

## Unit IV

### Normalization of Database

Database Normalizations is a technique of organizing the data in the database. Normalization is a systematic approach of decomposing tables to eliminate data redundancy and undesirable characteristics like Insertion, Update and Deletion Anomalies. It is a multi-step process that puts data into tabular form by removing duplicated data from the relation tables.

Normalization is used for mainly two purpose,

Eliminating redundant(useless) data.

Ensuring data dependencies make sense i.e. data is logically stored.

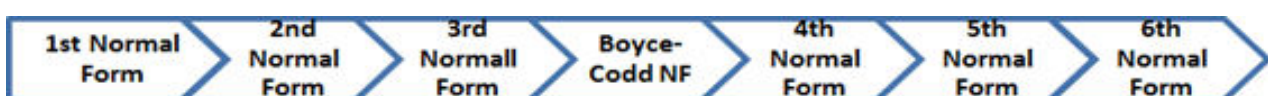
### Problem Without Normalization

Without Normalization, it becomes difficult to handle and update the database, without facing data loss. Insertion, Updating and Deletion Anomalies are very frequent if Database is not Normalized. To understand these anomalies let us take an example of **Student** table.

S_id	S_Name	S_Address	Subject_opted
401	Adam	Noida	Bio
402	Alex	Panipat	Maths
403	Stuart	Jammu	Maths
404	Adam	Noida	Physics

- **Updating Anomaly:** To update address of a student who occurs twice or more than twice in a table, we will have to update **Address** column in all the rows, else data will become inconsistent.
- **Insertion Anomaly:** Suppose for a new admission, we have a Student id(S\_id), name and address of a student but if student has not opted for any subjects yet then we have to insert **NULL** there, leading to Insertion Anomaly.
- **Deletion Anomaly:** If (S\_id) 401 has only one subject and temporarily he drops it, when we delete that row, entire student record will be deleted along with it.

Theory of Data Normalization in Sql is still being developed further. For example, there are discussions even on 6th Normal Form. But in most practical applications normalization achieves its best in 3rd Normal Form. The evolution of Normalization theories is illustrated below-



### Functional Dependencies

A functional dependency is a relationship between two attributes. Typically, between the PK and other non-key attributes with in the table. For any relation R, attribute Y is functionally dependent on attribute X (usually the PK), if for every valid instance of X, that value of X uniquely determines the value of Y.

$$X \longrightarrow Y$$

The left-hand side of the FD is called the determinant, and the right-hand side is the dependent.

Examples:

$$\text{SIN} \longrightarrow \text{Name, Address, Birthdate}$$

SIN determines names and address and birthdays. Given SIN, we can determine any of the other attributes within the table.

$$\text{Sin, Course} \longrightarrow \text{Date-Completed}$$

Sin and Course determine date completed. This must also work for a composite PK.

$$\text{ISBN} \longrightarrow \text{Title}$$

ISBN determines title.

Various Types of Functional Dependencies are –

- Single Valued Functional Dependency
- Fully Functional Dependency
- Partial Functional Dependency
- Transitive Functional Dependency
- Trivial Functional Dependency
- Non-Trivial Functional Dependency
  - Complete Non-Trivial Functional Dependency
  - Semi Non-Trivial Functional Dependency

### Single Valued Functional Dependency –

Database is a collection of related information in which one information depends on another information. The information is either single-valued or multi-valued. For example, the name of the person or his / her date of birth are single valued facts. But the qualification of a person is a multivalued fact.

A simple example of single value functional dependency is when A is the primary key of an entity (eg. SID) and B is some single valued attribute of the entity (eg. Sname). Then,  $A \rightarrow B$  must always hold.

CID	SID	Sname
C1	S1	A
C1	S2	A
C2	S1	A
C3	S1	A

SID $\rightarrow$ Sname	Sname $\rightarrow$ SID
S1 A	A S1
S1 A	A S2
S1 A	

For every SID, there should be unique name ( $X \rightarrow Y$ )

Definition: Let R be the relational schema and X, Y be the set of attributes over R. t1, t2 be the any tuples of R.  $X \rightarrow Y$  exists in relation R only if  $t1.X = t2.X$  then  $t1.Y = t2.Y$

If condition fails – then dependency is not there.

### Fully Functional Dependency

In a relation R, an attribute Q is said to be fully functional dependent on attribute P, if it is functionally dependent on P and not functionally dependent on any proper subset of P. The dependency  $P \rightarrow Q$  is left reduced, there being no extraneous attributes in the left-hand side of the dependency.

If  $AD \rightarrow C$ , is fully functional dependency, then we cannot remove A or D. I.e. C is fully functional dependent on AD. If we are able to remove A or D, then it is not fully functional dependency.

Another Example, Consider the following Company Relational Schema,

**EMPLOYEE**

ENAME	SSN(P.K)	BDATE	ADDRESS	DNUMBER
-------	----------	-------	---------	---------

**DEPARTMENT**

DNAME	<u>DNUMBER (P.K)</u>	DMGRSSN(F.K)
-------	----------------------	--------------

**DEPT\_LOCATIONS**

DNUMBER (P.k)	DLOCATION (P.K)
---------------	-----------------

**PROJECT**

PNAME	PNUMBER	PLOCATION	DNUM
-------	---------	-----------	------

**WORKS\_ON**

SSN(P.K)	PNUMBER(P.K)	HOURS
----------	--------------	-------

{SSN, PNUMBER}  $\rightarrow$  HOURS are a full FD since neither  $SSN \rightarrow HOURS$  nor  $PNUMBER \rightarrow HOURS$  hold

{SSN, PNUMBER}  $\rightarrow$  ENAME is not a full FD (it is called a partial dependency) since  $SSN \rightarrow ENAME$  also holds.

**Partial Functional Dependency –**

A Functional Dependency in which one or more non-key attributes are functionally depending on a part of the primary key is called partial functional dependency. or where the determinant consists of key attributes, but not the entire primary key, and the determined consist of non-key attributes.

For example, consider a Relation R (A, B, C, D, E) having FD:  $AB \rightarrow CDE$  where PK is AB.

Then, { $A \rightarrow C$ ;  $A \rightarrow D$ ;  $A \rightarrow E$ ;  $B \rightarrow C$ ;  $B \rightarrow D$ ;  $B \rightarrow E$ } all are Partial Dependencies.

**Transitive Dependency –**

Given a relation R (A, B, C) then dependency like  $A \rightarrow B$ ,  $B \rightarrow C$  is a transitive dependency, since  $A \rightarrow C$  is implied.

In the above Figure

$SSN \twoheadrightarrow DMGRSSN$  is a transitive FD

{since  $SSN \twoheadrightarrow DNUMBER$  and  $DNUMBER \twoheadrightarrow DMGRSSN$  hold}

$SSN \twoheadrightarrow NAME$  is non-transitive FD since there is no set of attributes X where  $SSN \twoheadrightarrow X$  and  $X \twoheadrightarrow ENAME$ .

**Trivial Functional Dependency –**

Some functional dependencies are said to be trivial because they are satisfied by all relations. Functional dependency of form  $A \rightarrow B$  is trivial if  $B \subseteq A$ . or

A trivial Functional Dependency is the one where RHS is a subset of LHS.

Example,  $A \rightarrow A$  is satisfied by all relations involving attribute A.

$SSN \rightarrow SSN$

$PNUMBER \rightarrow PNUMBER$

$SSN \ PNUMBER \rightarrow PNUMBER$

$SSN \ PNUMBER \rightarrow SSN \ PNUMBER$

### Non-Trivial Functional Dependency –

Non-Trivial Functional Dependency can be categorized into –

- Complete Non-Trivial Functional Dependency
- Semi Non-Trivial Functional Dependency

Complete Non-Trivial Functional Dependency –

A Functional Dependency is completely non-trivial if none of the RHS attributes are part of the LHS attributes.

Example,  $SSN \rightarrow Ename$ ,

$PNUMBER \rightarrow PNAME$

$PNUMBER \rightarrow BDATE$  X

**Semi Non-Trivial Functional Dependencies** – A Functional Dependency is semi non-trivial if at least one of the RHS attributes are not part of the LHS attributes.

{TRIVIAL + NONTRIVIAL}

Question 1 :

A	B	C
1	1	1
1	2	1
2	1	2
2	2	3

Identify Non-Trivial Functional Dependency?

Solution:

S.NO	Dependencies	Non-Trivial FD?
1	$A \rightarrow B$	x
2	$A \rightarrow C$	x
3	$A \rightarrow BC$	x
4	$B \rightarrow A$	x
5	$B \rightarrow C$	x
6	$B \rightarrow AC$	x
7	$C \rightarrow A$	√
8	$C \rightarrow B$	x
9	$C \rightarrow AB$	x
10	$AB \rightarrow C$	√
11	$BC \rightarrow A$	√
12	$AC \rightarrow B$	x

$A \rightarrow B$  is not a non-trivial FD because, for 2, it has two outputs. i.e  $2 \rightarrow 2$  and  $2 \rightarrow 3$ .

for  $AB \rightarrow C$ ,  $11 \rightarrow 1$ ,  $12 \rightarrow 1$ ,  $21 \rightarrow 2$ ,  $22 \rightarrow 3$ , so Non-trivial.

Question 2: R (A B C D) AB {Candidate Key}  $A \rightarrow C$   $B \rightarrow D$ . Where is the redundancy existing?

Solution: (A C) and (B D) is suffering from redundancy.

Question 3: Consider a relation with schema R (A, B, C, D) and FDs { $AB \rightarrow C$ ,  $C \rightarrow D$ ,  $D \rightarrow A$ }. a. What are some of the nontrivial FDs that can be inferred from the given FDs?

Some examples:

C → ACD  
D → AD  
AB → ABCD  
AC → ACD  
BC → ABCD  
BD → ABCD  
CD → ACD  
ABC → ABCD  
ABD → ABCD  
BCD → ABCD

#### Inference Rules for Functional Dependencies –

Let  $S$  be the set of functional dependencies that are specified on relation schema  $R$ . Numerous other dependencies can be inferred or deduced from the functional dependencies in  $S$ .

Example:

Let  $S = \{A \rightarrow B, B \rightarrow C\}$

A multivalued dependency occurs when the presence of one or more rows in a table implies the presence of one or more other rows in that same table. Put another way, two attributes (or columns) in a table are independent of one another, but both depend on a third attribute. A multivalued dependency prevents the normalization standard Fourth Normal Form (4NF).

#### FUNCTIONAL DEPENDENCY VS. MULTIVALUED DEPENDENCY

To understand this, let's revisit what a functional dependency is.

Remember that if an attribute  $X$  uniquely determines an attribute  $Y$ , then  $Y$  is functionally dependent on  $X$ . This is written as  $X \rightarrow Y$ . For example, in the Students table below, the Student\_Name determines the Major:

#### Students

Student_Name	Major
Ravi	Art History
Beth	Chemistry

This functional dependency can be written: Student\_Name → Major. Each Student\_Name determines exactly one Major, and no more.

Now, perhaps we also want to track the sports these students take. We might think the easiest way to do this is to just add another column, Sport:

#### Students

Student_Name	Major	Sport
Ravi	Art History	Soccer
Ravi	Art History	Volleyball
Ravi	Art History	Tennis
Beth	Chemistry	Tennis
Beth	Chemistry	Soccer

The problem here is that both Ravi and Beth play multiple sports. We need to add a new row for every additional sport.

This table has introduced a multivalued dependency because the major and the sport are independent of one another but both depend on the student.

Note that this is a very simple example and easily identifiable — but this could become a problem in a large, complex database.

A multivalued dependency is written  $X \twoheadrightarrow Y$ . In this case:

Student\_Name  $\twoheadrightarrow$  Major  
Student\_Name  $\twoheadrightarrow$  Sport

This is read as "Student\_Name multidetermined Major" and "Student Name multidetermined Sport."

A multivalued dependency always requires at least three attributes because it consists of at least two attributes that are dependent on a third.

### **Multivalued dependency and normalization**

A table with a multivalued dependency violates the normalization standard of Fourth Normal Form (4NF) because it creates unnecessary redundancies and can contribute to inconsistent data. To bring this up to 4NF, we can break this into two tables.

The table below now has a functional dependency of Student\_Name  $\rightarrow$  Major, and no multidependencies:

#### **Students & Majors**

Student_Name	Major
Ravi	Art History
Ravi	Art History
Ravi	Art History
Beth	Chemistry
Beth	Chemistry

While this table also has a single functional dependency of Student\_Name  $\rightarrow$  Sport:

#### **Students & Sports**

Student_Name	Sport
Ravi	Soccer
Ravi	Volleyball
Ravi	Tennis
Beth	Tennis
Beth	Soccer

It's clear that normalization is often addressed by simplifying complex tables so that they contain information related to a single idea or theme, rather than trying to make a single table contain too much disparate information.

**Numerical on Functional Dependency: -**

1. Let  $R = (A, B, C, D, E, F)$  be a relation scheme with the following dependencies:  $C \rightarrow F$ ,  $E \rightarrow A$ ,  $EC \rightarrow D$ ,  $A \rightarrow B$ . Which of the following is a key for  $R$ ?

- (a) CD      (b) EC      (c) AE      (d) AC

Ans: option (b)

Explanation:

Find the closure set of all the options give. If any closure covers all the attributes of the relation  $R$  then that is the key.

2. Consider a relation scheme  $R = (A, B, C, D, E, H)$  on which the following functional dependencies hold:  $\{A \rightarrow B, BC \rightarrow D, E \rightarrow C, D \rightarrow A\}$ . What are the candidate keys of  $R$ ?

- (a) AE, BE  
(b) AE, BE, DE  
(c) AEH, BEH, BCH  
(d) AEH, BEH, DEH

Ans: option (d)

Explanation:

As explained in question 1, if any closure includes all attributes of a table then it becomes the candidate key.

Closure of AEH = AEHB { $A \rightarrow B$ }  
= AEHBC { $E \rightarrow C$ }  
= AEHBCD { $BC \rightarrow D$ }

GATE-2005(IT)

5. In a schema with attributes  $A, B, C, D$  and  $E$ , following set of functional dependencies are given:

- $A \rightarrow B$   
 $A \rightarrow C$   
 $CD \rightarrow E$   
 $B \rightarrow D$   
 $E \rightarrow A$

Which of the following functional dependencies is NOT implied by the above set?

- (a)  $CD \rightarrow AC$       (b)  $BD \rightarrow CD$       (c)  $BC \rightarrow CD$       (d)  $AC \rightarrow BC$

Ans: option (b)

Explanation:

For every option given, find the closure set of left side of each FD. If the closure set of left side contains the right side of the FD, then the particular FD is implied by the given set.

Option (a): Closure set of  $CD = CDEAB$ . Therefore  $CD \rightarrow AC$  can be derived from the given set of FDs.

Option (c): Closure set of  $BC = BCDEA$ . Therefore  $BC \rightarrow CD$  can be derived from the given set of FDs.

Option (d): Closure set of  $AC = ACBDE$ . Therefore  $AC \rightarrow BC$  can be derived from the given set of FDs.

Option (b): Closure set of  $BD = BD$ . Therefore  $BD \rightarrow CD$  cannot be derived from the given set of FDs.

## Normalization

### First Normal Form

First Normal Form is defined in the definition of relations (tables) itself. This rule defines that all the attributes in a relation must have atomic domains. The values in an atomic domain are indivisible units.

Course	Content
Programming	Java,C++
Web	HTML,PHP,ASP

We re-arrange the relation (table) as below, to convert it to First Normal Form.

### Programming

Course	Content
Programming	JAVA
Programming	C++
Web	HTML
Web	PHP
Web	ASP

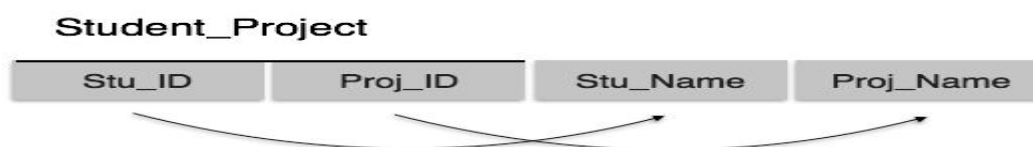
### Second Normal Form

Before we learn about the second normal form, we need to understand the following –

**Prime attribute** – An attribute, which is a part of the prime-key, is known as a prime attribute.

**Non-prime attribute** – An attribute, which is not a part of the prime-key, is said to be a non-prime attribute.

If we follow second normal form, then every non-prime attribute should be fully functionally dependent on prime key attribute. That is, if  $X \rightarrow A$  holds, then there should not be any proper subset  $Y$  of  $X$ , for which  $Y \rightarrow A$  also holds true.



We see here in Student\_Project relation that the prime key attributes are Stu\_ID and Proj\_ID. According to the rule, non-key attributes, i.e. Stu\_Name and Proj\_Name must be dependent upon both and not on any of the prime key attribute individually. But we find that Stu\_Name can be identified by Stu\_ID and Proj\_Name can be identified by Proj\_ID independently. This is called **partial dependency**, which is not allowed in Second Normal Form.

### Student

Stu_ID	Stu_Name	Proj_ID
--------	----------	---------

### Project

Proj_ID	Proj_Name
---------	-----------

We broke the relation in two as depicted in the above picture. So there exists no partial dependency.

### Third Normal Form

For a relation to be in Third Normal Form, it must be in Second Normal form and the following must satisfy –

No non-prime attribute is transitively dependent on prime key attribute.

For any non-trivial functional dependency,  $X \rightarrow A$ , then either –


X is a superkey or,

A is prime attribute.



## STUDENT\_DETAILS

Stu_ID	Stu_Name	City	Zip
--------	----------	------	-----



We find that in the above Student\_detail relation, Stu\_ID is the key and only prime key attribute. We find that City can be identified by Stu\_ID as well as Zip itself. Neither Zip is a superkey nor is City a prime attribute. Additionally,  $\text{Stu\_ID} \rightarrow \text{Zip} \rightarrow \text{City}$ , so there exists transitive dependency.

To bring this relation into third normal form, we break the relation into two relations as follows –.

### Student\_Details

Stu_ID	Stu_Name	Zip
--------	----------	-----

### Zip Codes

Zip	City
-----	------

### Boyce-Codd Normal Form

Boyce-Codd Normal Form (BCNF) is an extension of Third Normal Form on strict terms. BCNF states that for any non-trivial functional dependency,  $X \rightarrow A$ , X must be a super-key.

In the above image, Stu\_ID is the super-key in the relation Student\_Detail and Zip is the super-key in the relation ZipCodes. So,

$\text{Stu\_ID} \rightarrow \text{Stu\_Name, Zip}$

and

$\text{Zip} \rightarrow \text{City}$

Which confirms that both the relations are in BCNF.

### Fourth Normal Form (4NF)

When attributes in a relation have multi-valued dependency, further Normalization to 4NF and 5NF are required. Let us first find out what multi-valued dependency is.

A multi-valued dependency is a typical kind of dependency in which each and every attribute within a relation depends upon the other, yet none of them is a unique primary key.

We will illustrate this with an example. Consider a vendor supplying many items to many projects in an organization. The following are the assumptions:

A vendor is capable of supplying many items.

A project uses many items.

A vendor supplies to many projects.

An item may be supplied by many vendors.

A multi valued dependency exists here because all the attributes depend upon the other and yet none of them is a primary key having unique value.

Vendor Code	Item Code	Project No.
V1	I1	P1
V1	I2	P1
V1	I1	P3
V1	I2	P3

V2	I2	P1
V2	I3	P1
V3	I1	P2
V3	I1	P3

The table can be expressed as the two 4NF relations given as following. The fact that vendors are capable of supplying certain items and that they are assigned to supply for some projects is independently specified in the 4NF relation.

**Vendor-Supply**

Vendor Code	Item Code
V1	I1
V1	I2
V2	I2
V2	I3
V3	I1

**Vendor-Project**

Vendor Code	Project No.
V1	P1
V1	P3
V2	P1
V3	P2

**Fifth Normal Form (5NF)**

These relations still have a problem. While defining the 4NF we mentioned that all the attributes depend upon each other. While creating the two tables in the 4NF, although we have preserved the dependencies between Vendor Code and Item code in the first table and Vendor Code and Project No. in the second table, we have lost the relationship between Item Code and Project No. If there were a primary key then this loss of dependency would not have occurred. In order to revive this relationship, we must add a new table like the following. Please note that during the entire process of normalization, this is the only step where a new table is created by joining two attributes, rather than splitting them into separate tables.

Project No.	Item Code
P1	I1
P1	I2
P2	I1
P3	I1
P3	I3