

Week 4

SQL – Ubiquitous Database Format/Language

Applied Data Science

Columbia University - Columbia Engineering

- ❖ Week 1: Python Basics: How to Translate Procedures into Codes
- ❖ Week 2: Intermediate Python — Data Structures for Your Analysis
- ❖ Week 3: Relational Databases — Where Big Data is Typically Stored
- ❖ **Week 4: SQL — Ubiquitous Database Format/Language**
- ❖ Week 5: Statistical Distributions — The Shape of Data
- ❖ Week 6: Sampling — When You Can't or Won't Have ALL the Data
- ❖ Week 7: Hypothesis Testing — Answering Questions about Your Data
- ❖ Week 8: Data Analysis and Visualization — Using Python's NumPy for Analysis
- ❖ Week 9: Data Analysis and Visualization — Using Python's Pandas for Data Wrangling
- ❖ Week 10: Text Mining — Automatic Understanding of Text
- ❖ Week 11: Machine Learning — Basic Regression and Classification
- ❖ Week 12: Machine Learning — Decision Trees and Clustering

➡ The language for relational databases

➡ SQL is a *declarative language*

➡ SQL queries act on tables

➡ SQL queries return tables

- A **table** is a logically connected set of data organized as a set of columns and rows

A **record** is a row of information about a single item in a table

- A **field** is an individual data item in a record
- Each field has a precise **format** (e.g., number, string, date, etc.)
- A **key** field is a field whose value uniquely identifies a record in a table

SQL Servers

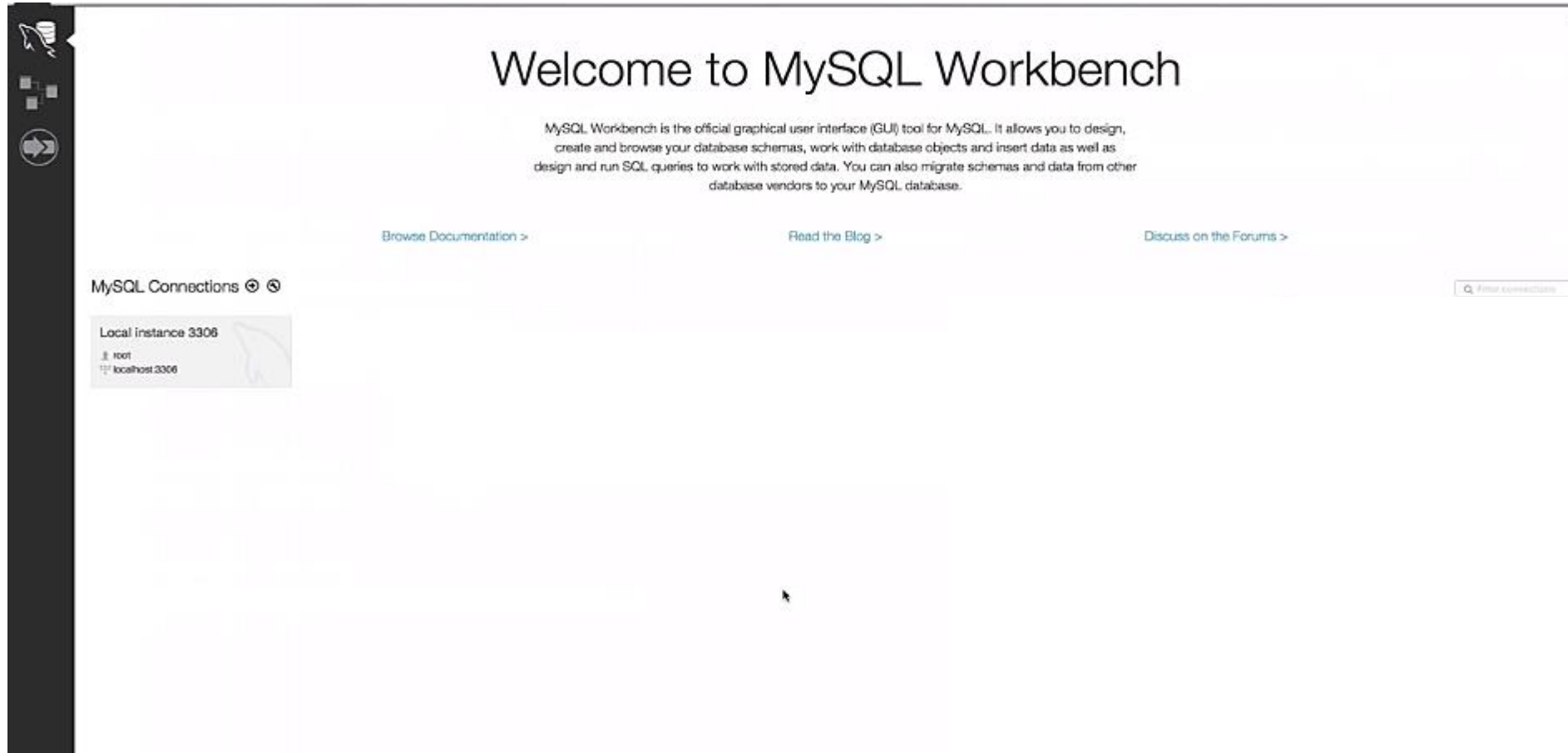
- ➡ The program that manages the database
- ➡ Provides database services to client programs
- ➡ Usually modeled as a **client server** system
- ➡ Numerous choices:
 - ➡ Microsoft SQL Server
 - ➡ Oracle
 - ➡ DB2
 - ➡ PostgreSQL
 - ➡ MySQL

MySQL

- An open source relational DBMS
- MySQL server is where the database resides
- MySQL clients are the programs that talk to the server
- One server can have many clients

- Data definitional
create, alter, drop
- Data manipulation
insert, delete, update
- Data querying
select
select where
select join
- List existing databases:
 - `SHOW databases;`
- Create new database:
 - `CREATE database IF NOT EXISTS SchoolDB;`
- Choose a DB to work on
 - `USE SchoolDB;`
- Delete a database:
 - `DROP DATABASE SchoolDB;`

MySQL Workbench is a GUI client for MySQL



Student

ssn	f name	l name	phone	city	zip
111-22-3333	John	Childs	646.123.1212	New York	10025
123-12-1234	Mary	Arias		New York	10011
555-11-7777	Roberto	Perez	917.333.5479	San Francisco	94110
222-33-4455	Lila	Pennington	425.123.1212	Seattle	98105



Course

number	name	room
c1	Data analytics	1127
c2	Python	303
c3	corp fin	331
c4	prod. mgmt	1127
c5	Ethics	303
c6	leadership	303
c7	bus analytics	1127

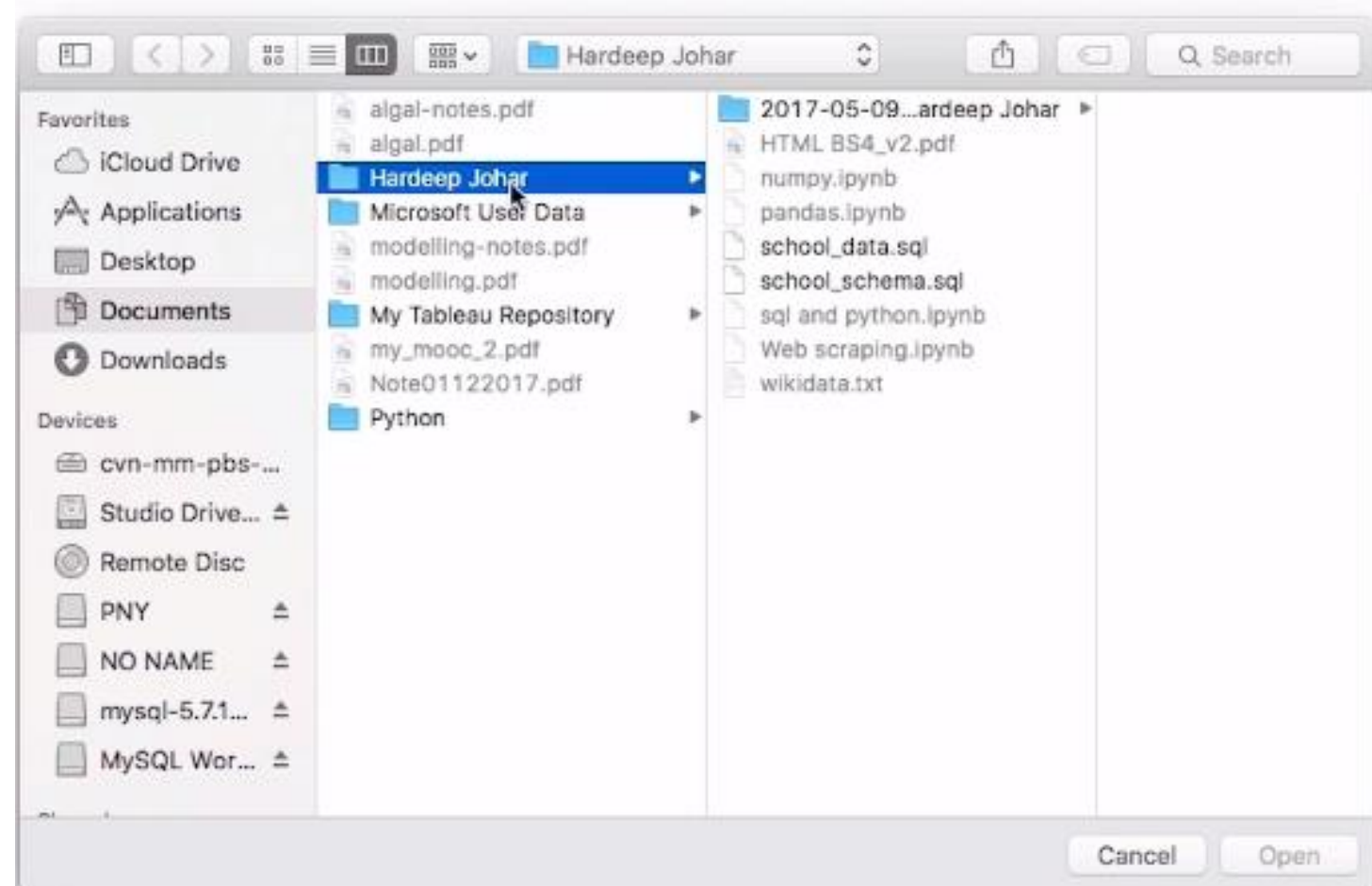
Enrolls-in

ssn	class	score
111-22-3333	c1	93
123-12-1234	c1	87
111-22-3333	c2	95
222-33-4455	c3	44
555-11-7777	c1	36

Setup the Database

- Download the files
 - school_schema.sql
 - school_data.sql
- Open mysql workbench (see installation instructions)
- Open school_schema.sql (click the  icon)
- Execute the script (click the  icon)
- Repeat with the file school_data.sql

```
1 DROP DATABASE IF EXISTS schooldb;
2 CREATE DATABASE IF NOT EXISTS schooldb;
3 USE schooldb;
4 CREATE TABLE IF NOT EXISTS Student (
5     ssn VARCHAR(11) NOT NULL PRIMARY KEY,
6     f_name VARCHAR(20) NOT NULL,
7     l_name VARCHAR(20) NOT NULL,
8     phone VARCHAR(10),
9     city VARCHAR(20) NOT NULL,
10    zip VARCHAR(5) NOT NULL
11 );
12
13 CREATE TABLE IF NOT EXISTS Course (
14     number VARCHAR(3) NOT NULL PRIMARY KEY,
15     name VARCHAR(30) NOT NULL,
16     room INT NOT NULL
17 );
18
19 CREATE TABLE IF NOT EXISTS Enrolls_in (
20     ssn VARCHAR(11) NOT NULL,
21     class VARCHAR(3) NOT NULL,
22     score FLOAT,
23     CONSTRAINT pk_enroll PRIMARY KEY (ssn,class)
24 );
```



Database Example



COLUMBIA | ENGINEERING
EXECUTIVE EDUCATION

1 • DESCRIBE enroll_in;

Limit to 1000 rows

Result Grid

Field	Type	Null	Key	Default	Extra
ssn	varchar(11)	NO	PRI		
class	varchar(3)	NO	PRI		
score	float	YES			

100% 9:1

Export: [Icon]

Field Type Null Key Default Extra

ssn varchar(11) NO PRI

class varchar(3) NO PRI

score float YES

2017-05-09...ardeep Johar

HTML BS4_v2.pdf

numpy.ipynb

pandas.ipynb

school_data.sql

school_schema.sql

sql and python.ipynb

Web scraping.ipynb

wikidata.txt

school_data.sql

SQL File - 1 KB

Created Today, 10:50 AM

Modified Today, 10:50 AM

Last opened Today, 10:50 AM

Add Tags...

Cancel Open

school_data

```
1 • INSERT IGNORE INTO Student VALUES ("111-22-3333", "JOHN", "CHILDS", "6461231212", "New York", "10025");
2 • INSERT IGNORE INTO Student (SSN, f_name, l_name, city, zip) VALUES (
3   "123-12-1234", "Mary", "Arias", "New York", "10011");
4 • INSERT IGNORE INTO Student VALUES ("555-11-7777", "Roberto", "Perez", "9173335479", "San Francisco", "94110");
5 • INSERT IGNORE INTO Student VALUES ("222-33-4455", "Lila", "Pennington", "4251231212", "Seattle", "98105");
6
7 • INSERT IGNORE INTO Course VALUES ("c1", "Data analytics", 1127);
8 • INSERT IGNORE INTO Course VALUES ("c2", "Python", 303);
9 • INSERT IGNORE INTO Course VALUES ("c3", "corp fin", 331);
10 • INSERT IGNORE INTO Course VALUES ("c4", "prod. mgmt", 1127);
11 • INSERT IGNORE INTO Course VALUES ("c5", "Ethics", 303);
12 • INSERT IGNORE INTO Course VALUES ("c6", "leadership", 303);
13 • INSERT IGNORE INTO Course VALUES ("c7", "bus analytics", 1127);
14
15 • INSERT IGNORE INTO Enrolls_in VALUES ("111-22-3333", "c1", 93.00);
16 • INSERT IGNORE INTO Enrolls_in VALUES ("123-12-1234", "c1", 87.00);
17 • INSERT IGNORE INTO Enrolls_in VALUES ("111-22-3333", "c2", 95.00);
18 • INSERT IGNORE INTO Enrolls_in VALUES ("222-33-4455", "c3", 44.00);
19 • INSERT IGNORE INTO Enrolls_in VALUES ("555-11-7777", "c1", 36.00);
```

100% 5:7

Export: [Icon]

Field Type Null Key Default Extra

ssn varchar(11) NO PRI

class varchar(3) NO PRI

score float YES

Limit to 1000 rows

```
1 • INSERT IGNORE INTO Student VALUES ("111-22-3333", "JOHN", "CHILDS", "6461231212", "New York", "10025");
2 • INSERT IGNORE INTO Student (SSN, f_name, l_name, city, zip) VALUES (
3   "123-12-1234", "Mary", "Arias", "New York", "10011");
4 • INSERT IGNORE INTO Student VALUES ("555-11-7777", "Roberto", "Perez", "9173335479", "San Francisco", "94110");
5 • INSERT IGNORE INTO Student VALUES ("222-33-4455", "Lila", "Pennington", "4251231212", "Seattle", "98105");
6
7 • INSERT IGNORE INTO Course VALUES ("c1", "Data analytics", 1127);
8 • INSERT IGNORE INTO Course VALUES ("c2", "Python", 303);
9 • INSERT IGNORE INTO Course VALUES ("c3", "corp fin", 331);
10 • INSERT IGNORE INTO Course VALUES ("c4", "prod. mgmt", 1127);
11 • INSERT IGNORE INTO Course VALUES ("c5", "Ethics", 303);
12 • INSERT IGNORE INTO Course VALUES ("c6", "leadership", 303);
13 • INSERT IGNORE INTO Course VALUES ("c7", "bus analytics", 1127);
14
15 • INSERT IGNORE INTO Enrolls_in VALUES ("111-22-3333", "c1", 93.00);
16 • INSERT IGNORE INTO Enrolls_in VALUES ("123-12-1234", "c1", 87.00);
17 • INSERT IGNORE INTO Enrolls_in VALUES ("111-22-3333", "c2", 95.00);
18 • INSERT IGNORE INTO Enrolls_in VALUES ("222-33-4455", "c3", 44.00);
19 • INSERT IGNORE INTO Enrolls_in VALUES ("555-11-7777", "c1", 36.00);
```

school_data

Limit to 1000 rows

1 • SELECT * FROM Enrolls_in;

100% 26:1

Result Grid

ssn	class	score
111-22-3333	c1	93
111-22-3333	c2	95
123-12-1234	c1	87
222-33-4455	c3	44
555-11-7777	c1	36

Filter Rows: [Search]

Export: [Icon]

- SELECT returns a set of records from one or more tables in a database
- SELECT describes the result set - processing is left to the server

Select

1. Get everything in a table

```
SELECT * FROM Student;
```

2. Get all students from New York

```
SELECT * FROM Student WHERE City = "New York";
```

3. Get only the names of students from New York

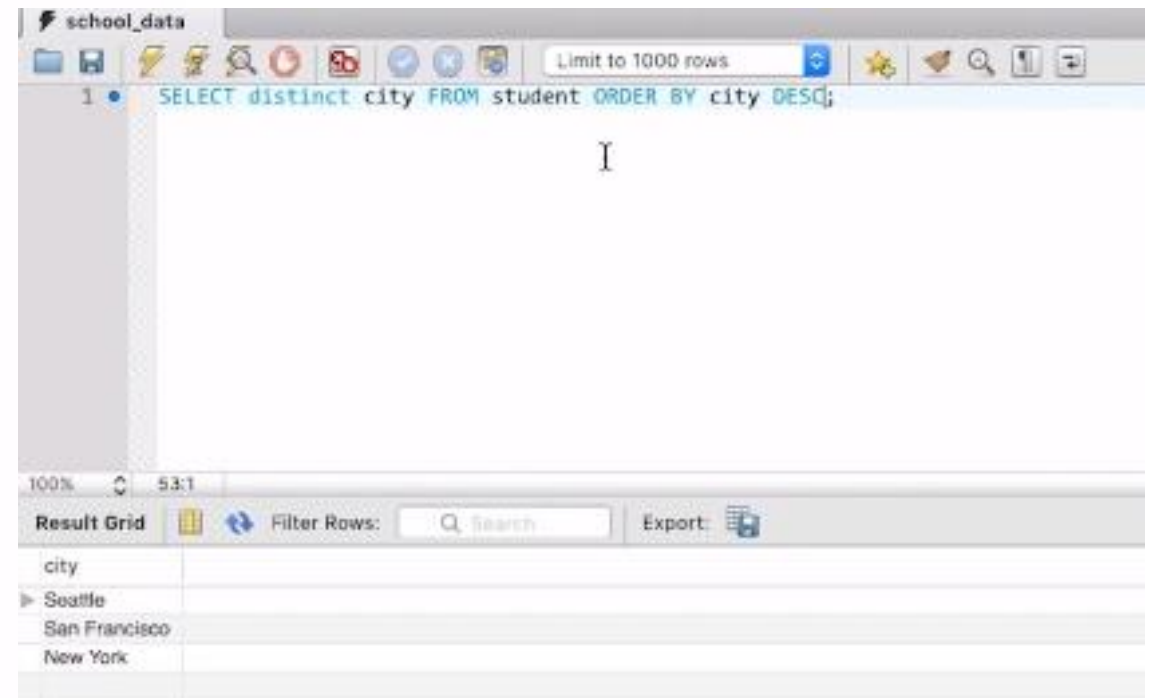
```
SELECT f_name, l_name FROM Student WHERE city = "New York";
```

4. Get the ssn of students ordered by f_name

```
SELECT ssn FROM Student ORDER BY f_name ASC;  
SELECT ssn FROM Student ORDER BY f_name DESC;
```

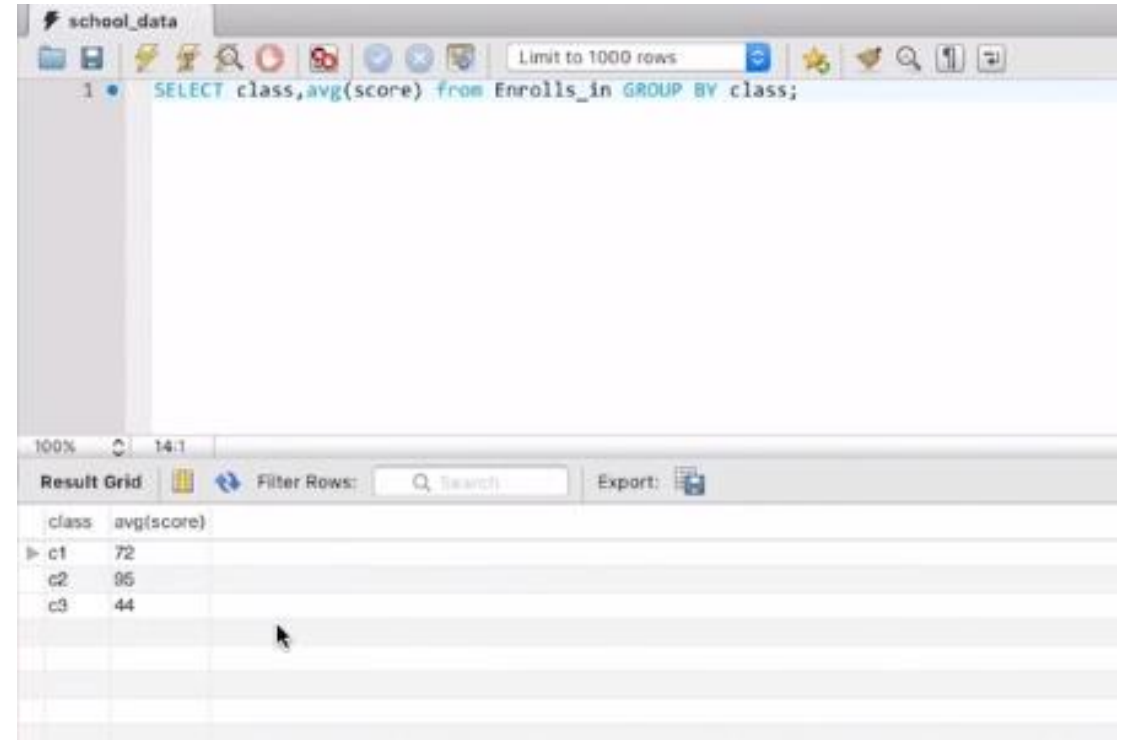
5. List the set of unique cities in the Student table

```
SELECT DISTINCT city FROM Student;
```



Select

1. Get the average score for all students
`SELECT avg(score) FROM Enrolls_in;`
2. Get the average score in class c1
`SELECT avg(score) FROM Enrolls_in WHERE class = "c1";`
3. Get the average score for each class
`SELECT class, avg(score) FROM Enrolls_in GROUP BY class;`
4. Get the average score for each class with more than one student
`SELECT class, avg(score) FROM Enrolls_in GROUP BY class
HAVING count(score) > 1;`



The screenshot shows a SQL query editor window titled "school_data". The query entered is `SELECT class, avg(score) from Enrolls_in GROUP BY class;`. The results are displayed in a table with two columns: "class" and "avg(score)". The table contains three rows of data.

class	avg(score)
c1	72
c2	95
c3	44

Select

1. Get the average score for each class and name the column average

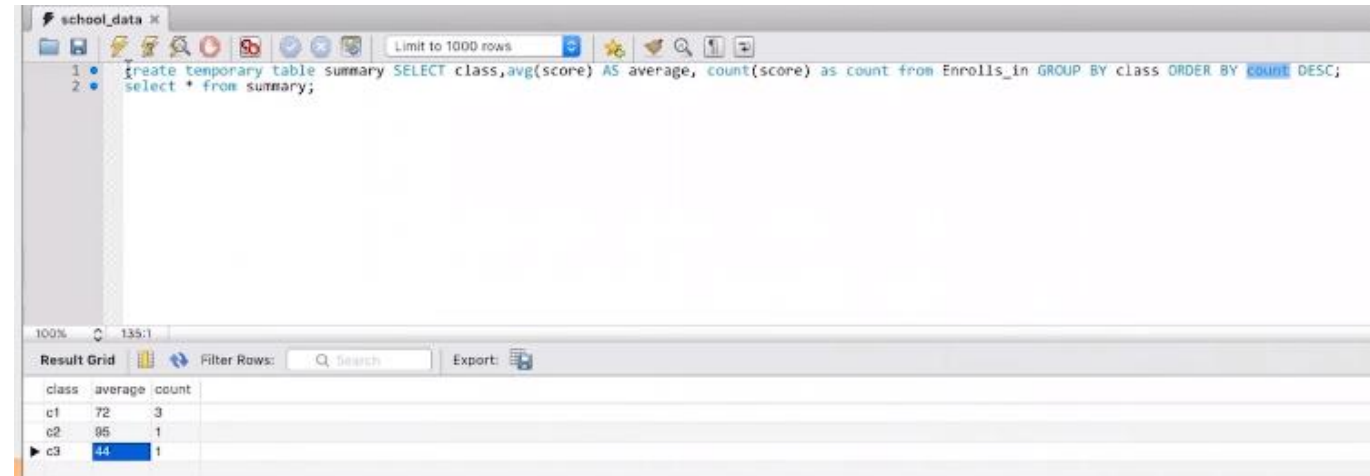
```
SELECT class, avg(score) AS average FROM Enrolls_in  
GROUP BY class;
```

2. Get the average score for each class, name the column average, and store the result set in a temporary table "averages"

```
CREATE TEMPORARY TABLE averages SELECT class,  
avg(score) AS average FROM Enrolls_in GROUP BY class;
```

3. Get the average score and number of students for each class, name the columns average and count, sort them in ascending order of count, and store the result set in a temporary table "summary"

```
CREATE TEMPORARY TABLE IF NOT EXISTS summary  
SELECT class, avg(score) AS average, count(score) AS count  
FROM Enrolls_in GROUP BY class ORDER BY count;
```



```
1 • create temporary table summary SELECT class, avg(score) AS average, count(score) as count from Enrolls_in GROUP BY class ORDER BY count DESC;  
2 • select * from summary;
```

class	average	count
c1	72	3
c2	95	1
c3	44	1

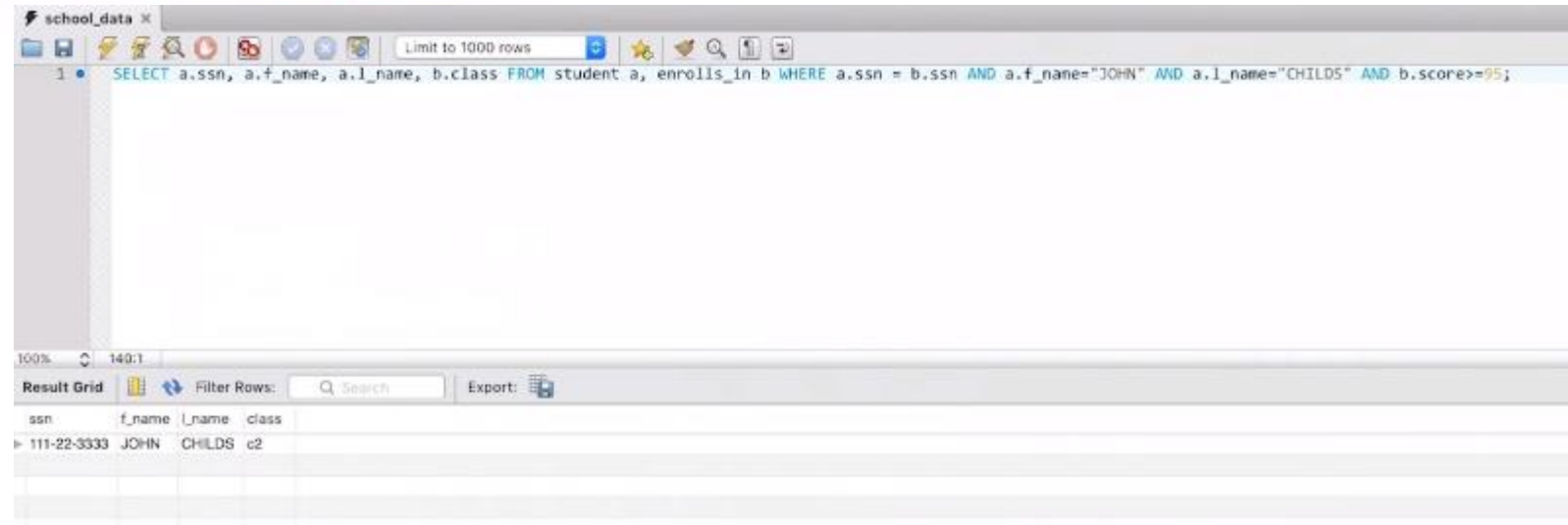
1. Get a list of student names and classes

```
SELECT f_name,l_name,class FROM Student a, Enrolls_in b  
WHERE a.ssn = b.ssn;
```

2. Get the list of classes that John Childs is enrolled in and his grade is 95 or greater

```
SELECT class FROM Student a, Enrolls_in b WHERE  
a.f_name = "JOHN" AND a.l_name = "CHILDS" AND a.ssn =  
b.ssn AND b.score >=95;
```

3. Get the names of the courses that John Childs is enrolled in
SELECT name from Student a, Enrolls_in b, Course c WHERE
a.f_name = "JOHN" AND a.l_name = "CHILDS" AND a.ssn =
b.ssn AND b.class = c.number;



The screenshot shows a database query tool interface. The query editor at the top contains the following SQL query:

```
1 • SELECT a.ssn, a.f_name, a.l_name, b.class FROM student a, enrolls_in b WHERE a.ssn = b.ssn AND a.f_name="JOHN" AND a.l_name="CHILDS" AND b.score>=95;
```

Below the query editor, the results are displayed in a table. The table has four columns: `ssn`, `f_name`, `l_name`, and `class`. The first row of data shows the student John Childs enrolled in class c2.

ssn	f_name	l_name	class
111-22-3333	JOHN	CHILDS	c2

Join

The process of combining rows from two or more tables

Course

number	name	room
c1	Data analytics	1127
c2	Python	303
c3	corp fin	331
c4	prod. mgmt	1127
c5	Ethics	303
c6	leadership	303
c7	bus analytics	1127

Enrolls-in

ssn	class	score
111-22-3333	c1	93
123-12-1234	c1	87
111-22-3333	c2	95
222-33-4455	c3	44
555-11-7777	c1	36

rows in Enrolls-in are matched with rows in course where class has the same value as number

We can combine rows in Course with rows in Enrolls_in using an "INNER JOIN"

```
SELECT number, name, room, ssn, score FROM course  
INNER JOIN Enrolls_in ON course.number = enrolls_in.class;
```

Joins can be explicit

```
SELECT number, name, room, ssn, score FROM course  
INNER JOIN Enrolls_in ON course.number = enrolls_in.class;
```

or implicit

```
SELECT number, name, room, ssn, score FROM  
course,enrolls_in WHERE course.number = enrolls_in.class;
```

Get the names of the courses that John is enrolled in

```
SELECT course.name FROM student  
INNER JOIN enrolls_in ON student.ssn = enrolls_in.ssn  
INNER JOIN course ON course.number = enrolls_in.class  
WHERE f_name = "JOHN";
```

DB API

- ➡ Set of standards for Python Database API
- ➡ Import the API module
- ➡ Acquire a connection with the database server
- ➡ Issue sql commands or call stored procedures
- ➡ Close the connection

pymysql (mysql)

- ➡ Import the API module
 - * `import pymysql`
- ➡ Acquire a connection with the database server
 - * `db = pymysql.connect("localhost","root","None")`
- ➡ Prepare a cursor object
 - * `cursor = db.cursor()`
 - * `cursor.close()`

First import the python module containing the API

```
In [ ]: import pymysql
```

Set up a connection and create a cursor object

```
In [ ]: db = pymysql.connect("localhost","root","None",database="schooldb")
        cursor = db.cursor()
```

Execute a query and get the results

```
In [ ]: cursor.execute('show tables;')
        cursor.fetchall()
```

```
In [ ]: query = """
        SELECT course.name FROM student
        INNER JOIN enrolls_in ON student.ssn = enrolls_in.ssn
        INNER JOIN course ON course.number = enrolls_in.class
        WHERE f_name = "JOHN";
        """
        cursor.execute(query)
        cursor.fetchall()
```

pymysql (mysql)

➡ Issue sql commands

- * cursor.execute('show tables;')

➡ Get results

- * cursor.fetchone()

- * cursor.fetchall()

➡ Close the connection

- * cursor.close()

