

Relations

Matt Warner

Contents

1	Anomalies	3
2	Decomposition	5
3	Keys	5
4	Domain	6
5	Order Independence	6
6	Constraints	6
	6.1 Entity contegrity constraint	6
	6.2 Referential Integrity Constraint	7
7	Functional Dependencies	7
	7.1 Armstrong's Axioms	8
	7.2 Revisiting Keys	9
	7.3 How to determine the functional dependencies	9
8	Normalization	10
	8.1 First Normal Form ${}_1NF$	10
	8.2 Second Normal Form ${}_2NF$	11

Overview

- **Anomalies**
 - Insert anomaly
 - Update anomaly
 - Deletion Anomaly
- **Constraints**
 - referential integrity constraints
 - Entity Integrity Constraint
- **Functional Dependency**
 - Armstrongs axioms
- **Keys**
 - Primary key
 - candidate key
 - Super key
 - foreign key
- **Key attributes**
- **non-key attributes**
- **Domain Order Independence**
- **Decomposition**
- **Normalization**
 - $1NF$
 - $2NF$
 - $3NF$

1 Anomalies

If our database is a single relation with schema ***SP*** (SuppName, SuppAddr, Item, Price).

With out instance data:

SuppName	SuppAddr	Item	Price
John	10 Main	Apple	\$2.00
John	10 Main	Orange	\$2.50
Jane	20 State	Grape	\$1.25
Jane	20 State	Apple	\$2.25
Frank	30 Elm	Apple	\$6.00

There are some common things that we might want to do that would cause issues. We refer to these as **anomalies**, and there are split into three categories.

Insertion Anomaly

Let's say we want to add a new vendor, "Sally", and store her address, "40 Pine", but she is not selling anything yet. Can this be inserted into the relation SP?

Sally	40 Pine	???	???
-------	---------	-----	-----

The answer to this question is **NO**. The *primary key* is (SuppName, Item), but we only have SuppName. The *entity integrity constraint* is violated if we try to insert the data as a tuple in this relation. It cannot fit. We call this an *insertion anomaly*.

Deletion Anomaly

This time, let's say that Frank no longer sells Mango. We want to take that out of the database so nobody can order a mango that is not available. Our new tuple would look like this:

Frank	30 Elm	???	???
-------	--------	-----	-----

Can this tuple remain in the relation with the Mango information removed?

No, it cannot. The *primary key* is (SuppName, Item), and the Item is going away. The *entity integrity constraint* is violated if we remove the data from the tuple in this relation. We can either keep the whole tuple, advertising fake mango, or delete the whole tuple and lose the information on Frank, which doesn't exist in any other tuples. We call this a *deletion anomaly*.

Update Anomaly

Next, let's say that John is moving to a different address. We would want to change it once for every item John is selling.

John	10 Main	Apple	\$2.00
John	10 Main	Apple	\$2.00

This isn't a big deal with only two items, but as John's list of supplied items grows, so does the amount of database work that needs to be done every time he moves. If any of the SuppAddr values for John don't agree, then it may not be clear which is the right address for John. This is an *update anomaly*.

In summary, we have:

- Insertion anomalies
 - When a piece of data cannot be inserted because it violates some *constraint* of the relation.
 - Usually is the *entity integrity constraint* being violated, but not always.
- Deletion anomalies
 - When deleting some piece of data, a *deletion anomaly* is when more data is lost than intended.
 - Usually this is caused when the data removed is part of the *primary key*, which would cause a violation of the entity integrity constraint.
- Update anomalies
 - When updating a single value requires changes to multiple tuples, this is an *update anomaly*.
 - This is caused by unnecessary redundancies in the data.
 - These cause inefficiency, and potential inconsistencies.

2 Decomposition

Here, we represent the original data in two relations, rather than the one.

$SP(\underline{\text{SuppName}}, \underline{\text{Item}}, \text{Price})$

SuppName	Item	Price
John	Apple	\$2.00
John	Orange	\$2.50
Jane	Grape	\$ 1.25
Jane	Apple	\$2.25
Frank	Mango	\$6.00

$SP(\underline{\text{SuppName}}, \text{SuppAddr})$

SuppName	SuppAddr
John	10 Main
Jane	20 State
Frank	30 Elm

Now, we can make changes to insert Sally and his address without needing to care whether or not he is selling anything. We can update johns address in one spot instead of multiple, and we can delete Jane from the top table when she is no longer selling anything.

3 Keys

Keys are one of the basic requirements of a relational database model. It is widely used to identify the tuples uniquely in the table. We also use keys to set up relations amongst various columns and tables of a relational database.

Types of Keys

- Super key
- Candidate Key
- Primary key
- Foreign key

Super key

A super key is an attribute or set of attributes whose values can uniquely identify any tuple.

Every relation has at least one - the set of all attributes in the relation (since duplicate tuples are considered to be the same tuple)

There can potentially be many available, some more useful than others.

Candidate Key

This is a minimal super key. It is the minimal set of attributes that can uniquely identify a tuple. For example, `Student_ID` in `Student` relation.

StudentID	RollNo	Name	MobileNo	EmailID
A1	1	Matt	9120	a@gmail.com
A2	2	John	8732	b@gmail.com
A3	3	Luke	8344	c@gmail.com

In this table: `StudentID` \rightarrow `RollNo`, `Name`, `MobileNo`, `EmailID`. Therefore, it is a candidate key.

Primary Key

The **primary key** for a relation is chosen by the database designer from among the relation's candidate keys. It becomes the "official" key that is used to reference tuples within the relation. There can be only one.

Once a primary key is chosen, each of the attributes in the relation will be either **prime** or **non-prime** with respect to the relation.

- A **prime** attribute is one of the attributes that can be found in any of the candidate keys.
- A **non-prime** attribute is one of the attributes *not found* in any of the candidate keys.

Once a primary key is chosen for it, the **schema** of a relation is written with the primary key's attributes underlined:

$$\text{Relation_Name}(\underline{A_1}, A_2, A_3, \dots, A_n)$$

Foreign Keys

A **foreign key** is a tool used to link relations within a database. Since every relation has a primary key that uniquely identifies each tuple, the values of those key attributes can be used from another relation to reference individual tuples.

The relation whose primary key is being used is the **home relation**.

4 Domain

The **domain** of an *attribute* is the set of all possible values it may hold.

The **domain** of a *set of attributes* is the set of all possible combinations of values for the attributes in the set.

5 Order Independence

In relations, the order things appear doesn't matter. There are ways to force them to sort later when we're working with SQL, but the relation itself has no order for either rows or attributes.

It doesn't matter what order the attributes appear in, if two relational schemas have the same name, the same attributes, and the same primary key, then they are equivalent.

So, all of these are equivalent:

$$R(\underline{A}, B, C, D)$$

$$R(D, C, B, \underline{A})$$

$$R(\underline{A}, D, B, C)$$

Tuples are stored unordered. If you need to have them appear in some order later, you will be able to sort based on the values inside of them using SQL.

6 Constraints

6.1 Entity contegrity constraint

The entity integrity constraint applies to all relations. It states that no tuple may exist within a relation that has null value for any of attributes that make up the primary key.

This is a consequence of the primary key being a candidate key, which is minimal and cannot do its job with less data.

6.2 Referential Integrity Constraint

The referential integrity constraint applies to all foreign keys. It constrains the values of foreign keys in relations to values that actually exist as primary keys for tuples within the home relation.

If the foreign key is otherwise allowed to be NULL, then that is also an acceptable value.

7 Functional Dependencies

A *functional dependency* is a statement about which attributes can be inferred from other attributes. If we take X and Y as *sets* of attributes, we can write:

$$X \rightarrow Y$$

If, whenever unique values for **all** of the attributes in X are known, unique values for **each** of the attributes of Y are guaranteed to be possible to look up or to infer using those values. This is read either as:

X functionally determines Y , *or*

Y is functionally dependent upon X

Definition

For a functional dependency to exist between two attributes, x and y , that is:

$$x \rightarrow y$$

Then the following must be true:

$$\begin{array}{ll} \text{if} & t_1.x = t_2.x \\ \text{then} & t_1.y = t_2.y \end{array}$$

They are statements about the operational data. Later on, we will see how to read them off of ER diagrams, though they may come from elsewhere as well.

Real-life Examples

$$ZID \rightarrow \text{StudentFirstName, StudentLastName, Birthday}$$

If I identify a student using their ZID, that student has *one* first name, last name, and birthday.

$$\text{StudentFirstName} \rightarrow ZID$$

The first name is not enough to determine a single ZID, as there are multiple students with the same first name.

$$ZID, \text{CourseID}, \text{Semester} \rightarrow \text{Grade}$$

If i know which student, which course, and which semester, I can find a single grade.

Keep in mind that *functional dependencies* are constraints present within the operational data your database models. They don't necessarily describe how things work in the real world, but they do have to accurately describe any data you will store in your database.

Additionally, *functional dependencies* **must** hold for all possible data values. Attempts to add data that does not obey the functional dependencies will result in anomalies.

Furthermore, functional dependencies **can** be enforced during insertion if the database is set up properly.

7.1 Armstrong's Axioms

Armstrong's Axioms are a set of rules for operations that are permissible when manipulating *functional dependencies*.

Primary Rules:

Axiom of reflexivity:

If $Y \subseteq X$ then $X \rightarrow Y$

Axiom of augmentation:

If $X \rightarrow Y$, then $XZ \rightarrow YZ$ for any Z

Axiom of transitivity:

If $X \rightarrow Y$ and $Y \rightarrow Z$, then $X \rightarrow Z$

Secondary Rules

Decomposition:

If $X \rightarrow YZ$ then $X \rightarrow Y$ and $X \rightarrow Z$

Composition:

If $X \rightarrow Y$ and $A \rightarrow B$ then $XA \rightarrow YB$

Union (Notation):

If $X \rightarrow Y$ and $Y \rightarrow Z$ then $X \rightarrow YZ$

Pseudo-transitivity:

If $X \rightarrow Y$ and $YZ \rightarrow W$ then $XZ \rightarrow W$

Self-determination:

$I \rightarrow I$ for any I

Example: Relation with FDs

Lets say we have the following relation:

EmpProj(EmpID,Project, Supv, Dept, Case)

EmpID	Project	Supv	Dept	Case
e1	p1	s1	d1	c1
e2	p2	s2	d2	c2
e1	p3	s1	d1	c3
e3	p3	s1	d1	c3

Our functional dependencies are:

- EmpID, Project \rightarrow Supv, Dept, Case
- EmpID \rightarrow Supv, Dept
- Supv \rightarrow Dept

As written, there are some anomalies present. We will use *normalization* to move toward a better design.

7.2 Revisiting Keys

When we talked about *keys*, we talked about how their purpose is to uniquely identify a tuple within a relation. Another way of stating this, now that we know about *functional dependencies*, is **the attributes of a superkey must functionally determine all of the attributes of the relation**.

Candidate keys and *primary keys* **are** *super keys*, so this is true of them as well, and they also satisfy additional requirements. As an example, say we have the relation $R(\underline{a}, b, c, d, e, f)$

$$a \rightarrow a, b, c, d, e, f$$

But, since it is always the case that $a \rightarrow a$ because of the self-determination axiom, we usually omit the left hand side from the righthand side. So we would usually write this instead:

$$a \rightarrow b, c, d, e, f$$

7.3 How to determine the functional dependencies

Lets say we have the following attributes, and want to identify all the functional dependencies:

- shipmentID
- shipmentDate
- origin
- destination
- shipID
- shipName
- CaptainID
- capatinName
- ItemID
- Description
- Weight
- quantity

In terms of functional dependencies we have:

shipmentID \rightarrow shipmentDate, origin, Destination, shipID, ShipName, CaptainID, CaptainName

shipID \rightarrow shipName, captainID, CaptainName

captainID \rightarrow CaptainName

itemID \rightarrow description, weight

itemID, ShipmentID \rightarrow quantity

The first determinate is the shipmentID. If we know what the shipmentID is, then we also know its ShipmentDate, Origin, destination, shipID, ShipName, CaptainID and captainName. Regarding the captain, we are making an assumption that a ship has only one captain. We cannot include itemID, description, weight, or quantity because these all refer to the individual items, and a shipment ID cannot determine an item, because there can be many items.

We can ignore shipmentDate. Usually, you shouldn't think twice about skipping dates as the date alone cant really determine anything. The next attributes are origin and destination, we are also going to skip over these since they cant really determine anything either.

Our next determinate is shipID. This one is pretty obvious... a shipID determines shipName, captainID (again assuming there is one captain per ship), and CaptainName.

Next up is captainID, if we know the captainsID, we know the captains name, since each each Captain is linked to one captainID.

ItemID gives us its description and weight, and itemID + the shipmentID gives us the quantity of the item.

Our schema for this can be seen as such:

R(shipmentID, shipmentDate, origin, destination, shipID, shipName, CaptainID, CaptainName, ItemID, Description, Weight, Quantity)

At this point we should probably select our primary key. Since we select our primary key from our list of candidate keys, we need to first assess all our candidate keys. Recall that a candidate key is the minimum set of attributes necessary to uniquely identify a tuple.

The first attribute we should look at is obviously ShipmentID, since it can determine the most amount of attributes within the tuple. Since it does not determine all of our attributes, it by itself is not a candidate key, so we need another attribute to pair it with. We need an attribute or a set of attributes that can functional determine itemID, description, weight, and quantity. itemID determines description and weight and when it is paired with shipmentID, it also determines quantity. Therefore, our first candidate key is: {ShipmentID ItemID}.

Now, our prime attributes are ShipmentID and ItemID and our non-prime attributes are ShipmentData, origin, destination, shipID, ShipName, CaptainID, CaptainName, Description, Weight and quantity.

8 Normalization

8.1 First Normal Form $1NF$

The requirement for a relation to be in $1NF$ is that all of the values must be **atomic**.

What this usually looks like is a table with multiple values in a single cell. A non- $1NF$ relation would not even technically count as a relation. This table has a cell that is non-atomic,

X	Y	Z
x1	y1	z1
		z2
		z3
x2	y2	z4
x3	y2	z5

It looks X *would* have been the primary key, but it's not doing its job of uniquely determining Z , which is showing as a *repeating group* so X can't be a key.

The notation for this “pseudo-relation”. like the one above would be to use **inner parenthesis** on the repeating group, i.e.

$R(\underline{X}, Y, (Z))$

This is not $1NF$, and has functional dependencies:

$X \rightarrow Y$

$X, Z \rightarrow Z$ **but** $X \nrightarrow Z$

To move this pseudo-relation into an actual relation that doesn't violate $1NF$, we need to choose a *real* primary key that meets the requirements. We do that using the FDs. In this case, (X,Z) works.

Changing the primary key yields: - $R(\underline{X}, Y, \underline{Z})$

X	Y	Z
x1	y1	z1
x1	y1	z2
x1	y1	z3
x2	y2	z4
x3	y2	z5

Now everything is atomic, and we are in $1NF$. Notice that this did introduce a new update anomaly, but the other normal forms will take care of it. It is more important to get into $1NF$ for now. As another example, consider the following unnormalized pseudo-relation:

$R(\underline{A}, B, C, (d_1, d_2, d_3), E, F)$

$A \rightarrow B, C, E, F$

$A, d_1 \rightarrow d_2, d_3$

Notice that (d_1, d_2, d_3) is a repeating group. A is not enough to form a primary key, it needs d_1 to be able to determine d_2 and d_3 . So, the actual primary key in this case should be (A, d_1) , making the $1NF$ relation.

$R_{1NF}(\underline{A}, B, C, d_1, d_2, d_3, E, F)$

8.2 Second Normal Form $2NF$

Second Normal Form ($2NF$) has to do with the concept of *full dependence*.

Given two sets of attributes, X and Y , we can say that Y is *fully dependent* on X , if (and only if)

- $X \rightarrow Y$
- No subset of X determines Y

A relation is in $2NF$ if:

- It already meets the requirements of $1NF$,
- All *non-prime* attributes of the relation are *fully dependent* upon the **entire** primary key.

What breaks $2NF$ is when attributes are dependent upon only **part** of the primary key.

To fix $2NF$ violations once we're in $1NF$, *decomposition* is the solution.

Going back to our early example: **EmpProj**(EmpID, Project, Supv, Dept, Case)

EmpID	Project	Supv	Dept	Case
e1	p1	s1	d1	c1
e2	p2	s2	d2	c2
e1	p3	s1	d1	c3
e3	p3	s1	d1	c3

EmpID, Project \rightarrow Supv, Dept, Case

EmpID \rightarrow Supv, Dept

Supv \rightarrow Dept

A quick glance confirms all values are atomic, so $1NF$ is confirmed.

There is a $2NF$ violation caused by (EmpID \rightarrow Supv, Dept) because the primary key is (EmpID, Project), but only EmpID is on the LHS.

Observing the instance Data, you should easily see that the attributes of the RHS cause update anomalies in this table. We also can't insert a new employee with no project (insertion anomaly). These are symptoms of the $2NF$ violation.