

→ Database Design

- Informal design Guidelines
- Functional dependencies
- Normal form based on Primary keys
- General definition of 2nd & 3rd NF.

→ Informal design Guidelines for Relational Schema

- Schematics of the attributes
- Reducing the redundant values in tuples
- Reducing null values in tuples
- Disallowing spurious tuples

1. Schematics of the attributes

Guideline 1: Design a relational schema so that it is easy to explain its meaning. Do not combine attributes from multiple entity types and relationship entity types into a single entity.

If a relational schema corresponds to one entity type or one relationship type the meaning tends to be clear. otherwise it tends to be a mixture of multiple entities and relationship and hence semantically not correct.

2. Reducing the redundant values in tables.

one of the goal of the Schema design is to minimize the storage space that the base relation occupies.

Serious Problem with redundancy is the problem of Update Anomalies.

This is classified in to

- Insertion Anomalies
- Deletion Anomalies
- Modification —

for ex: Suppose a manufacturing Company stores the employee information in a table as follows:

emp-id	emp-name	emp-address	emp-Department
101	Rick	Delhi	D 001
101	Rick	Delhi	D 002
123	anoggie	Agra.	D 890
166	Glenn	Chennai	D 900
166	Glenn	Chennai	D 004

- **Update Anomaly :** In the above table we have 2 rows for employee Rick as he belongs to two department of Company. If we want to update the address of Rick we have to update the same in two rows otherwise the data will become inconsistent.

Some how the correct address gets updated in one department but not in other as per database Rick could be having two different address which is not correct which leads to data inconsistency.

b) Insert Anomaly : Suppose a new employee joins the Company who is under the training and ~~is~~ currently not assigned to any department, then we could not be able to insert data into the table if employee department field does not allow null values.

c) Delete Anomaly : Suppose if at a point of time the company closes the department D890 then deleting rows that are having employee department D890 would also delete the information of employee Maggie. Since she is assigned only to this department.

Guideline 2 :- Design the Relational Schema so that no insertion, deletion, modification anomalies occur in the relation.

If any anomalies are present note them clearly so that the program that update database will operate correctly.

3. Reducing the null values in tuples

In some schema designs we may group many attributes together in to a base relation.

Null values can have multiple interpretation. Such as

- Attribute does not apply to this tuple.
- Attribute Value for this is unknown.
- The value is known but absent it has not been recorded yet.

Guideline 3: As far as possible avoid placing attributes in base relations whose values may be null. If nulls are unavoidable make sure that they apply in exceptional cases only and do not apply to a majority of tuples in the relation.

4. Disallowing Spurious tuples:

Spurious tuples may occur after join operation. Good example is Cartesian product.

Guideline No. 4: Design relational Schema so that they can be joined with equality condition on attributes, that are either on primary keys and foreign keys

Functional dependencies :-

Definition :-

A functional dependency is defined between 2 sets of attributes from database.

- X functionally determines $-Y$ in a relation schema, if and only if whenever 2 tuples $R(X, Y)$ agree on their X values, they must necessarily agree on $-Y$ value.

Given a table T containing at least 2 attributes designated by X and Y , we say that $X \rightarrow Y$,

X functionally determines Y

OR

Y is functionally dependent on X

Eg:

<u>Employee</u>	
SSN	fname

$SSN \leftarrow fname$

works-on

SSN	Pno	Hrs
-----	-----	-----

$\{ SSN, Pno \} \leftarrow hours$

Inference rules for functional dependencies

→ Inference rules are also known as Armstrong's Axiom. are published by Armstrong.

→ Armstrong Properties are

1. Reflexivity Property.

$x \rightarrow y$ if y is subset of x

2. Augmentation Property.

If $x \rightarrow y$ is true, then

$xz \rightarrow yz$ is also true.

3. Transitivity Property.

If $x \rightarrow y$ and $y \rightarrow z$ then

$x \rightarrow z$ is true.

4. Union Property.

If $x \rightarrow y$ and $x \rightarrow z$ are true

then $x \rightarrow yz$ is also true.

This property indicates that right hand side of F.D contains many attributes then F.D exists for each of them

5. Decomposition Property

If $x \rightarrow y$ is implied and z is subset of y then $x \rightarrow z$ is implied.
This property is reverse of Union property

6. Pseudo transitivity Property

$x \rightarrow y$ and $wy \rightarrow z$ are given

then $xw \rightarrow z$ is true.

Normalization:

Normalization is the process of removing redundant data from the tables in order to improve storage efficiency, data integrity, and scalability.

There are two goals of the normalization process

- Eliminating redundant data.
- Ensuring data dependencies.

why do we need normalization :

Normalisation is the aim of well designed Relational DBs Management System. It is step by step set of rules by which data is put in its simplest form. we normalize the RDBMS because of the following reasons:

- Minimize data redundancy.
- To make DB structure flexible i.e. it should be possible to add new data values and rows without reorganizing the database structure.
- Data should be consistent throughout the DB i.e. it should not suffer from following anomalies
 - Insert Anomaly
 - Delete Anomaly
 - Update Anomaly
- Complex ~~querier~~ queries required by the user should be easy to handle.

- on decomposition of a relation up to smaller relations with fewer attributes on normalization the resulting relations whenever joined must result in the same relation without any extra rows. This is known as lossless join decomposition.

Advantages of Normalization

- more efficient data structure.
- Avoid redundant fields or columns.
- more flexible data structure.
- Better Understanding of data
-

Different forms of Normalization:

- first normal form
- Second — —
- Third — —
- Boyce - codd Normal form
- fourth normal form,
- fifth — —

First normal form :- Its purpose is to eliminate repeating groups of attributes in an entity.

1NF disallow Composite attributes, multivalued attributes, nested relations.

Attributes whose values for an individual tuple are atomic.

Consider the following table

Department

a)	Dname	Dnumber	Dmgr-ssn	Dlocation
----	-------	---------	----------	-----------

Relational schema not in 1NF

b)	Dname	Dnumber	Dmgr-ssn	Dlocation
	Research	5	333	Bangalore, Mangalore, Calcutta.
	Administrative	3	555	Delhi
	Headquarters	6	777	Bombay

c)	Dname	Dnumber	Dmgr-ssn	Dlocation
	Research	5	333	Bangalore
	Research	5	333	Mangalore
	Research	5	333	Calcutta
	Administrative	3	555	Delhi
	Headquarters	6	777	Bombay

Second normal form : 2NF uses the concept of functional dependencies and primary key.

A table is said to be in 2NF if both the following conditions hold :

- Table is in 1NF
- All non key attributes are fully functional dependent on the primary key.

Consider the following example :

Customer Id.	Store Id.	Purchase location
1	1	Los Angeles
1	3	San Francisco
2	1	Los Angeles
3	2	New York
4	3	San Francisco

This table has a composite primary key, [Customer Id, Store Id]. The non-key attribute is [Purchase location]. In this case [Purchase location] ~~is~~ only depends on [Store Id], which is only a part of primary key. Therefore the table is not in 2NF.

- we break the table into two tables and.
Now we have the following

Table - purchase

Customer id	Store id
1	1
1	3
2	1
3	2
4	3

Table - store

Store id	Purchase location
1	Los Angeles
2	New York
3	San Francisco

- Know the table is in 2NF. The [Purchase location] is fully dependent on the Primary key of the table which is [Store id].

3rd Normal Form:

A Relational schema is in third NF if it satisfies the following conditions :

- It is in Second Normal Form
- There is no transitive functional dependency,

Transitive functional dependency :

A is functionally dependent on B and.
B is functionally dependent on C. ~~Intra~~
In this case, -C- is transitively dependent
on -A- via B.

~~3rd~~ Example :

Consider the following example :

Book-id	Gener-id	Gener-type	Price
1	1	Gardening	25.99
2	2	Sports	14.99
3	1	Gardening	10.00
4	3	Travel	12.99
5	2	Sports	17.99

In this table

[Book-id] determines [Gener-id]
and [Gener-ID] determines [Gener-type].

Therefore [Book-ID] determines [Gener-type]
via [Gener-id].

we have transitive functional dependency and
this structure does not satisfy third NF.

To bring this table to third NF, we split
table into two as follows :

Table-book			Table-Gener	
Book-id	Gener-id	Price	Gener-id	Gener-type
1	1	25.99	1	Gardening
1	2	14.99	2	Sports
3	1	10.00	3	Travel
4	3	12.99		
5	2	17.99		

- . Now all non key attributes are fully functional dependent only on the primary key.
- . In [Table - Book], both [Genre Id] and [price] are only dependent on [Book ID].
- . In [Table - Genre], [Genre type] is only dependent on [Genre Id].

Boyce-Codd Normal Form (BCNF) was proposed as a simpler form of 3NF, but it was found to be stricter than 3NF. That is, every relation in BCNF is also in 3NF, however a relation in 3NF is not necessarily in BCNF.

Definition: A relational schema R is in BCNF if whenever a nontrivial function dependency $X \rightarrow A$ holds in R, then X is a superkey of R.

The formal definition of BCNF slightly differs from the definition of 3NF. The only difference b/w the two is that in 3NF a relational schema R is in 3NF if every nonprime attribute of R is non transitively dependent on every key of R, but this condition does not exist in BCNF.

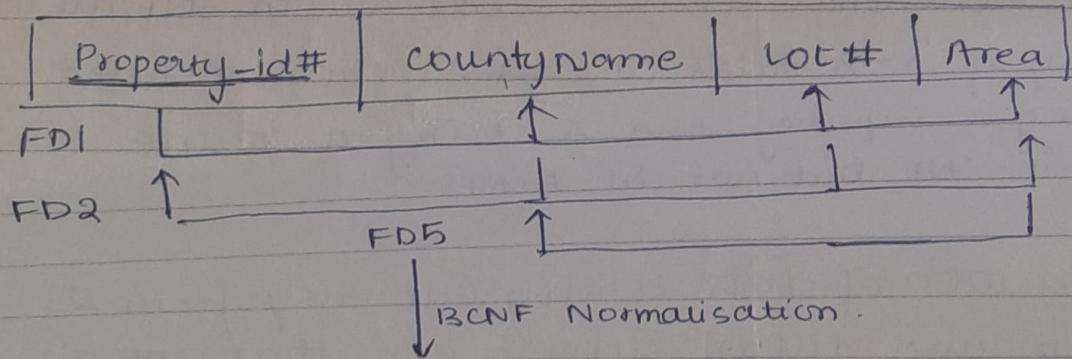
In practice, most relational schemas that are in 3NF are also in BCNF. Only if $X \rightarrow A$ holds in ~~not~~ relational schema R with X not being a superkey and A being a prime attribute will be in 3NF but not in BCNF.

Suppose we have 1000 of LOTS coming but only from 2 countries: DeKalb and Fulton. Suppose also that lot sizes in DeKalb country are only 0.5, 0.6, 0.7, 0.8, 0.9, 1.0 acres whereas lot sizes in Fulton country restricted to 1.1, 1.2, ..., 20 acres. In such situation we could have the additional functional dependency FDS : Area \rightarrow country_Name. If we add this to the other dependencies, the relational schema LOTSIA still is in 3NF because country_Name is a prime attribute.

The area of a lot that determines the county, a specified key FD5 can be represented by 16 tuples in separate relation R (Area, County-name), since there are only 16 possible Area values. This representation reduces the redundancy of repeating the same information in the thousand of LOTSIA tuples. BCNF is stronger normal form that would disallow LOTSIA and suggest need of decomposing it.

FD5 violates BCNF in LOTSIA because Area is not the superkey of LOTSIA. Note that, FD5 satisfies 3NF in LOTSIA because county-name is a prime attribute (condition b), but this condition does not exist in definition of BCNF. We can decompose LOTSIA into 2 BCNF relation LOTSIA_X and LOTSIA_Y. This decomposition loses the functional dependency FD₂ because its attributes no longer coexist in the same relation after decomposition.

LOTSIA

LOTSIA_X

<u>Property-id#</u>	<u>Area</u>	<u>LOT#</u>

LOTSIA_Y

<u>Area</u>	<u>County-Name</u>

Multivalued Dependencies and Fourth Normal Form :-

In database design there are certain cases where relation have constraints that cannot be specified as functional dependency.

Multivalued dependencies are a consequence of INF which disallowe an attribute in a tuple to have a set of values.

If we have two or more multivalued independent attributes in a tuple the same relation schema, we get into a problem of having to repeat every value of one of the attributes with each value of the other attributes to keep the relation state consistent and to maintain the independence among the attributes involved. This constraint is specified by multivalued dependency.

EMP

Ename	Pname	Sname
Smith	X	John
Smith	Y	Anna
Smith	X	Anna
Smith	Y	John

$\text{Ename} \rightarrow\!\! \rightarrow \text{Pname}$

$\text{Ename} \rightarrow\!\! \rightarrow \text{Sname}$

* Here Ename works on Pname and has dependent Sname

* An employee may work on several projects and may have several dependents

* Emp's project and Emp's dependent are independent of one another

To keep the relation state consistent we must have separate tuples to represent every combination of employee's project and employee's dependent.

This constraint is termed as Multivalued dependency.

Informally, when two or more independent 1:N relationships A:B and A:C are mixed in same relation a MVD arises.

Formal Definition of MVD :-

A MVD $X \rightarrow\!\!> Y$ specified on relation schema R where X and Y are both subset of R, specifies the following constraints on any relation state σ of R : If two tuples t_1 and t_2 exists in σ such that $t_1[X] = t_2[X]$ then two tuples t_3 and t_4 should also exist in σ with the following properties —

$$t_3[X] = t_4[X] = t_1[X] = t_2[X]$$

$$t_3[Y] = t_1[Y] \text{ and } t_4[Y] = t_2[Y]$$

$$t_3[Z] = t_2[Z] \text{ and } t_4[Z] = t_1[Z]$$

where $Z = R - (X \cup Y)$.

* If MVD $X \rightarrow\!\!\! \rightarrow Y$ is said to be a trivial MVD if it satisfies the below two conditions -

- ① Y is subset of X and set is closed
- ② $X \cup Y = R$

* A MVD $X \rightarrow\!\!\! \rightarrow Y$ is said to be a non-trivial MVD if it does not satisfies the above conditions.

Fourth Normal Form

Defⁿ

A relation schema is in 4NF with respect to a set of dependencies F (Functional and multivalued) if for each non-trivial multivalued dependency $X \rightarrow\!\!\! \rightarrow Y$ in F , X is a superkey for R .

Converting the relation EMP in 4NF

Since relation EMP is a all-key relation hence it is by default in ~~4NF~~ BCNF

* EMP relation is not in 4NF because of

(i) Ename $\rightarrow\!\!\! \rightarrow$ Pname

(ii) Ename $\rightarrow\!\!\! \rightarrow$ Sname

(iii) Ename is not a superkey.

20

The relation EMP is decomposed into two relations in order to make it in 4NF based on the above MVDs (i) and (ii)

<u>Emp-Project</u>		<u>Emp-dependent</u>	
Ename	Prname	Ename	Dname
Smith	X	Enam	John
Smith	Y	Smith	Anna

This decomposition helps

- (i) eliminate increased storage need
- (ii) eliminate the update anomaly.

Joint dependency

The formal definition for Joint dependency is denoted by $JD(R_1, R_2 \dots R_n)$ specified on relation schema R .

R specifies a constraint on the states τ of R . Every legal state τ of R should have a non-additive join decomposition into $R_1, R_2 \dots R_n$.

These dependency is very peculiar semantic constraint i.e difficult to detect in practice.

e.g. consider a schema R having three attributes supp, parts + project.

		R		parts	project
		s ₁	s ₂	P ₁	P ₂
		r ₁	r ₂	r ₁	r ₂
s ₁					
s ₂					

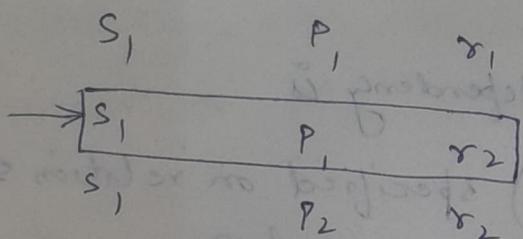
Now decompose R into three other schema such that $R \rightarrow (R_1, R_2, R_3)$

be	R ₁		R ₂		R ₃	
	supp	parts	supp	proj	parts	proj
s ₁	P ₁		s ₁	r ₁	P ₁	r ₁
s ₁	P ₂		s ₁	r ₂	P ₂	r ₂
s ₂	P ₁		s ₂	r ₁	P ₁	r ₂
s ₂	P ₁		s ₂	r ₂		

Now let us merge the above decomposed relations $R_1, R_2 + R_3$ to form a schema R . Here we are merging for supplier S_1 only.

supp

parts proj



After merging the decomposed schema R, R_1, R_2 , we get the new schema above for supp S_1 . In this schema we get a new tuple which was not there in original schema R . Such tuple is called additive Join.

If we don't get any new tuple after merging the decomposed schema then it is called Non-additive loss-less join.

For a relation schema R to be join dependent

A relation schema R is join dependent if it is non-additive loss-less join.

8th *

(a, b, c) is a tuple in R with some attributes missing

a	b	c	d	e	f	g	h	i	j	k	l	m	n	o	p	q	r	s	t	u	v	w	x	y	z
6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	
1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	
2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	
3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	

5th Normal Form or Project join normal form (PJNF) (23)

A relation schema R is said to be in 5th NF if

it satisfies following conditions

1) R must be in 4NF

2) if Joint Dependency (JD), ^{not} exists

R is in 5NF

else if JD exists

check for trivial JD if yes

R is in 5NF

else if $\forall R_i$ is superkey

$\parallel R_i \in R_1, R_2, R_3, \dots, R_n$

R is in 5NF

else

R is not in 5NF

So, if above two conditions are satisfied then we say

the relation schema R is in 5NF.