

Table of Contents

1. Business Rules and Table Design.....	1
2. Relationships and Referential Integrity.....	4
2.1. Type of relationships.....	4
2.2. Referential Integrity rule.....	4
3. Normalization.....	5
3.1. Goals of normalization	5
3.2. Anomalies.....	5
4. Examples of data that is not normalized.....	6

Before we cover the syntax for creating tables and setting constraints, it is time to talk a bit more about relations and relationships. Relations is the formal term for how data is structured; we visualize relations as tables. Relationships are the way that tables relate to each other and are essential to understanding how a database works. So we will also look at a few more definitions and considerations for relationships.

As we talk about the design and structure of tables we are discussing metadata. Metadata is the data that refers to the structure of the database, such as the names of the tables and their attributes, the data types of the attributes, and the primary keys. It does not refer to the values (such as the last name of a customer is Jones); instead it refers to the attribute that we use to store that value (such as the last name can store up to 25 characters and cannot be null).

1. Business Rules and Table Design

The ultimate design of a database and its tables comes from the business rules of the company. Business rules are statements about the way a company structures its data and controls its operations. For example, in our vets database there is a rule that each animal is the responsibility of a single client. This business rule dictates that each animal row in the animal table is associated with one client row and, therefore, `cl_id` is an attribute of the animals table.

Many poorly designed databases have poorly designed tables. Most problems arise from having tables that are too complex—that have too many attributes that do not belong to that table. This is often a result of thinking of the data as a manual table or an Excel spreadsheet. Since manual tables for storing data do not have automatic look-up features, we tend to put all of the data into one big table so that we can find it. This same attitude is often found when people create spreadsheets. We might have a spreadsheet where each row for an animal includes the client id but also their name and address. We would not want to bring that design directly into a database table.

Suppose we had the following single table for animals and clients.

An_id	Type	Name	dob	Cl_id	Last Name	First Name	Address	City	State	Zip
11015	snake	Kenny	2005-10-23	4534	Montgomery	Wes	POB 345	Dayton	OH	43784
11029	bird		2005-10-01	4534	Montgomery	Wes	POB 345	Dayton	OH	43874
12035	bird	Mr Peanut	1995-02-28	3560	Monk	Theo	345 Post Street	New York	NY	10006
12038	cat	Gutsy	2007-04-29	3560	Monk	Theo	3405 Post Street	New York	NY	10006
15001	turtle	Big Mike	2008-02-02	25	Harris	Eddie	2 Marshall Ave	Big Rock	AR	71601
15341	turtle	Pete	2007-06-02	93	Wilson	Sam	123 Park Place	Big Rock	AR	71601
15002	turtle	George	2008-02-02	25	Harris	Eddie	2 Marshall Ave	Big Rock	AR	71601

16002	cat	Fritz	2009-05-25	25	Harris	Eddie	2 Marshall Ave	Big Rock	AR	71601
21002	snake	Edger	2002-10-02	25	Harris	Eddie	2 Marshall Ave	Big Rock	AR	71601
21004	snake	Gutsy	2001-05-12	25	Harris	Eddie	2 Marshall Ave	Big Rock	AR	71601
21005	cat	Koshka	2004-06-06	25	Harris	Eddie	2 Marshall Ave	Big Rock	AR	71601
21006	snake		1995-11-02	25	Harris	Eddie	2 Marshall Ave	Big Rock	AR	71601
				5686	Biederbecke	NULL	50 Phelan Ave	San Francisco	CA	94112
				5689	Biederbecke	NULL	50 Phelan Ave	San Francisco	CA	94112
				5698	Biederbecke	Sue	50 Phelan Ave	San Francisco	CA	94112

The spreadsheet approach leads to two problems:

- having rows of data with a lot of empty cells
- data redundancy, having the same data value stored more than once.

In the above table we have several rows which have no data values for the animal columns because these clients don't have any animals. But if we are storing our animal and client data in this table we must have rows for every client even if they do not have any animals.

We also are repeating a lot of values which leads to more problems- did you notice the two errors in the data in the table?

Data redundancy leads to problems. One of the first problems we tend to think of is that we are wasting storage space- we do not need to store client 25's name and address 7 times just because he has 7 animals. But storage space is pretty cheap. The more important problems are the logical ones.

- update anomalies where one copy of the data value is changed but other values of the same data item are not changed
- delete anomalies where deleting one data value results in too much data being deleted
- insert anomalies where the user cannot insert the data values they wish to store.

Suppose client 25 moves- how many places do we have to update that data? Are you certain that every application that updates data will update every one of those rows?

If animal 15341 dies and we want to remove its row, we might also remove the only row we have for client 93.

Suppose we need a rule that the animal id and type and dob could not be null. How do we enforce this rule and still allow the table to include clients with no animals?

These are problems that were a major consideration with traditional file processing and database systems were suppose to help solve these problems.

With a DBMS you can split the data into smaller, more tightly focused tables and then reassemble the big collections of data when needed through queries. The general problem in going from a big table to properly designed tables is to create a collection of related tables without losing any of the information contained in the original collection of data. Note that we have to preserve the table information-not just the table data. Table information includes facts imbedded in the table design.

At this point in the semester, these should seem natural- the client table includes the client name and address and each client gets one name and one address. And the animal table gets one row for each animal.

Cl id	Last Name	First Name	Address	City	State	Zip
25	Harris	Eddie	2 Marshall Ave	Big Rock	AR	71601
93	Wilson	Sam	123 Park Place	Big Rock	AR	71601

3560	Monk	Theo	3405 Post Street	New York	NY	10006
4534	Montgomery	Wes	POB 345	Dayton	OH	43784
5686	Biederbecke	NULL	50 Phelan Ave	San Francisco	CA	94112
5689	Biederbecke	NULL	50 Phelan Ave	San Francisco	CA	94112
5698	Biederbecke	Sue	50 Phelan Ave	San Francisco	CA	94112

An_id	Type	Name	dob	Cl_id
11015	snake	Kenny	2005-10-23	4534
11029	bird		2005-10-01	4534
12035	bird	Mr Peanut	1995-02-28	3560
12038	cat	Gutsy	2007-04-29	3560
15001	turtle	Big Mike	2008-02-02	25
15341	turtle	Pete	2007-06-02	93
15002	turtle	George	2008-02-02	25
16002	cat	Fritz	2009-05-25	25
21002	snake	Edger	2002-10-02	25
21004	snake	Gutsy	2001-05-12	25
21005	cat	Koshka	2004-06-06	25
21006	snake		1995-11-02	25

You might notice that often well designed tables have fewer columns.

If it were important to allow more than one address for a client, then we could split the client table into multiple tables. One table for the client name (each client gets one name) and a second related table for the client addresses. But those tables would also each need the client id to link back to the main client table. This requires understanding the business rules and policies reflected in the data collection.

Now some clients can have multiple addresses

Cl_id	Last Name	First Name
25	Harris	Eddie
93	Wilson	Sam
3560	Monk	Theo
4534	Montgomery	Wes
5686	Biederbecke	NULL
5689	Biederbecke	NULL
5698	Biederbecke	Sue

Cl_id	Address	City	State	Zip
25	2 Marshall Ave	Big Rock	AR	71601
25	2957 Grover Ave	Springfield	IL	61258
93	123 Park Place	Big Rock	AR	71601
93	56 Meadowland Ln	Springfield	NY	10027
3560	3405 Post Street	New York	NY	10006
4534	POB 345	Dayton	OH	43784
5686	50 Phelan Ave	San Francisco	CA	94112
5689	50 Phelan Ave	San Francisco	CA	94112
5698	50 Phelan Ave	San Francisco	CA	94112

2. Relationships and Referential Integrity

The following are rules and considerations that you would need to think about when you set up tables and constraints.

2.1. Type of relationships

There are three types of relationships:

- One to One
- One to Many
- Many to Many

The relationship types indicate how many rows in one table can be associated with a row in another table.

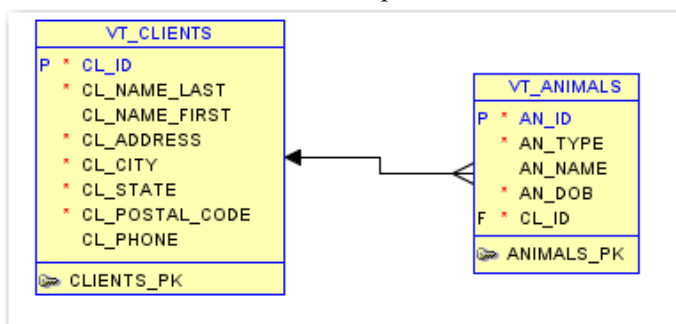
One-to-One- this means that a row in table A is related to at most one row in table B. Suppose we wanted to keep a photo of each member of our staff. Assuming we had a way to store a photo in our database, we might want to set up a second table vt_staffPhoto (staff_id, Photo) since a photo would take a lot of space and we would not access it very often. Keeping that data in another table makes our regular staff table more efficient to access. We might also decide that we want to keep some of the employee data confidential; we could store the confidential data in a second table and restrict access to that table to only a few users. We would include the staff id attribute in the second table to link the two tables in a 1:1 relationship.

One-to-Many- this means that a row in table A can be related to 0, 1, or multiple rows in table B and each row in table B is related to one row in table A. This is the more common relationship. A client row in the clients table can be related to 0, 1, or multiple rows in the animals table and each row in the animals table can be related to 0, 1, or multiple rows in the exams table.

Many-to-Many – this means that a row in table A can be related to 0, 1, or multiple rows in table B and each row in table B is related to 0, 1, or multiple rows in table A. This is a relationship that we can think of as a logical relationship but we seldom see it in a database as most dbms do not support this directly. If we think about medications and exams- a medication can be used on multiple exams and an exam can include multiple medications. One of the purposes of the exam details table is to serve as a bridge to implement this many-to-many relationship.

2.2. Referential Integrity rule

We will look at these in terms of the vets set of tables. This is a diagram showing the clients table and the animals table and the relationship between them.



We have referential integrity set between the Clients table and the Animal table. The pk for the clients table is cl_id and since it is a pk it is not nullable. The pk for the animals table is an_id. The animals table also contains an attribute cl_id which is the fk to the client table. We can navigate through the relationship from the animal table back to the client table to find the name of the client for an animal.

Referential Integrity Rule: The database must not contain any non-null unmatched foreign keys. If there is a value in a foreign key attribute, then that value must match a value in the related attribute in the parent table.

In terms of our tables this means that we cannot add a row to the animals table that has a value for the client id that does not match an existing client row. If we want to add a new animal for a new client, we have to add the client first.

The referential integrity rule does allow for null foreign keys. Your business rules might call for a non-null foreign key.

In our animal table the cl_id field is not nullable. We are not allowed to have an animal row without a related client row. But that is not required by the referential integrity rule. It would be possible to set up the animals table with a nullable cl_id attribute. That would let us add animals with no associated client. This would make it difficult to know who to charge for an exam done for the animal- and the vet probably would not like this. What the referential integrity rule says is that an animal row cannot have a value for cl_id that is not a valid cl_id in the clients table- you can't just make up a client id for an animal.

You will not be able to set referential integrity if there already is data present in the tables that violates the rule, or if the attributes involved do not have the same data types.

3. Normalization

Normalization is the formal process of decomposing tables that have design problems into well-structured relations. Good design is sometimes expressed as

- Storing one fact in one place
- Having all of the attributes in a table depend on the primary key, the whole key and nothing but the key.

3.1. Goals of normalization

Normalization helps you design table to

- preserve the integrity of the data
- identify a unique identifier for each record
- reduce redundancy. Redundancy is storing the same data value in more than one table. Redundancy leads to errors.
- minimize the space needed for the tables which also reduces the volume of data that must be transferred to the client.
- store the data you need without having to create dummy data

No design technique will guarantee that you have well designed tables. Having normalized the tables will help.

3.2. Anomalies

Errors in a table that are a result of a user attempts to modify the data are called anomalies.

Update Anomaly — having inconsistent data because the data was updated in one row but not in all rows. For example, in the spreadsheet style table for the animal and client data, the zip code for client 4534 has two different values.

Insertion Anomaly — the inability to add a new data to a table. For example, we cannot add a new client row to the table unless we have data for an animal for that client (or leave these as null values).

Deletion Anomaly — loss of data because we wanted to remove other data values. For example, if we want to remove the data for animal 15342, we could lose all information for client 93 (unless we leave several attributes with null values).

4. Examples of data that is not normalized

For this we will work with variations on the tables in vets database. The examples will focus on mistakes we might have made in the design and problems these might cause.

Design change 1) In the clients table we have a phone number column. We have now found that our clients have more than one phone number and want to give us both their home number and their office number. If the vet has to keep Fluffy overnight and make decisions about Fluffy's treatment, it might be important to be able to reach the client quickly. So we decide to add a new column to the client table and rename the existing column- we now have a column named `home_phone` and a column named `office_phone`. This might have worked in the past when people had one home phone and one job (M_F 9 to 5). This design might still work- but it would be harder to get a list of all of the phone numbers in a single column in the output of a query. And a query that tries to find out which client has a certain phone number has to query two columns. This table is now no longer considered to be good design.

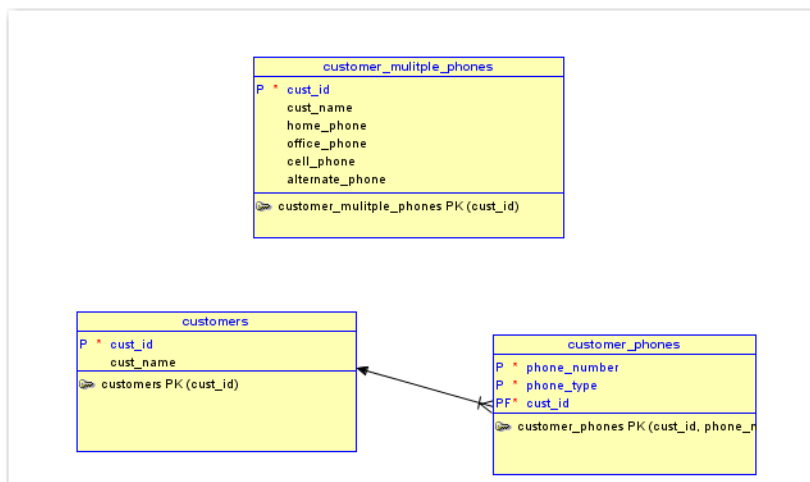
There are several bigger problems: (1) we cannot always distinguish between a home phone and an office phone- what do you do if someone works from home? (2) some people work multiple jobs and have more than one office phone; some people have more than one home phone. (3) how do we classify a cell phone?

Another approach that does not work well is to put all of the phone numbers into a single column separated by commas. Searching is difficult and it is difficult to enforce any format for a legitimate phone number

So if we are going to store multiple phone numbers for a customer, then we need another table- `ClientsPhones`. This table could have two columns: `cl_id` and `cl_phone`. The `cl_id` relates back to the customers table so we can identify a phone number as belonging to a specific client and we can go from a client to their phone numbers. We might decide to have more columns in that table- phone type- (home, office, cell, fax, weekendHome), maybe a column for the best time to call.

The relationship between clients and phone numbers now follows the same pattern as the relationship between clients and animals. A client could have 0, 1, or more phone numbers. This is a one-to-many relationship. We implement it by making the `cl_id` in `ClientsPhones` be a fk to the clients table. We also avoid having a column in the client table which is nullable (The `vt.client.cl_phone` attribute was a nullable column).

This shows two design possibilities- the first has multiple columns for phone numbers and the second design uses two tables and a relationship.



Design change 2) Suppose we had designed the exam headers and the exam details tables so that the exam date was in the exam details table. Perhaps we knew that we would need to write many queries about services that were done on specific dates. By putting the exam date in the exam details table we could avoid a join. But this is a poor decision. The PK of the exam details table is (`ex_id`, `srv_id`) but the exam date depends on the exam id only - it does not depend on the service id (Remember, every piece of data in the table should depend on the

key, the whole key and nothing but the key) From a practical point of view, putting the exam date in the exam detail table would be that we would have to repeat that data value in multiple rows- if an exam involved 5 services we would have to repeat the exam date five times. This is redundant data. There also is a problem that the exam date value might be entered with different values for different rows. This could be due to a data entry error or a change in the system time value while the rows are being entered. If the date needed to be corrected, then the app would need to correct all of the rows consistently. This could be handled at the app level- but it does make this harder.

Design change 3) Suppose we decide that it makes sense to store the staff name in the exam table- so we would have (ex_id, an_id, stf_id, stf_name_last, ex_date). That would save us a join when we wanted to display the details of who was responsible for an exam. But this is a transitive dependency. The pk is ex_id which determines the stf_id which determines the stf_name_last. So in a sense the ex_id determines the stf_name_last but the dependency goes through the stf_id. This design would also cause us to store redundant data and open up the possibility that we spell the person's last name in different ways in different rows. We would also have a problem when a staff person changed their name- we would have to update it in multiple tables.