

University of Warwick
Department of Computer Science

CS258

Database Systems



Cem Yilmaz
April 17, 2023

Contents

1	Introduction to Data	2
2	Relational Model	2
2.1	Key	2
2.2	Table	2
2.3	Formal Definition	2
2.4	Key Constraints	3
2.5	Entity Integrity Constraints	3
2.6	Referential Integrity Concerns cross-table relationships	4
2.7	Referential Integrity Constraints	4
2.8	Possible Violations for each operation	4
2.9	Functional Dependencies	5
2.10	Schemas	5
2.11	Constraints	5
3	Queries	6
3.1	Basic Queries	6
3.1.1	Ordering	6
3.1.2	Combining	6
3.2	More Complicated Queries	6
3.2.1	String Matching	6
3.2.2	NULL	7
4	Database Modifications	7
4.1	Insertion	7
4.2	Deletion	7
4.3	Updating	7
5	Programming DB	7
5.1	JDBC: Java Data Base Connectivity	7
5.2	Programming	8
5.2.1	Connection	8
5.2.2	Statements and Queries	9
5.2.3	Updates	9
5.3	Errors	10
5.3.1	Connections	10
5.3.2	Statements	10
5.3.3	Updating INSERT	10
5.3.4	Updating TUPLE	11
5.3.5	Querying	11
6	Advanced SQL	11
6.1	Multiple Tables	11
6.2	NATURAL JOIN	12
6.3	JOIN	13
6.4	INNER JOIN	13
6.5	OUTER JOIN	13

1 Introduction to Data

Data comes into two different kinds: structured and unstructured.

- Structured - e.g., tables with predefined columns
- Unstructured - e.g., web pages, text docs, images, videos...

Semi-structured data also exists and is found within XML and JSON. Traditionally, these are structured databases and were defined to be so.

Definition 1.1. Relation

Informally, a relation is a table of values having:

- A set of rows. The data elements in each row represent certain facts that correspond to a real-world entity or relationship. In the formal model, rows are called tuples.
- Each column represents a characteristic / attribute of interest of that entity. Has a column header that gives an indication of the meaning of the data items in that column. In formal model, column header is called an attribute name.

2 Relational Model

2.1 Key

Each row must be uniquely identifiable in the table. The key of the row does this. Sometimes row-ids or sequential numbers are assigned as keys to identify the rows in a table. This is called an artificial key or surrogate key.

2.2 Table

A table is just an acceptable visual representation of the mathematical notion of relation. To formulate queries, we specify that table name(s), and attributes names of interest and special constraints (aka predicates) that need to be satisfied in order for a data item (=row) to be of interest.

2.3 Formal Definition

The schema of a relation is denoted by

$$R(A_1, A_2, \dots, A_n)$$

where R is the name of the relation

The attributes / columns of the relations are denoted by A_1, A_2, \dots, A_n . For example,

$$CUSTOMER(CUST - ID, CUST - NAME, ADDRESS, PHONENO)$$

Each Attribute/column has a domain or a set of valid values. For example, CUST-ID is a 7 digit number. More about this is discussed in section 2.10.

A tuple (aka row) of a relation is an ordered set of values enclosed in angled brackets $\langle \dots \rangle$. Each value is derived from an appropriate domain. For example,

$$\langle 632895, \text{"PETER T."}, \text{"2 Main St. Warwick"}, \text{"(024)894 - 2000"} \rangle$$

A domain has a logical definition in the real world. A domain also has data format. For example, USA_phone_numbers may have a format: $(ddd)ddd - dddd$ where each d is a decimal digit. The attribute name designates the role played by a domain in a relation. Used to interpret the meaning of the data elements corresponding to that attribute. Example, the domain Date may be used to define two attributes named "invoice-date" and "payment-date" with different meanings.

Definition 2.1. Relation State

The relation state R is the set that contains all the set of tuples in the relation. The relation state is a subset of the Cartesian product of the domains of its attributes. To put it all together,

$R(A_1, A_2, A \dots, A_n)$ is the schema of the relation

R is the name of the relation

A_1, A_2, \dots, A_n are the attributes (columns) of the relation

$r(R)$ is a specific state or instance of relation

R is an actual set of tuples (rows)

$r(R) = \{t_1, t_2, \dots, t_n\}$ where each t_i is an n -tuple

$t_i = \langle v_1, v_2, \dots, v_n \rangle$ where each v_j comes from $dom(A_j)$

$r(R) \subset dom(A_1) \times dom(A_2) \times \dots \times dom(A_n)$

To create a table, we would write

```
1 CREATE TABLE Drinkers (  
2   name CHAR(31) Primary Key,  
3   Addr CHAR(50) DEFAULT '123 Sesame St',  
4   phone CHAR(16) NOT NULL  
5 );
```

Listing 1: Creating a table

The constraints include the fact that there are key constraints, entity integrity constraints and referential integrity constraints. Also there exist domain constraints, values in an attribute in a tuple must come from the domain of that attribute. Values could also be NULL if allowed.

2.4 Key Constraints

Superkey of R : a subset of attributes of R , SK , such that:

In any valid state for $r(R)$: for all two distinct tuples t_1 and $t_2 \in r(R)$, where $t_1[SK] \neq t_2[SK]$, where SK is the superkey. For example, let us define superkey SK with the attributes

$$SK = \{\text{Country, PhoneNumber}\}$$

then, no two entries country and phone number coincides.

A candidate key of R is a minimal super for any key K if:

- The removal of any attribute from K results in a set of attributes that is no longer a super key.

Definition 2.2. Superkey

Superkey is for all two tuples t_n and t_k in $r(R)$, there does not exist $t_n[SK] = t_k[SK]$.

Definition 2.3. Candidate Key

Otherwise known as the minimal key, it is by definition a super key that for which, if any attribute is removed from its definition of SK , it is no longer a super key.

2.5 Entity Integrity Constraints

The primary key PK cannot be NULL in any tuple of $r(R)$. This is because primary key values are used to identify the individual tuples. If PK has several attributes, null is not allowed as a value in any of these attributes. Also, any attribute of R even non-key may not be allowed to be NULL.

2.6 Referential Integrity Concerns cross-table relationships

A table students lists data about students, such as age, grades, etc. Students has an attribute STUDENTID, and suppose it is the primary key. Then, another table student-courses has also an attribute STUDENT ID, listing, for each module, which student are taking it, at which term etc. One could argue that all this info should be in one table, however, that is not a good idea and will be explained later. A logical DB Design spreads this information across tow tables, it is fundamental to R-DBMS design. Tables are "linked" according to references between them. Care must be exercised which such cross-table relationships are exercise.d.

2.7 Referential Integrity Constraints

A set of attributes FK from relation R_1 is a foreign key references relation R_2 : Attributes in the FK from R_1 have the same domain as the attributes of in the primary key PK of R_2 , and a value of FK in a tuple t_1 of R_1 must either:

- refer to a value of the PK of some tuple t_2 in R_2 ,
- be NULL

Formally,

- t_1 in R_1 references t_2 in R_2 if $t_1[FK] = t_2[PK]$ or
- t_1 in R_1 makes no reference if $t_1[FK] = NULL$

Table 1: R_1

PK	FK
1	1
2	NULL
3	2

Table 2: R_2

PK	data
1	alpha
2	beta
3	gamma

Where in this case. $1 \rightarrow 1, 2 \rightarrow 2$ from R_1 FK to R_2 PK . All of these constraints are expressed in create table in SQL. There also exist "semantic attribute integrity constraints" where, for example, things like bank balance > 0 are expressed.

2.8 Possible Violations for each operation

- Domain - one of the attribute values for new tuple is not in the attribute domain
- Key - the value of a key attribute in new tuple already exists
- Referential integrity - a foreign key value in new tuple references a primary key value that does not exist in referenced relation
- Entity integrity - if primary key value is null in new tuple.
- Delete - may violate only referential integrity, if primary key value of the tuple being deleted is referenced from other tuples in other relations. Some option must be specified during database design for each foreign key constraint on how to handle such deletions leading to relational integrity violations.
- UPDATE - obviously the domain and NOT NULL constraints may be violated on an attribute being modified. Furthermore, updating the primary key duplicates, updating a foreign key may violate referential integrity. Updating an ordinary attribute (not PK or FK) can violate only domain constraints.

2.9 Functional Dependencies

We need functional dependencies as they are a formal tool that allow us derivation of 'good' DB designs (relations and their schemata). Ultimately, a good design depends on the dependencies between attributes between to be in the same relation. Assume, X, Y, Z represent sets of attributes. A, B, C represent single attributes. Notation ABC implies $\{A, B, C\}$. $X \rightarrow Y$ is an assertion about a relation R and two sets of attributes from R .

If $X \rightarrow Y$ the values of the Y component of a tuple depend on values of the X component. i.e., the values of the X component of a tuple uniquely (functionally) determine those of the Y component.

So, $X \rightarrow Y$ specifies that whenever two tuples R agree on values of all the attributes of X , they must also agree on values of attribute of Y . Formally, $t_1[X] = t_2[X] \rightarrow t_1[Y] = t_2[Y]$.

2.10 Schemas

A db contains one or more schemas. Each schema contains one or more tables. To Create a schema, we type

```
CREATE SCHEMA company AUTHORIZATION 'Jsmith';
```

Creating a schema "company" owned by "Jsmith". We can create tables for each schema. The command is

```
CREATE TABLE
```

which creates a relation R . However, we must also specify the name of the table, its attributes (columns) along with their types, and any constraints. All created tables belong in a schema and belong to a special "public" schema if no other has been created. Relations and its tuples are actually created and stored as a file by the DBMS. Virtual relations/tables which are created using VIEW do not correspond to a physical file. The attribute/constraint format for a relation R is as the following:

```
(< Att1Name >< Att1DataType >< AttConstraints >, < AttName2 >< Att2DataType >< AttConstraints >, ...)
```

Example 2.1.

```
CREATE TABLE students(studentID INTEGER PRIMARY KEY, studentName VARCHAR(30), courseID INTEGER);
```

2.11 Constraints

Some basic constraints include:

- Key constraint: A primary key is unique
- Entity integrity: a primary key value cannot be null
- Referential integrity constraints: a "foreign key" must have a value that is already present as a primary key, or it may be null
- Plus attribute constraints

Additional restrictions include NOT NULL, default value of an attribute i.e., DEFAULT< value >, CHECK i.e., CHECK(Age > 0 AND Age < 125), PRIMARY KEY, UNIQUE. The foreign key constraint, as we know, must have that it is a subset of the primary key it is referencing to. As such, with updates and deletions, we may have errors. To remedy this, we have the commands

```
ON DELETE and ON UPDATE
```

which can be paired with the commands

```
DEFAULT, CASCADE, SET NULL
```

Cascade refers to the deletion of child tables foreign keys where $FK = PK$.

3 Queries

3.1 Basic Queries

A query is done using the command

```
SELECT < attribute1, ... attributei > FROM < table list > WHERE < condition >
```

One can use the special operator `*` in the attribute which retrieves all attributes available. If there are tables with the same name, then we can specify the table and attribute by adding the table name with a dot before the attribute. For example,

```
SELECT students.studentID from students;
```

You can further alias names of tables or even attributes/columns, for example, below we will rename the table

```
SELECT s.studentID FROM students AS S;
```

which means that the name *S* is now used for any qualification for the table name (including in `WHERE` clauses). The logical comparison operators that exist are:

`=, <, <=, >, >=, <>`

However, when using `=` to filter by strings, it is important to put the string around `'`. One could also use

```
BETWEEN low-value AND high-value
```

When querying, SQL first uses the `FROM` clause and selects the table. Then, it applies `WHERE` to find the correct rows in the table. Finally, it applies `SELECT` to filter the columns/attributes.

3.1.1 Ordering

You can further order your final result by adding

```
ORDER BY attribute(s) [ASC|DESC]
```

in the end of your query. You can also mathematical expressions within the query. You could, for example, retrieve the price in `SELECT price*147` which would multiply all prices by 147.

3.1.2 Combining

You can further complicate and define new columns, such as

```
SELECT drinker, 'likes Earl Grey' AS whoLikesEarlGrey FROM Likes WHERE tea = 'Earl Grey';
```

This would define a new column `whoLikesEarlGrey` with the string `'likes Early Grey'` for all rows. You could create a more sophisticated column by placing functions such as string count for the drinker name.

3.2 More Complicated Queries

3.2.1 String Matching

the `WHERE` clause can have conditions for strings, in which a string is compared with a pattern to see if it matches. The general form is

```
< StringAttribute > LIKE < pattern >
< StringAttribute > NOT LIKE < pattern >
```

The pattern is a quoted string with special characters where

```
% = "any string" of zero or more characters
_ = "any character", exactly one character
```

For example, the query

```
SELECT name FROM drinkers WHERE phone LIKE '%1926_____'
```

Would create a query that finds drinkers whose phone has the dialling code 1927. The `%` accommodates for possible country codes e.g. `+44`, `+90` and the specific number of `_` ensures the length of the phone has exactly 6 characters after 1926.

3.2.2 NULL

We need to be able to test if a value is NULL. The syntax is as follows:

`< Attribute > IS (NOT) NULL`

The (NOT) was placed to emphasise that one could place the word NOT if desired.

4 Database Modifications

4.1 Insertion

`INSERT INTO < tablename > VALUES (< val1 >, < val2 >, ... < valn >);`

Each list refers to a record (tuple) being inserted. Each value belongs to the attribute of the same position in the attribute list of the schema. You can also insert partial records by specifying which attributes you would like to insert. That is,

`INSERT INTO < tablename > < A1, ..., Aj > VALUES (< val1 >, < val2 >, ... < valn >);`

where `< A1, ..., Aj >` is a partial attribute list

4.2 Deletion

`DELETE FROM < tableName > WHERE < expression >;`

For example,

`DELETE FROM students WHERE studentID = 2000;`

would delete the student with the studentID 2000. Note that it has to pass referential integrity checks.

4.3 Updating

`UPDATE < relation > SET < list of (attribute assignment) pairs > WHERE < condition on tuples >;`

For example,

`UPDATE drinkers SET phone = '09126 - 554555' WHERE name = 'Peter';`

5 Programming DB

5.1 JDBC: Java Data Base Connectivity

The JDBC is an API that calls SQL commands from Java. It is implemented as a driver for a particular DBMS with methods, objects and classes. One Java program can have connections to several databases and possibly even connected to different DBMSs. These are called data sources.

The driver manager class is used to handle multiple connections to different DBs i.e., managing their drivers, and mapping clients to right drivers. The methods include, but are not limited to:

- `getDriver`
- `registerDriver`
- `deregisterDriver`

There are several driver types for different types of support and we mostly refer to type 4 drivers, which are pure Java driver for direct-to-database. Fundamentally, there are 3 key tasks carried out with the help of important new Classes/ objects.

1. The Connection class

- A connection object is needed per DB connection
- Created using DriverManager's getConnection()
- Has functions to connect to specified DB sources using URLs

2. The Statement class with two subclasses

- PreparedStatement and CallableStatement which have methods to execute an SQL command.
- Need a statement object per SQL query / update you wish to issue

3. The ResultSet class

- Query results are returned into objects of this type
- It is akin to tables, where result tuples are rows and columns are attributes

5.2 Programming

5.2.1 Connection

The import the JDBC class library we write

```
6 import java.sql.*
```

Listing 2: JDBC import

We then load explicitly the JDBC driver: use a generic Java function for loading class e.g.,

```
7 Class.forName("oracle.jdbc.driver.OracleDriver"); // for Oracle
8 Class.forName("org.postgresql.Driver"); // for Postgresql
```

Listing 3: Loading JDBC drivers

This register the driver with the driver manager and makes it available to the program. We now create vars for connection credentials, that is

```
9 String dbacct, passwd, ssn, name;
10 dbacct = readentry("Enter database account:");
11 passwd = readentry("Enter password:");
```

Listing 4: Creating connection credentials

Although this can also be replaced with the Scanner class to read any input from the console. Finally we use a connection object to connect to DB, that is

```
12 Connection conn = DriverManager.getConnection("jdbc:oracle:oci8:" + dbacct + "/" + passwd);
```

Listing 5: Connection to DB

After connection, it might be wise to check whether the Postgres Driver can be loaded by catching. In particular,

```
13 try {
14     Class.forName("org.postgresql.Driver");
15 }
16 catch (ClassNotFoundException x) {
17     System.out.println("Driver could not be loaded");
18 }
```

Listing 6: Catching if driver can be loaded

Example 5.1. Connecting to a local DB server

```
19 String username = "me";
20 String password = "mypassword";
21 String url = "jdbc:postgresql://127.0.0.1:5432/mydb?user=" + username + "&password=" +
    password;
22 Connection conn = DriverManager.getConnection(url);
```

Listing 7: connecting to a local DB server

5.2.2 Statements and Queries

We now have loaded drivers and set up the DB connection. We now use a Statement object to hold a query. We use PreparedStatement if the query will be issued many times, that is, it is prepared, checked and compiled only once to save processing. It is used as such:

```
23 String stmt1 = "SELECT course, year FROM students WHERE Name = ? AND age = ?";
24 PreparedStatement p = conn.prepareStatement(stmt1);
25 p.setString(1, "John"); // says that the first ? is John as string
26 p.setInt(2, 20); // says that the second ? is 2 as int
```

Listing 8: Statement Query

After preparation, we can now run our statement using `executeQuery`. For every stmt object, there are 2 functions: `executeQuery` which returns an object of type `ResultSet` for the `SELECT` command and `executeUpdate` which is for `INSERT`, `UPDATE`, `DELETE`, `CREATE` and `DROP`, returning an integer for the number of affected tuples. The command to run our earlier statement is

```
27 ResultSet r = p.executeQuery()
```

Listing 9: executing the query

The object type of `ResultSet` contains the set of result. After executing `r` call, `r` refers to an initial position before the 1st tuple count in result. If `r` is of type `ResultSet`, we can scroll through the results using `r.next()`. This will return a `NULL` in the end as the last result. In particular

```
28 while (r.next()) {
29     course = r.getString(1); // The 1 refers to the attribute position within the tuple as
        defined in our query
30     year = r.getInt(2); // The 2 refers to the attribute position within the tuple as defined in
        our query
31     System.out.println(course + "," + year );
32 }
```

Listing 10: Iterating through ResultSet

5.2.3 Updates

For creating tables, inserting and deleting tuples, etc. we first need a statement object, as per queries, over an established connection `conn`

```
33 Statement stmt = conn.createStatement();
```

Listing 11: Connection

We then use the `executeUpdate` method, i.e.,

```

34 String sql = "CREATE TABLE Students " +
35             "(ID INTEGER NOT NULL, " +
36             "Name VARCHAR(25), " +
37             "gpa REAL DEFAULT 0.0, " +
38             "PRIMARY KEY (ID))";
39 stmt.executeUpdate(sql)

```

Listing 12: Table Creation

You can further do other commands, such as INSERT by modifying the sql string

```

40 String sql = "INSERT INTO STAFF VALUES " + "(1234567, 'Peter')";

```

Listing 13: Insert

It is also possible to use a prepared statement for the actions described above. If done so, instead of *ResultSet* $r = p.executeQuery()$ would be *introws* = $p.executeUpdate()$.

5.3 Errors

As specified before, one of the best errors we can implement is to see if the driver can be loaded earlier in the notes. However, managing errors with JDBC can get more complicated for different functions

5.3.1 Connections

When creating connections, we can check if it has been established with

```

41 try {
42     connection conn = DriverManager.getConnection(URLstring);
43 }
44 catch (SQLException se) {
45     System.out.println("Could not open connection with connection string " + URLstring);
46     se.printStackTrace();
47 }

```

Listing 14: Connection Error Catching

5.3.2 Statements

To catch errors using statements we use

```

48 try {
49     Statement stmt = conn.createStatement();
50 }
51 catch (SQLException se) {
52     System.out.println("Could not create statement ");
53     se.printStackTrace();
54 }

```

Listing 15: Statements Error Catching

5.3.3 Updating INSERT

```

55 try {
56     String sql = "INSERT INTO Students VALUES " + "(1234567, 'Andrew H.')";
57     stmt.executeUpdate(sql);
58 }
59 catch (SQLException se) {

```

```

60 System.out.println("Could not insert tuple for Andrew H.");
61 se.printStackTrace();
62 }

```

Listing 16: Update INSERT Error Catching

5.3.4 Updating TUPLE

```

63 try {
64     int count = stmt.executeUpdate("UPDATE Students SET Name = 'Andrew Hague' WHERE ID = 1234567"
65 );
66 catch (SQLException se) {
67     System.out.println("Could not update Andrew H.");
68     se.printStackTrace();
69 }

```

Listing 17: Update UPDATE Error Catching

5.3.5 Querying

```

69 try {
70     sql = "SELECT * FROM Students WHERE ID < 5555555";
71     rs = stmt.executeQuery(query);
72 }
73 catch (SQLException se) {
74     System.out.println("Could not get result of query " + sql);
75     se.printStackTrace();
76 }

```

Listing 18: Query Error Catching

6 Advanced SQL

6.1 Multiple Tables

If we have multiple tables in FROM clause, then we compute the cartesian product of the tables, which produces a new table. And then we scan this table, one a time. You can think of this as a nested loop: for each tuple in the 1st table, concatenate it with each tuple of the 2nd table, in turn we achieve an operation called CROSS JOIN.

```
SELECT foods.item_name,foods.item_unit,
company.company_name,company.company_city
FROM foods
CROSS JOIN company;
```

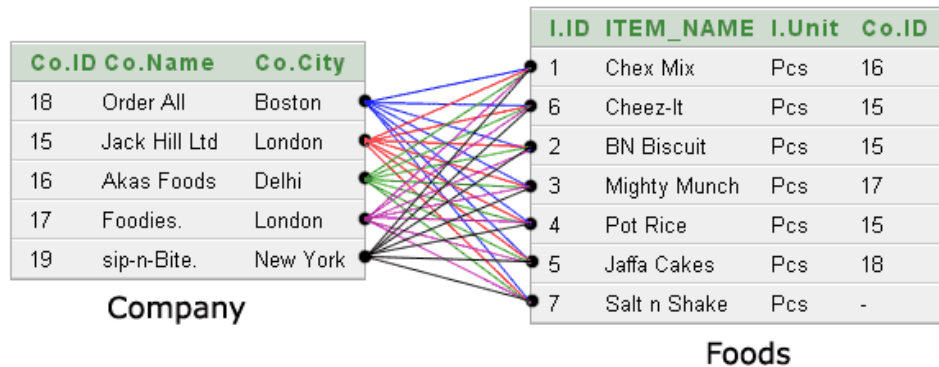


Figure 1: Cross Join Table Example

Follow the coloured lines above to see how it cross joins. Each edge is a newly created tuple in the new table. To specify CROSS JOIN between two tables R and S , we write

$R \text{ CROSS JOIN } S$

And similarly, for NATURAL JOIN

$R \text{ NATURAL JOIN } S$

Recall that queries create a table, therefore, we can also write

$query_1 \text{ CROSS JOIN } query_2$
 $query_1 \text{ NATURAL JOIN } query_2$

6.2 NATURAL JOIN

$R \text{ NATURAL JOIN } S$ is equivalent to formatting the cross product of R and S and then looking for attributes of R that have the same name as attributes in S . Keep only tuples from product whose same-name-type attributes have same value. Only one of each set of duplicate columns is kept.



Figure 2: NATURAL JOIN

6.3 JOIN

This is an extension to natural joins and it allows us to combine constraints (in WHERE) and a product in 1 operations. The syntax is

$R \text{ JOINS } S \text{ ON } \langle \text{condition} \rangle$

$\langle \text{condition} \rangle$ is used for selecting tuples. If operator is = this is called an equi-join. For example,

Drinkers D JOIN Frequents F ON D.drinker = F.drinker

returns (d, a, d, b) quadruples such that drinker d lives at address a and frequents teaRoom b . The default join operator we have seen thus far is called inner join, for which there exist an explicit operator

SELECT $\langle \text{attList} \rangle$ FROM $\langle \text{table1} \rangle$ INNER JOIN $\langle \text{table2} \rangle$ ON $\langle \text{expr} \rangle$;

An inner join will combine the rows from $\langle \text{table1} \rangle$ with the rows from $\langle \text{table2} \rangle$ based on satisfying the ON boolean expression. The expression usually takes the form of a comparison between a column from table 1 and a column from table 2.

6.4 INNER JOIN

The inner join has the same functionality as

SELECT * FROM student S, studentCourses C WHERE S.studentID = C.studentID

Explicitly using "Join" makes formatting more readable. Using 'ON' separates clauses that deal with sources of data i.e., belonging to FROM away from CLAUSES that deal with selection criteria

6.5 OUTER JOIN

Inner joins can "lose" information, become a tuple that doesn't join with any other relation disappears. Instead, outer join keeps this information and fills tuples that do not match with nulls if it is not matched when joining. You can further specify if you want to keep information on the left side or the right side, more specifically if you wanna keep information on the left table, then you type LEFT OUTER JOIN. You can further use the syntax *NATURAL* to make the join natural, or use *ON* syntax after the OUTER JOIN to specify when to join.