabori
Rockstar

Informatics Forum — Wednesday 1745–2000 — mycareerhub.ed.ac.uk

BCSWomen Lovelace Colloquium

Sheffield Hallam / University of Sheffield
Thursday 31 March 2016

<https://tinyurl.com/bcs-lovelace>



Sydney Padua: 2D Goggles

One-day conference for women students of computing

Enter the first-year student poster contest!

(Closing date for 250-word abstract Saturday 6 February 2016)

Read This



Inside Google Spanner, the Largest Single Database on Earth

Wired 2012-11-26

<http://www.wired.com/2012/11/google-spanner-time/all/>

Do This

Tuple-relational calculus can be quite tricky to understand, and it's not always obvious to follow what a query means. So, homework this time is to read the lectures slides again, going through the examples to see how each query works.

Top Hat TRC Question

Match the following statements in the tuple-relational calculus to their descriptions.

- X. $\{ R \mid \exists S . S \in \text{Student} \wedge S.\text{name} = R.\text{name} \wedge S.\text{matric} = \text{"s0378435"} \}$
- Y. $\{ R \mid \exists T \in \text{Takes} . R.\text{matric} = T.\text{matric} \wedge T.\text{code} = \text{"dbs"} \}$
- Z. $\{ R \mid \exists T1, T2 \in \text{Takes} . T1.\text{matric} = R.\text{matric} \wedge T1.\text{code} = T2.\text{code} \wedge T2.\text{matric} = \text{"s0378435"} \}$

Data Representation

This first course section starts by presenting two common **data representation models**.

- The *entity-relationship (ER)* model
- The *relational* model

Data Manipulation

This is followed by some methods for manipulating data in the relational model and using it to extract information.

- *Relational algebra*
- The *tuple-relational calculus*
- The query language **SQL**

SQL: Structured Query Language

- SQL is the standard language for interacting with relational database management systems
- Substantial parts of SQL are **declarative**: code states what should be done, not necessarily how to do it.
- When actually querying a large database, database systems take advantage of this to plan, rearrange, and optimize the execution of queries.
- Procedural parts of SQL do contain **imperative** code to make changes to the database.
- While SQL is an international standard (ISO 9075), individual implementations have notable idiosyncrasies, and code may not be entirely portable.

MySQL : PostgreSQL : Oracle : SQL Server : DB2 : SQLite : Sybase

SQL Data Manipulation Language

In an earlier lecture we saw the SQL **Data Definition Language (DDL)**, used to declare the schema of relations and create new tables.

This lecture introduces the **Data Manipulation Language (DML)** which allows us to:

- Insert, delete and update rows in existing tables;
- Query the database.

Note that “query” here covers many different scales: from extracting a single statistic or a simple list, to building large tables that combine several others, or creating *views* on existing data.

SQL is a large and complex language. Here we shall only see some of the basic and most important parts. For a much more extensive coverage of the topic, sign up for the *Database Systems* course in Year 3.

Inserting Data into a Table

```
CREATE TABLE Student (  
    matric VARCHAR(8),  
    name   VARCHAR(20),  
    age    INTEGER,  
    email  VARCHAR(25),  
    PRIMARY KEY (matric) )
```

The following adds a single record to this table:

```
INSERT  
    INTO Student (matric, name, age, email)  
    VALUES ('s1428751', 'Bob', 19, 'bob@sms.ed.ac.uk')
```

For multiple records, repeat; or consult your RDBMS manual.

Strictly, SQL allows omission of the field names; but if we include them, then the compiler will check them against the schema for us.

Update and Delete Rows in a Table

Update

This command changes the **name** recorded for one student:

```
UPDATE Student  
  SET name = 'Bobby'  
  WHERE matric = 's1428571'
```

Delete

This deletes from the table all records for students named “Bobby”:

```
DELETE  
  FROM Students  
  WHERE name = 'Bobby'
```

Simple Query

Extract all records for students older than 19.

```
SELECT *  
FROM Student  
WHERE age > 19
```

Returns a new table, with the same schema as **Student**, but containing only some of its rows.

Student

matric	name	age	email
s0456782	John	18	john@inf
s0378435	Helen	20	helen@phys
s0412375	Mary	18	mary@inf
s0189034	Peter	22	peter@math

Simple Query

Extract all records for students older than 19.

```
SELECT *  
FROM Student  
WHERE age > 19
```

Returns a new table, with the same schema as **Student**, but containing only some of its rows.

Student

matric	name	age	email
s0456782	John	18	john@inf
s0378435	Helen	20	helen@phys
s0412375	Mary	18	mary@inf
s0189034	Peter	22	peter@math

Simple Query

Extract all records for students older than 19.

```
SELECT *  
FROM Student  
WHERE age > 19
```

Returns a new table, with the same schema as **Student**, but containing only some of its rows.

matric	name	age	email
s0378435	Helen	20	helen@phys
s0189034	Peter	22	peter@math

Simple Query

Extract all records for students older than 19.

```
SELECT *  
FROM Student  
WHERE age > 19
```

Returns a new table, with the same schema as **Student**, but containing only some of its rows.

Tuple-Relational Calculus

SQL is similar in form to the **comprehensions** of tuple-relational calculus:

$$\{ S \mid S \in \text{Student} \wedge S.\text{age} > 19 \}$$

Working out how to implement this efficiently through relational algebra operations is the job of an SQL compiler and database query engine.

Simple Query

Extract all records for students older than 19.

```
SELECT *  
FROM Student  
WHERE age > 19
```

Returns a new table, with the same schema as **Student**, but containing only some of its rows.

Variations

We can explicitly name the selected fields.

```
SELECT matric, name, age, email  
FROM Student  
WHERE age > 19
```


Simple Query

Extract all records for students older than 19.

```
SELECT *  
FROM Student  
WHERE age > 19
```

Returns a new table, with the same schema as `Student`, but containing only some of its rows.

Variations

We can identify which table the fields are from.

```
SELECT Student.matric, Student.name, Student.age, Student.email  
FROM Student  
WHERE Student.age > 19
```

Simple Query

Extract all records for students older than 19.

```
SELECT *  
FROM Student  
WHERE age > 19
```

Returns a new table, with the same schema as *Student*, but containing only some of its rows.

Variations

We can locally abbreviate the table name with an *alias*.

```
SELECT S.matric, S.name, S.age, S.email  
FROM Student AS S  
WHERE S.age > 19
```

Simple Query

Extract all records for students older than 19.

```
SELECT *  
FROM Student  
WHERE age > 19
```

Returns a new table, with the same schema as **Student**, but containing only some of its rows.

Variations

We can save ourselves a very small amount of typing.

```
SELECT S.matric, S.name, S.age, S.email  
FROM Student S  
WHERE S.age > 19
```

Anatomy of an SQL Query

```
SELECT field-list  
FROM table-list  
[ WHERE qualification ]
```

- The **SELECT** keyword starts the query.
- The list of fields specifies *projection*: what columns should be retained in the result. Using *** means all fields.
- The **FROM** clause lists one or more tables from which to take data.
- An optional **WHERE** clause specifies *selection*: which records to pick out and return from those tables.

Anatomy of an SQL Query

```
SELECT field-list  
FROM table-list  
[ WHERE qualification ]
```

The *table-list* in the **FROM** clause is a comma-separated list of tables to be used in the query:

```
...  
FROM Student, Takes, Course  
...
```

Each table can be followed by an alias **Course AS C**, or even just **Course C**.

Anatomy of an SQL Query

```
SELECT field-list  
FROM table-list  
[ WHERE qualification ]
```

The *field-list* after **SELECT** is a comma-separated list of (expressions involving) names of fields from the tables in **FROM**.

```
SELECT name, age
```

```
...
```

```
...
```

Field names can be referred to explicitly using table names or aliases: such as **Student.name** or **C.title**.

Anatomy of an SQL Query

```
SELECT field-list  
FROM table-list  
[ WHERE qualification ]
```

The *qualification* in the **WHERE** clause is a logical expression built from tests involving field names, constants and arithmetic expressions.

...

...

```
WHERE age > 18 AND age < 65
```

Expressions can involve a range of numeric, string and date operations.

Simple Query with Multiset Result

Extract all student ages.

```
SELECT age  
FROM Student
```

Returns a new table, similar to *Student*, but containing only some of its columns.

Student

matric	name	age	email
s0456782	John	18	john@inf
s0378435	Helen	20	helen@phys
s0412375	Mary	18	mary@inf
s0189034	Peter	22	peter@math

Simple Query with Multiset Result

Extract all student ages.

```
SELECT age  
FROM Student
```

Returns a new table, similar to `Student`, but containing only some of its columns.

Student

matric	name	age	email
s0456782	John	18	john@inf
s0378435	Helen	20	helen@phys
s0412375	Mary	18	mary@inf
s0189034	Peter	22	peter@math

Simple Query with Multiset Result

Extract all student ages.

```
SELECT age  
FROM Student
```

Returns a new table, similar to `Student`, but containing only some of its columns.

age
18
20
18
22

Aside: Multisets

The relational model given in earlier lectures has tables as *sets* of rows: so the ordering doesn't matter, and there are no duplicates.

Actual SQL does allow duplicate rows, with a **SELECT DISTINCT** operation to remove duplicates on request.

Thus SQL relations are not sets but *multisets* of rows. A multiset, or *bag*, is like a set but values can appear several times. The number of repetitions of a value is its *multiplicity* in the bag.

The following are distinct multisets:

$\{2, 3, 5\}$ $\{2, 3, 3, 5\}$ $\{2, 3, 3, 5, 5, 5\}$ $\{2, 2, 2, 3, 5\}$

Ordering still doesn't matter, so these are all the same multiset:

$\{2, 2, 3, 5\}$ $\{2, 3, 2, 5\}$ $\{5, 2, 3, 2\}$ $\{3, 2, 2, 5\}$

Simple Query with Set Result

Extract all student ages.

```
SELECT DISTINCT age  
FROM Student
```

Returns a new table, similar to *Student*, but containing only some elements from some of its columns.

Student

matric	name	age	email
s0456782	John	18	john@inf
s0378435	Helen	20	helen@phys
s0412375	Mary	18	mary@inf
s0189034	Peter	22	peter@math

Simple Query with Set Result

Extract all student ages.

```
SELECT DISTINCT age  
FROM Student
```

Returns a new table, similar to `Student`, but containing only some elements from some of its columns.

Student

matric	name	age	email
s0456782	John	18	john@inf
s0378435	Helen	20	helen@phys
s0412375	Mary	18	mary@inf
s0189034	Peter	22	peter@math

Simple Query with Set Result

Extract all student ages.

```
SELECT DISTINCT age  
FROM Student
```

Returns a new table, similar to *Student*, but containing only some elements from some of its columns.

age
18
20
22

Further Aside: Quotation Marks in SQL Syntax

SQL uses alphanumeric **tokens** of three kinds:

- Keywords: **SELECT**, **FROM**, **UPDATE**, ...
- Identifiers: **Student**, **matric**, **age**, **S**, ...
- Strings: **'Bobby'**, **'Informatics 1'**, ...

Each of these kinds of token has different rules about **case sensitivity**, the use of **quotation marks**, and whether they can contain **spaces**.

While programmers can use a variety of formats, and SQL compilers should accept them, programs that *generate* SQL code may be very cautious in what they emit and use apparently verbose formats.

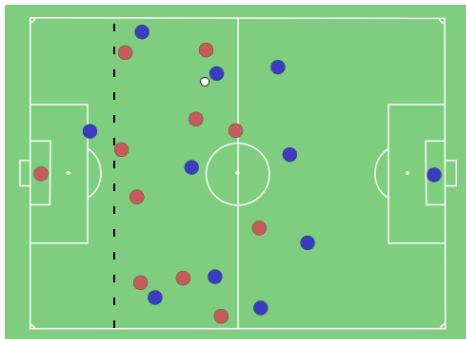
Further Aside: Know Your Syntax

		Case sensitive?	Spaces allowed?	Quotation character?	Quotation Required?
Keywords	FROM	No	Never	None	No
Identifiers	Student	Maybe	If quoted	"Double"	If spaces
Strings	'Bob'	It depends	Yes	'Single'	Always

For example:

```
select matric  
from Student as "Student Table"  
where "Student Table".age > 19 and "name" = 'Bobby Tables'
```

It's always safe to use only uppercase keywords and put quotation marks around all identifiers. Some tools will do this automatically.



NielsF, Wikimedia Commons

The blue forward on the left of the diagram is in an offside position as he is in front of both the second-to-last defender (marked by the dotted line) and the ball. Note that this does not necessarily mean he is committing an *offside offence*; it only becomes an offence if the ball were to be played to him at this moment, whether or not he is in an offside position when he receives the ball, as he could receive the ball in an *onside position* but he'd still have committed an *offside offence*.

(FIFA guidelines 2003; IFAB Law XI 2005; Clarified 2010; Explained by Wikipedia)

Bring your solutions, or your attempts at them. You will need to be able to show these to your tutor and exchange them with other students.

Come to tutorials prepared. Students who have not even attempted the exercises will be sent away to do them elsewhere and return later.

Even so, if you feel you are falling behind and cannot complete the work **do not skip the tutorial**. Instead, do what you can and then explain at the tutorial where you could make no more progress. Tutors are there to help you at whatever level you have reached.

You can also take your exercises to **InfBASE** to ask for help.

UK Labour Market, January 2016



Coverage: **UK**

Date: **20 January 2016**

Geographical Area: **UK**

Theme: **Labour Market**

Main points for September to November 2015

- There were 31.39 million people in work, 267,000 more than for June to August 2015 and 588,000 more than for a year earlier.
- There were 22.96 million people working full-time, 436,000 more than for a year earlier. There were 8.43 million people working part-time, 152,000 more than for a year earlier.
- The employment rate (the proportion of people aged from 16 to 64 who were in work) was 74.0%, the highest since comparable records began in 1971.

- Full-time study is a full-time job. During all weeks of semester, plan to spend at least 40 hours each week across your courses.
- Within each week, balance time between courses according to their relative credit point weights out of 60.
- For a 10-point course like Inf1-DA, that suggests at least 6–7 hours of study per week.
- Inf1-DA has three contact hours per week, which means spending **at least the same again** in independent study.
- This rule of thumb applies to many courses: for every directly taught hour (lecture, tutorial, lab), add at least one hour for self-study.

Some activities are very different, such as field trips or research projects

- Some students also undertake paid or voluntary work outside their studies.
- The University recommends an absolute maximum of 15 hours each week for this.
- This makes a baseline of 55 hours each week on study + outside work. That's possible, but it is a lot of hours.

Simple Query

Extract all records for students older than 19.

```
SELECT *  
FROM Student  
WHERE age > 19
```

Returns a new table, with the same schema as **Student**, but containing only some of its rows.

$$\{ S \mid S \in \text{Student} \wedge S.\text{age} > 19 \}$$

matric	name	age	email
s0378435	Helen	20	helen@phys
s0189034	Peter	22	peter@math

Students and Courses

Student

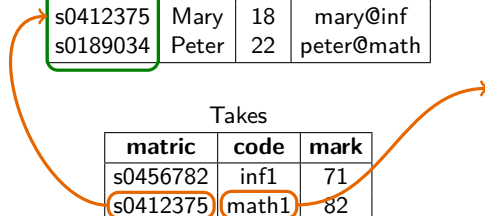
matric	name	age	email
s0456782	John	18	john@inf
s0378435	Helen	20	helen@phys
s0412375	Mary	18	mary@inf
s0189034	Peter	22	peter@math

Takes

matric	code	mark
s0456782	inf1	71
s0412375	math1	82
s0412375	geo1	64
s0189034	math1	56

Course

code	title	year
inf1	Informatics 1	1
math1	Mathematics 1	1
geo1	Geology 1	1
db1	Database Systems	3
adbs	Advanced Databases	4



Example Query

Find the names and email addresses of all students taking Mathematics 1.

```
SELECT Student.name, Student.email
FROM Student, Takes, Course
WHERE Student.matric = Takes.matric
      AND Takes.code = Course.code
      AND Course.title = 'Mathematics 1'
```

Take rows from all three tables at once,

Student				Takes			Course		
matric	name	age	email	matric	code	mark	code	title	year
s0456780	John	18	john@inf	s0456780	inf1	71	inf1	Informatics 1	1
s0378435	Helen	20	helen@phys	s0412375	math1	82	math1	Mathematics 1	1
s0412375	Mary	18	mary@inf	s0412375	geo1	64	geo1	Geology 1	1
s0189034	Peter	22	peter@math	s0189034	math1	56	dbb	Database Systems	3
							adbs	Advanced Databases	4

Example Query

Find the names and email addresses of all students taking Mathematics 1.

```
SELECT Student.name, Student.email
FROM Student, Takes, Course
WHERE Student.matric = Takes.matric
      AND Takes.code = Course.code
      AND Course.title = 'Mathematics 1'
```

Take rows from all three tables at once, pick out only those row combinations which match the test,

Student				Takes			Course		
matric	name	age	email	matric	code	mark	code	title	year
s0456780	John	18	john@inf	s0456780	inf1	71	inf1	Informatics 1	1
s0378435	Helen	20	helen@phys	s0412375	math1	82	math1	Mathematics 1	1
s0412375	Mary	18	mary@inf	s0412375	geo1	64	geo1	Geology 1	1
s0189034	Peter	22	peter@math	s0189034	math1	56	db1	Database Systems	3
							adbs	Advanced Databases	4

Example Query

Find the names and email addresses of all students taking Mathematics 1.

```
SELECT Student.name, Student.email
FROM Student, Takes, Course
WHERE Student.matric = Takes.matric
      AND Takes.code = Course.code
      AND Course.title = 'Mathematics 1'
```

Take rows from all three tables at once, pick out only those row combinations which match the test,

Student				Takes			Course		
matric	name	age	email	matric	code	mark	code	title	year
s0456780	John	18	john@inf	s0456780	inf1	71	inf1	Informatics 1	1
s0378435	Helen	20	helen@phys	s0412375	math1	82	math1	Mathematics 1	1
s0412375	Mary	18	mary@inf	s0412375	geo1	64	geo1	Geology 1	1
s0189034	Peter	22	peter@math	s0189034	math1	56	db1	Database Systems	3
							adbs	Advanced Databases	4

Example Query

Find the names and email addresses of all students taking Mathematics 1.

```
SELECT Student.name, Student.email
FROM Student, Takes, Course
WHERE Student.matric = Takes.matric
      AND Takes.code = Course.code
      AND Course.title = 'Mathematics 1'
```

Take rows from all three tables at once, pick out only those row combinations which match the test, and return the named columns.

Student				Takes			Course		
matric	name	age	email	matric	code	mark	code	title	year
s0456780	John	18	john@inf	s0456780	inf1	71	inf1	Informatics 1	1
s0378435	Helen	20	helen@phys	s0412375	math1	82	math1	Mathematics 1	1
s0412375	Mary	18	mary@inf	s0412375	geo1	64	geo1	Geology 1	1
s0189034	Peter	22	peter@math	s0189034	math1	56	db1	Database Systems	3
							adbs	Advanced Databases	4

Example Query

Find the names and email addresses of all students taking Mathematics 1.

```
SELECT Student.name, Student.email  
FROM Student, Takes, Course  
WHERE Student.matric = Takes.matric  
        AND Takes.code = Course.code  
        AND Course.title = 'Mathematics 1'
```

Take rows from all three tables at once, pick out only those row combinations which match the test, and return the named columns.

name	email
Mary	mary@inf
Peter	peter@math

Example Query

Find the names and email addresses of all students taking Mathematics 1.

```
SELECT Student.name, Student.email  
FROM Student, Takes, Course  
WHERE Student.matric = Takes.matric  
        AND Takes.code = Course.code  
        AND Course.title = 'Mathematics 1'
```

Take rows from all three tables at once, pick out only those row combinations which match the test, and return the named columns.

Expressed in tuple-relational calculus:

$$\{ R \mid \exists S \in \text{Student}, T \in \text{Takes}, C \in \text{Course} . \\ R.\text{name} = S.\text{name} \wedge R.\text{email} = S.\text{email} \wedge S.\text{matric} = T.\text{matric} \\ \wedge T.\text{code} = C.\text{code} \wedge C.\text{title} = \text{"Mathematics 1"} \}$$

Example Query

Find the names and email addresses of all students taking Mathematics 1.

```
SELECT Student.name, Student.email  
FROM Student, Takes, Course  
WHERE Student.matric = Takes.matric  
      AND Takes.code = Course.code  
      AND Course.title = 'Mathematics 1'
```

Take rows from all three tables at once, pick out only those row combinations which match the test, and return the named columns.

Implemented in relational algebra,

$$\pi_{\text{name,email}}(\sigma_{\text{Student.matric} = \text{Takes.matric} \wedge \text{Takes.code} = \text{Course.code} \wedge \text{Course.name} = \text{"Mathematics 1"}}(\text{Student} \times \text{Takes} \times \text{Course}))$$

Example Query

Find the names and email addresses of all students taking Mathematics 1.

```
SELECT Student.name, Student.email  
FROM Student, Takes, Course  
WHERE Student.matric = Takes.matric  
        AND Takes.code = Course.code  
        AND Course.title = 'Mathematics 1'
```

Take rows from all three tables at once, pick out only those row combinations which match the test, and return the named columns.

Implemented in relational algebra, in several possible ways:

$$\pi_{\text{name,email}}(\sigma_{\text{title}=\text{"Mathematics 1"}}(\text{Student} \bowtie \text{Takes} \bowtie \text{Course}))$$

$$\pi_{\text{name,email}}(\text{Student} \bowtie (\text{Takes} \bowtie (\sigma_{\text{title}=\text{"Mathematics 1"}}(\text{Course}))))$$

Query Evaluation

SQL **SELECT** queries are very close to a programming-language form for the expressions of the tuple-relational calculus, describing the information desired but not dictating how it should be computed.

To do that computation, we need something more like relational algebra. A single **SELECT** statement combines the operations of join, selection and projection, which immediately suggests one strategy:

- Compute the complete cross product of all the **FROM** tables;
- Select all the rows which match the **WHERE** condition;
- Project out only the columns named on the **SELECT** line.

Crucially, real database engines don't do that. Instead, they use relational algebra to rewrite that procedure into a range of different possible *query plans*, estimate the cost of each — looking at indexes, table sizes, selectivity, potential parallelism — and then execute one of them.

Find the names and email addresses of all students taking Mathematics 1.

```
SELECT Student.name, Student.email  
FROM Student JOIN Takes ON Student.matric=Takes.matric  
           JOIN Course ON Takes.code = Course.code  
WHERE Course.title = 'Mathematics 1'
```

This is **explicit JOIN** syntax.

It has exactly the same effect as **implicit JOIN** syntax:

```
SELECT Student.name, Student.email  
FROM Student, Takes, Course  
WHERE Student.matric = Takes.matric  
       AND Takes.code = Course.code  
       AND Course.title = 'Mathematics 1'
```



A *transaction* is a single coherent operation on a database. This might involve substantial amounts of data, or take considerable computation; but is intended to be an all-or-nothing action.

The features that characterise a reliable implementation of transactions are standardly initialized as the *ACID* properties.

Task

Find out what each letter **A C I D** stands for here, and what those four terms mean.

Homework (2/2): Try This

!

Try writing some SQL by hand using one of these web demonstrators.

w3schools.com <https://tinyurl.com/try-sql>
SQL Bolt <http://sqlbolt.com>
SQL Fiddle <http://sqlfiddle.com>

SQL Statement: Edit the SQL Statement, and click "Run SQL" to see the result.

```
SELECT * FROM Products WHERE Price>75;
```

Run SQL »

Result:

Number of Records: 4

ProductID	ProductName	SupplierID	CategoryID	Unit	Price
9	Mishi Kobe Niku	4	6	18 - 500 g pkgs.	97
20	Sir Rodney's Marmalade	8	3	30 gift boxes	81
29	Thüringer Rostbratwurst	12	6	50 bags x 30 saugs.	123.79
38	Côte de Blaye	18	1	12 - 75 cl bottles	263.5

Your Database: ?

Tablename	Records
Customers	91
Categories	8
Employees	10
OrderDetails	518
Orders	196
Products	77
Shippers	3
Suppliers	29

Restore Database