---

## Module 5

## Relational-Database Design

**Pitfalls in Relational-Database designs, Concept of normalization, Function Dependencies, First Normal Form, 2NF, 3NF, BCNF.**

## Normalization

Normalization is the process of minimizing redundancy from a relation or set of relations. Redundancy in relation may cause insertion, deletion, and update anomalies. So, it helps to minimize the redundancy in relations. Normal forms are used to eliminate or reduce redundancy in database tables. A large database defined as a single relation may result in data duplication. This repetition of data may result in:

- Making relations very large.
- It isn't easy to maintain and update data as it would involve searching many records in relation.
- Wastage and poor utilization of disk space and resources.
- The likelihood of errors and inconsistencies increases.

So to handle these problems, we should analyze and decompose the relations with redundant data into smaller, simpler, and well-structured relations that are satisfy desirable properties. Normalization is a process of decomposing the relations into relations with fewer attributes.

**What is Normalization?**

- Normalization is the process of organizing the data in the database.
- Normalization is used to minimize the redundancy from a relation or set of relations. It is also used to eliminate undesirable characteristics like Insertion, Update, and Deletion Anomalies.
- Normalization divides the larger table into smaller and links them using relationships.
- The normal form is used to reduce redundancy from the database table.

**Why do we need Normalization?**

The main reason for normalizing the relations is removing these anomalies. Failure to eliminate anomalies leads to data redundancy and can cause data integrity and other problems as the database grows. Normalization consists of a series of guidelines that helps to guide you in creating a good database structure.

**Data modification anomalies can be categorized into three types:**

- **Insertion Anomaly:** Insertion Anomaly refers to when one cannot insert a new tuple into a relationship due to lack of data.

- **Deletion Anomaly:** The delete anomaly refers to the situation where the deletion of data results in the unintended loss of some other important data.

- **Updatation Anomaly:** The update anomaly is when an update of a single data value requires multiple rows of data to be updated.

## Types of Normal Forms:

Normalization works through a series of stages called Normal forms. The normal forms apply to individual relations. The relation is said to be in particular normal form if it satisfies constraints.

**Following are the various types of Normal forms:**

| | 1NF | 2NF | 3NF | 4NF | 5NF |
|---|---|---|---|---|---|
| **Decomposition of Relation** | R | $R_{11}$ | $R_{21}$ | $R_{31}$ | $R_{41}$ |
| | | $R_{12}$ | $R_{22}$ | $R_{32}$ | $R_{42}$ |
| | | | $R_{23}$ | $R_{33}$ | $R_{43}$ |
| | | | | $R_{34}$ | $R_{44}$ |
| | | | | | $R_{45}$ |
| **Conditions** | Eliminate Repeating Groups | Eliminate Partial Functional Dependency | Eliminate Transitive Dependency | Eliminate Multi-values Dependency | Eliminate Join Dependency |

| Normal Form | Description |
|---|---|
| 1NF | A relation is in 1NF if it contains an atomic value. |
| 2NF | A relation will be in 2NF if it is in 1NF and all non-key attributes are fully functional dependent on the primary key. |
| 3NF | A relation will be in 3NF if it is in 2NF and no transition dependency exists. |
| BCNF | A stronger definition of 3NF is known as Boyce Codd's normal form. |
| 4NF | A relation will be in 4NF if it is in Boyce Codd's normal form and has no multi-valued dependency. |
| 5NF | A relation is in 5NF. If it is in 4NF and does not contain any join dependency, joining should be lossless. |

## Advantages of Normalization

- o Normalization helps to minimize data redundancy.
- o Greater overall database organization.
- o Data consistency within the database.
- o Much more flexible database design.
- o Enforces the concept of relational integrity.

## Disadvantages of Normalization

- o You cannot start building the database before knowing what the user needs.
- o The performance degrades when normalizing the relations to higher normal forms, i.e., 4NF, 5NF.
- o It is very time-consuming and difficult to normalize relations of a higher degree.
- o Careless decomposition may lead to a bad database design, leading to serious problems.

## Concept of Functional dependency:

The functional dependency is a relationship that exists between two attributes. It typically exists between the primary key and non-key attribute within a table.

1.  X  →  Y

The left side of FD is known as a determinant, the right side of the production is known as a dependent.

**For example:**

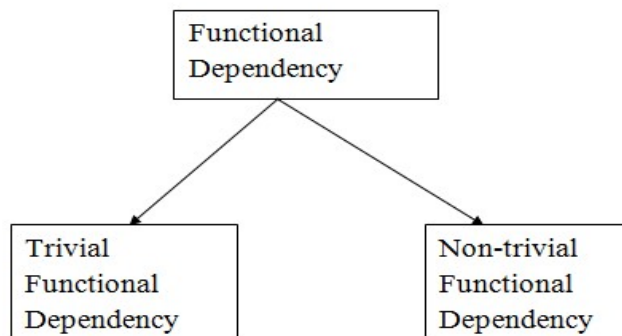Assume we have an employee table with attributes: Emp_Id, Emp_Name, Emp_Address.

Here Emp_Id attribute can uniquely identify the Emp_Name attribute of employee table because if we know the Emp_Id, we can tell that employee name associated with it.

Functional dependency can be written as:

1.  Emp_Id → Emp_Name

We can say that Emp_Name is functionally dependent on Emp_Id.

**Types of Functional dependency**



**1. Trivial functional dependency**

o   A → B has trivial functional dependency if B is a subset of A.

o   The following dependencies are also trivial like: A → A, B → B

**Example:**

1.  Consider a table with two columns Employee_Id and Employee_Name.

2. {Employee_id, Employee_Name}  →   Employee_Id is a trivial functional dependency as

3. Employee_Id is a subset of {Employee_Id, Employee_Name}.

4. Also, Employee_Id → Employee_Id and Employee_Name  →   Employee_Name are trivial dependencies too.

## 2. Non-trivial functional dependency

- A → B has a non-trivial functional dependency if B is not a subset of A.

- When A intersection B is NULL, then A → B is called as complete non-trivial.

**Example:**

1. ID  →   Name,
2. Name  →   DOB

## 1. 1NF: First Normal Form –

If a relation contain composite or multi-valued attribute, it violates first normal form or a relation is in first normal form if it does not contain any composite or multi-valued attribute. A relation is in first normal form if every attribute in that relation is **singled valued attribute**.

- **Example 1** – Relation STUDENT in table 1 is not in 1NF because of multi-valued attribute STUD_PHONE. Its decomposition into 1NF has been shown in table 2.

| STUD_NO | STUD_NAME | STUD_PHONE | STUD_STATE | STUD_COUNTRY |
|---------|-----------|------------|------------|--------------|
| 1 | RAM | 9716271721, 9871717178 | HARYANA | INDIA |
| 2 | RAM | 9898297281 | PUNJAB | INDIA |
| 3 | SURESH | | PUNJAB | INDIA |

**Table 1**

Conversion to first normal form

| STUD_NO | STUD_NAME | STUD_PHONE | STUD_STATE | STUD_COUNTRY |
|---------|-----------|------------|------------|--------------|
| 1 | RAM | 9716271721 | HARYANA | INDIA |
| 1 | RAM | 9871717178 | HARYANA | INDIA |
| 2 | RAM | 9898297281 | PUNJAB | INDIA |
| 3 | SURESH | | PUNJAB | INDIA |

**Table 2**

- **Example 2 –**

  ID  Name  Courses

```
------------------
1   A    c1, c2
2   E    c3
3   M    C2, c3
```

In the above table Course is a multi-valued attribute so it is not in 1NF.

Below Table is in 1NF as there is no multi-valued attribute

```
ID   Name   Course
------------------
1    A      c1
1    A      c2
2    E      c3
3    M      c2
3    M      c3
```

## 2.  2NF: Second Normal Form –

To be in second normal form, a relation must be in first normal form and relation must not contain any partial dependency.

A relation is in 2NF if it has **No Partial Dependency,** i.e., no non-prime attribute (attributes which are not part of any candidate key) is dependent on any proper subset of any candidate key of the table.

**Partial Dependency** – If the proper subset of candidate key determines non-prime attribute, it is called partial dependency.

**Example 1** – Consider table-3 as following below.

| STUD_NO | COURSE_NO | COURSE_FEE |
|---------|-----------|------------|
| 1       | C1        | 1000       |
| 2       | C2        | 1500       |
| 1       | C4        | 2000       |
| 4       | C3        | 1000       |
| 4       | C1        | 1000       |
| 2       | C5        | 2000       |

{Note that, there are many courses having the same course fee. }

Here,

COURSE_FEE cannot alone decide the value of COURSE_NO or STUD_NO;

COURSE_FEE together with STUD_NO cannot decide the value of COURSE_NO;

COURSE_FEE together with COURSE_NO cannot decide the value of STUD_NO;

Hence,

COURSE_FEE would be a non-prime attribute, as it does not belong to the one only candidate key {STUD_NO, COURSE_NO} ;

But, COURSE_NO -> COURSE_FEE, i.e., COURSE_FEE is dependent on COURSE_NO, which is a proper subset of the candidate key. Non-prime attribute COURSE_FEE is dependent on a proper subset of the candidate key, which is a partial dependency and so this relation is not in 2NF.

To convert the above relation to 2NF,

we need to split the table into two tables such as :

Table 1: STUD_NO, COURSE_NO

Table 2: COURSE_NO, COURSE_FEE

| Table 1 | | Table 2 | |
|---------|-----------|-----------|------------|
| STUD_NO | COURSE_NO | COURSE_NO | COURSE_FEE |
| 1 | C1 | C1 | 1000 |
| 2 | C2 | C2 | 1500 |
| 1 | C4 | C3 | 1000 |
| 4 | C3 | C4 | 2000 |
| 4 | C1 | C5 | 2000 |
| 2 | C5 | | |

**NOTE:** 2NF tries to reduce the redundant data getting stored in memory. For instance, if there are 100 students taking C1 course, we don't need to store its Fee as 1000 for all the 100 records, instead, once we can store it in the second table as the course fee for C1 is 1000.

**Example 2 –** Consider following functional dependencies in relation  R (A,  B , C,  D )

AB -> C  [A and B together determine C]

BC -> D  [B and C together determine D]

In the above relation, AB is the only candidate key and there is no partial dependency, i.e., any proper subset of AB doesn't determine any non-prime attribute.

## 3. 3NF:Third Normal Form –

A relation is in third normal form, if there is **no transitive dependency** for non-prime attributes as well as it is in second normal form. A relation is in 3NF if **at least one of the following condition holds** in every non-trivial function dependency X –> Y

1. X is a super key.
2. Y is a prime attribute (each element of Y is part of some candidate key).

| STUD_NO | STUD_NAME | STUD_STATE | STUD_COUNTRY | STUD_AGE |
|---------|-----------|------------|--------------|----------|
| 1 | RAM | HARYANA | INDIA | 20 |
| 2 | RAM | PUNJAB | INDIA | 19 |
| 3 | SURESH | PUNJAB | INDIA | 21 |

**Table 4**

**Transitive dependency** – If A->B and B->C are two FDs then A->C is called transitive dependency.

**Example 1** – In relation STUDENT given in Table 4,

FD set: {STUD_NO -> STUD_NAME, STUD_NO -> STUD_STATE, STUD_STATE -> STUD_COUNTRY, STUD_NO -> STUD_AGE}

Candidate Key: {STUD_NO}

For this relation in table 4,
 STUD_NO -> STUD_STATE and STUD_STATE -> STUD_COUNTRY are true.
So STUD_COUNTRY is transitively dependent on STUD_NO.
 It violates the third normal form. To convert it in third normal form, we will decompose the relation STUDENT (STUD_NO, STUD_NAME, STUD_PHONE, STUD_STATE, STUD_COUNTRY_STUD_AGE) as:
STUDENT (STUD_NO, STUD_NAME, STUD_PHONE, STUD_STATE, STUD_AGE)
STATE_COUNTRY (STATE, COUNTRY)

**Example 2 –** Consider relation R(A, B, C, D, E)

A -> BC,

CD -> E,

B -> D,

E -> A

All possible candidate keys in above relation are {A, E, CD, BC} All attributes are on right sides of all functional dependencies are prime.


## 4. BCNF: Boyce-Codd Normal Form–

A relation R is in BCNF if R is in Third Normal Form and for every FD, LHS is super key. A relation is in BCNF iff in every non-trivial functional dependency X –> Y, X is a super key.

**Example 1 –** Find the highest normal form of a relation R(A,B,C,D,E) with FD set as {BC->D, AC->BE, B->E}


Step 1. As we can see, (AC)+ ={A,C,B,E,D} but none of its subset can determine all attribute of relation, So AC will be candidate key. A or C can't be derived from any other attribute of the relation, so there will be only 1 candidate key {AC}.


Step 2. Prime attributes are those attributes that are part of candidate key {A, C} in this example and others will be non-prime {B, D, E} in this example.


Step 3. The relation R is in 1st normal form as a relational DBMS does not allow multi-valued or composite attribute.


The relation is in 2nd normal form because BC->D is in 2nd normal form (BC is not a proper subset of candidate key AC) and AC->BE is in 2nd normal form (AC is candidate key) and B->E is in 2nd normal form (B is not a proper subset of candidate key AC).
The relation is not in 3rd normal form because in BC->D (neither BC is a super key nor D is a prime attribute) and in B->E (neither B is a super key nor E is a prime attribute) but to satisfy 3rd normal for, either LHS of an FD should be super key or RHS should be prime

attribute.

So the highest normal form of relation will be 2nd Normal form.

**Example 2 –** For example consider relation R(A, B, C)

A -> BC,

B ->

A and B both are super keys so above relation is in BCNF.

**Key Points –**

- BCNF is free from redundancy.

- If a relation is in BCNF, then 3NF is also satisfied.

- If all attributes of relation are prime attribute, then the relation is always in 3NF.

- A relation in a Relational Database is always and at least in 1NF form.

- Every Binary Relation ( a Relation with only 2 attributes ) is always in BCNF.

- If a Relation has only singleton candidate keys( i.e. every candidate key consists of only 1 attribute), then the Relation is always in 2NF( because no Partial functional dependency possible).

- Sometimes going for BCNF form may not preserve functional dependency. In that case go for BCNF only if the lost FD(s) is not required, else normalize till 3NF only.

- There are many more Normal forms that exist after BCNF, like 4NF and more. But in real world database systems it's generally not required to go beyond BCNF.