Prepared By: Dr. H P Patra, GVPCE(A), Kommadi,9861012411

Syllabus:

**Schema refinement – Problems Caused by redundancy – Decompositions – Problem related to decomposition – reasoning about FDS – FIRST, SECOND, THIRD Normal forms – BCNF– Schema refinement in Data base Design – Multi valued Dependencies – FOURTH Normal Form.**

# Schema Refinement:

It is the process intended to address and refined the relation based on decomposition.

# The Problem of redundancy in Database

**Redundancy** means having multiple copies of same data in the database. This problem arises when a database is not normalized. Suppose a table of student details attributes are: student Id, student name, college name, college rank, course opted.

| Student_ID | Name | Contact | College | Course | Rank |
|---|---|---|---|---|---|
| 100 | Himanshu | 7300934851 | GEU | Btech | 1 |
| 101 | Ankit | 7900734858 | GEU | Btech | 1 |
| 102 | Aysuh | 7300936759 | GEU | Btech | 1 |
| 103 | Ravi | 7300901556 | GEU | Btech | 1 |

As it can be observed that values of attribute college name, college rank, course is being repeated which can lead to problems. Problems caused due to redundancy are: I

1. Insertion anomaly
2. Deletion anomaly
3. Updation anomaly.

1. **Insertion Anomaly –**
   If a student detail has to be inserted whose course is not being decided yet then insertion will not be possible till the time course is decided for student.

| Student_ID | Name | Contact | College | Course | Rank |
|---|---|---|---|---|---|
| 100 | Himanshu | 7300934851 | GEU | | 1 |

This problem happens when the insertion of a data record is not possible without adding some additional unrelated data to the record.
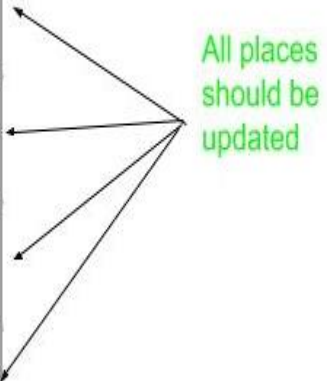
2. **Deletion Anomaly –**
   If the details of students in this table is deleted then the details of college will also get deleted which should not occur by common sense.
   This anomaly happens when deletion of a data record results in losing some unrelated information that was stored as part of the record that was deleted from a table.

3. **Updation Anomaly –**
   Suppose if the rank of the college changes then changes will have to be all over the database which will be time-consuming and computationally costly.

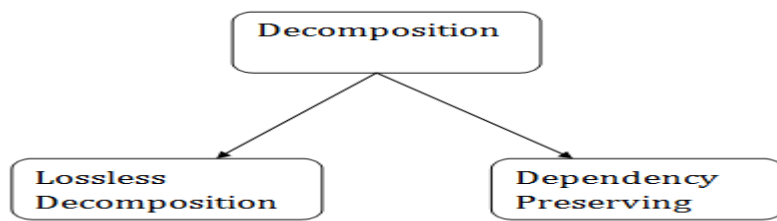| Student_ID | Name | Contact | College | Course | Rank |
|---|---|---|---|---|---|
| 100 | Himanshu | 7300934851 | GEU | Btech | 1 |
| 101 | Ankit | 7900734858 | GEU | Btech | 1 |
| 102 | Aysuh | 7300936759 | GEU | Btech | 1 |
| 103 | Ravi | 7300901556 | GEU | Btech | 1 |

All places should be updated

If updation do not occur at all places then database will be in inconsistent state

**Decompositions:**

Relational Decomposition

- o  When a relation in the relational model is not in appropriate normal form then the decomposition of a relation is required.
- o  In a database, it breaks the table into multiple tables.
- o  If the relation has no proper decomposition, then it may lead to problems like loss of information.
- o  Decomposition is used to eliminate some of the problems of bad design like anomalies, inconsistencies, and redundancy.

Types of Decomposition

Lossless Decomposition

- o   If the information is not lost from the relation that is decomposed, then the decomposition will be lossless.
- o   The lossless decomposition guarantees that the join of relations will result in the same relation as it was decomposed.
- o   The relation is said to be lossless decomposition if natural joins of all the decomposition give the original relation.

**Example:**

**EMPLOYEE_DEPARTMENT table:**

| EMP_ID | EMP_NAME | EMP_AGE | EMP_CITY | DEPT_ID | DEPT_NAME |
|--------|----------|---------|-----------|---------|-----------|
| 22 | Denim | 28 | Mumbai | 827 | Sales |
| 33 | Alina | 25 | Delhi | 438 | Marketing |
| 46 | Stephan | 30 | Bangalore | 869 | Finance |
| 52 | Katherine | 36 | Mumbai | 575 | Production |
| 60 | Jack | 40 | Noida | 678 | Testing |

The above relation is decomposed into two relations EMPLOYEE and DEPARTMENT

**EMPLOYEE table:**

| EMP_ID | EMP_NAME | EMP_AGE | EMP_CITY |
|--------|----------|---------|-----------|
| 22 | Denim | 28 | Mumbai |
| 33 | Alina | 25 | Delhi |
| 46 | Stephan | 30 | Bangalore |
| 52 | Katherine | 36 | Mumbai |
| 60 | Jack | 40 | Noida |

**DEPARTMENT table**

| DEPT_ID | EMP_ID | DEPT_NAME |
|---------|--------|-----------|
| 827 | 22 | Sales |
| 438 | 33 | Marketing |
| 869 | 46 | Finance |
| 575 | 52 | Production |
| 678 | 60 | Testing |

Now, when these two relations are joined on the common column "EMP_ID", then the resultant relation will look like:

**Employee ⋈ Department**

| EMP_ID | EMP_NAME | EMP_AGE | EMP_CITY | DEPT_ID | DEPT_NAME |
|--------|----------|---------|----------|---------|-----------|
| 22 | Denim | 28 | Mumbai | 827 | Sales |
| 33 | Alina | 25 | Delhi | 438 | Marketing |
| 46 | Stephan | 30 | Bangalore | 869 | Finance |
| 52 | Katherine | 36 | Mumbai | 575 | Production |
| 60 | Jack | 40 | Noida | 678 | Testing |

Hence, the decomposition is Lossless join decomposition.

Dependency Preserving

- o   It is an important constraint of the database.
- o   In the dependency preservation, at least one decomposed table must satisfy every dependency.
- o   If a relation R is decomposed into relation R1 and R2, then the dependencies of R either must be a part of R1 or R2 or must be derivable from the combination of functional dependencies of R1 and R2.
- o   For example, suppose there is a relation R (A, B, C, D) with functional dependency set (A->BC). The relational R is decomposed into R1(ABC) and R2(AD) which is dependency preserving because FD A->BC is a part of relation R1(ABC).

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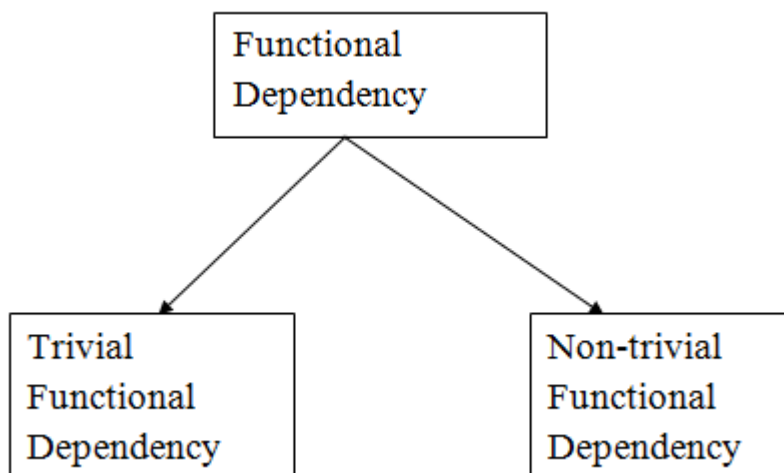
Here Emp_Id attribute can uniquely identify the Emp_Name attribute of employee table because if we know the Emp_Id, we can tell that employee name associated with it.

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

- o   A → B has a non-trivial functional dependency if B is not a subset of A.
- o   When A intersection B is NULL, then A → B is called as complete non-trivial.

**Example:**

1. ID  →  Name,
2. Name  →  DOB

# Inference Rule (IR):

- ○ The Armstrong's axioms are the basic inference rule.
- ○ Armstrong's axioms are used to conclude functional dependencies on a relational database.
- ○ The inference rule is a type of assertion. It can apply to a set of FD(functional dependency) to derive other FD.
- ○ Using the inference rule, we can derive additional functional dependency from the initial set.

The Functional dependency has 6 types of inference rule:

## 1. Reflexive Rule (IR$_1$)

In the reflexive rule, if Y is a subset of X, then X determines Y.

1. If X ⊇ Y then X  →  Y

**Example:**

1. X = {a, b, c, d, e}
2. Y = {a, b, c}

## 2. Augmentation Rule (IR$_2$)

The augmentation is also called as a partial dependency. In augmentation, if X determines Y, then XZ determines YZ for any Z.

1. If X  →  Y then XZ  →  YZ

**Example:**

1. For R(ABCD),  **if** A  →  B then AC  →  BC

## 3. Transitive Rule (IR$_3$)

In the transitive rule, if X determines Y and Y determine Z, then X must also determine Z.

1. If X  →  Y and Y  →  Z then X  →  Z

## 4. Union Rule (IR$_4$)

Union rule says, if X determines Y and X determines Z, then X must also determine Y and Z.

1. If X  →  Y and X  →  Z then X  →  YZ

**Proof:**

1. X → Y (given)
2. X → Z (given)
3. X → XY (using IR$_2$ on 1 by augmentation with X. Where XX = X)

4. $XY \rightarrow YZ$ (using $IR_2$ on 2 by augmentation with Y)
5. $X \rightarrow YZ$ (using $IR_3$ on 3 and 4)

## 5. Decomposition Rule ($IR_5$)

Decomposition rule is also known as project rule. It is the reverse of union rule.

This Rule says, if X determines Y and Z, then X determines Y and X determines Z separately.

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

    **Proof:**

    1. $X \rightarrow YZ$ (given)
    2. $YZ \rightarrow Y$ (using $IR_1$ Rule)
    3. $X \rightarrow Y$ (using $IR_3$ on 1 and 2)

## 6. Pseudo transitive Rule ($IR_6$)

In Pseudo transitive Rule, if X determines Y and YZ determines W, then XZ determines W.

1.  If $X \rightarrow Y$ and $YZ \rightarrow W$ then $XZ \rightarrow W$
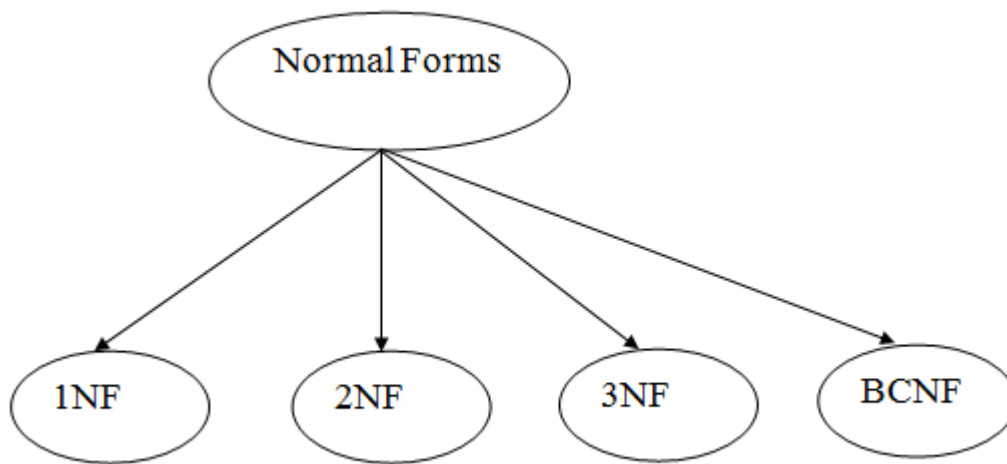
    **Proof:**

    1. $X \rightarrow Y$ (given)
    2. $WY \rightarrow Z$ (given)
    3. $WX \rightarrow WY$ (using $IR_2$ on 1 by augmenting with W)
    4. $WX \rightarrow Z$ (using $IR_3$ on 3 and 2)

# Normalization

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

## Types of Normal Forms

There are the four types of normal forms:

| 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. |
| 4NF | A relation will be in 4NF if it is in Boyce Codd normal form and has no multi-valued dependency. |
| 5NF | A relation is in 5NF if it is in 4NF and not contains any join dependency and joining should be lossless. |

# First Normal Form (1NF)

- o   A relation will be 1NF if it contains an atomic value.
- o   It states that an attribute of a table cannot hold multiple values. It must hold only single-valued attribute.
- o   First normal form disallows the multi-valued attribute, composite attribute, and their combinations.

**Example:** Relation EMPLOYEE is not in 1NF because of multi-valued attribute EMP_PHONE.

**EMPLOYEE table:**

| EMP_ID | EMP_NAME | EMP_PHONE | EMP_STATE |
|---|---|---|---|
| 14 | John | 7272826385, 9064738238 | UP |
| 20 | Harry | 8574783832 | Bihar |
| 12 | Sam | 7390372389, 8589830302 | Punjab |

The decomposition of the EMPLOYEE table into 1NF has been shown below:

| EMP_ID | EMP_NAME | EMP_PHONE | EMP_STATE |
|--------|----------|-----------|-----------|
| 14 | John | 7272826385 | UP |
| 14 | John | 9064738238 | UP |
| 20 | Harry | 8574783832 | Bihar |
| 12 | Sam | 7390372389 | Punjab |
| 12 | Sam | 8589830302 | Punjab |

# Second Normal Form (2NF)

- o    In the 2NF, relational must be in 1NF.
- o    In the second normal form, all non-key attributes are fully functional dependent on the primary key

**Example:** Let's assume, a school can store the data of teachers and the subjects they teach. In a school, a teacher can teach more than one subject.

**TEACHER table**

| TEACHER_ID | SUBJECT | TEACHER_AGE |
|------------|---------|-------------|
| 25 | Chemistry | 30 |
| 25 | Biology | 30 |
| 47 | English | 35 |
| 83 | Math | 38 |
| 83 | Computer | 38 |

In the given table, non-prime attribute TEACHER_AGE is dependent on TEACHER_ID which is a proper subset of a candidate key. That's why it violates the rule for 2NF.

To convert the given table into 2NF, we decompose it into two tables:

**TEACHER_DETAIL table:**

| TEACHER_ID | TEACHER_AGE |
|------------|-------------|
| 25 | 30 |
| 47 | 35 |

| 83 | 38 |
|----|----|

**TEACHER_SUBJECT table:**

| TEACHER_ID | SUBJECT |
|------------|---------|
| 25 | Chemistry |
| 25 | Biology |
| 47 | English |
| 83 | Math |
| 83 | Computer |

# Third Normal Form (3NF)

- o   A relation will be in 3NF if it is in 2NF and not contain any transitive partial dependency.
- o   3NF is used to reduce the data duplication. It is also used to achieve the data integrity.
- o   If there is no transitive dependency for non-prime attributes, then the relation must be in third normal form.

A relation is in third normal form if it holds atleast one of the following conditions for every non-trivial function dependency X → Y.

1.   X is a super key.
2.   Y is a prime attribute, i.e., each element of Y is part of some candidate key.

**Example:**

**EMPLOYEE_DETAIL table:**

| EMP_ID | EMP_NAME | EMP_ZIP | EMP_STATE | EMP_CITY |
|--------|----------|---------|-----------|----------|
| 222 | Harry | 201010 | UP | Noida |
| 333 | Stephan | 02228 | US | Boston |
| 444 | Lan | 60007 | US | Chicago |
| 555 | Katharine | 06389 | UK | Norwich |
| 666 | John | 462007 | MP | Bhopal |

**Super key in the table above:**

1.   {EMP_ID}, {EMP_ID, EMP_NAME}, {EMP_ID, EMP_NAME, EMP_ZIP}....so on

**Candidate key:** {EMP_ID}

**Non-prime attributes:** In the given table, all attributes except EMP_ID are non-prime.

Here, EMP_STATE & EMP_CITY dependent on EMP_ZIP and EMP_ZIP dependent on EMP_ID. The non-prime attributes (EMP_STATE, EMP_CITY) transitively dependent on super key(EMP_ID). It violates the rule of third normal form.

That's why we need to move the EMP_CITY and EMP_STATE to the new <EMPLOYEE_ZIP> table, with EMP_ZIP as a Primary key.

**EMPLOYEE table:**

| EMP_ID | EMP_NAME | EMP_ZIP |
|--------|----------|---------|
| 222 | Harry | 201010 |
| 333 | Stephan | 02228 |
| 444 | Lan | 60007 |
| 555 | Katharine | 06389 |
| 666 | John | 462007 |

**EMPLOYEE_ZIP table:**

| EMP_ZIP | EMP_STATE | EMP_CITY |
|---------|-----------|----------|
| 201010 | UP | Noida |
| 02228 | US | Boston |
| 60007 | US | Chicago |
| 06389 | UK | Norwich |
| 462007 | MP | Bhopal |

# Boyce Codd normal form (BCNF)

- o   BCNF is the advance version of 3NF. It is stricter than 3NF.
- o   A table is in BCNF if every functional dependency X → Y, X is the super key of the table.
- o   For BCNF, the table should be in 3NF, and for every FD, LHS is super key.

**Example:** Let's assume there is a company where employees work in more than one department.

**EMPLOYEE table:**

| EMP_ID | EMP_COUNTRY | EMP_DEPT | DEPT_TYPE | EMP_DEPT_NO |
|--------|-------------|----------|-----------|-------------|
| 264 | India | Designing | D394 | 283 |

| 264 | India | Testing | D394 | 300 |
| 364 | UK | Stores | D283 | 232 |
| 364 | UK | Developing | D283 | 549 |

**In the above table Functional dependencies are as follows:**

1. EMP_ID → EMP_COUNTRY
2. EMP_DEPT → {DEPT_TYPE, EMP_DEPT_NO}

**Candidate key: {EMP-ID, EMP-DEPT}**

The table is not in BCNF because neither EMP_DEPT nor EMP_ID alone are keys.

To convert the given table into BCNF, we decompose it into three tables:

**EMP_COUNTRY table:**

| EMP_ID | EMP_COUNTRY |
|---|---|
| 264 | India |
| 264 | India |

**EMP_DEPT table:**

| EMP_DEPT | DEPT_TYPE | EMP_DEPT_NO |
|---|---|---|
| Designing | D394 | 283 |
| Testing | D394 | 300 |
| Stores | D283 | 232 |
| Developing | D283 | 549 |

**EMP_DEPT_MAPPING table:**

| EMP_ID | EMP_DEPT |
|---|---|
| D394 | 283 |
| D394 | 300 |
| D283 | 232 |
| D283 | 549 |

**Functional dependencies:**

1. EMP_ID  →  EMP_COUNTRY
2. EMP_DEPT  →  {DEPT_TYPE, EMP_DEPT_NO}

**Candidate keys:**

**For the first table:** EMP_ID
**For the second table:** EMP_DEPT
**For the third table:** {EMP_ID, EMP_DEPT}

Now, this is in BCNF because left side part of both the functional dependencies is a key.

## Multivalued Dependency

- o  Multivalued dependency occurs when two attributes in a table are independent of each other but, both depend on a third attribute.
- o  A multivalued dependency consists of at least two attributes that are dependent on a third attribute that's why it always requires at least three attributes.

**Example:** Suppose there is a bike manufacturer company which produces two colors(white and black) of each model every year.

| BIKE_MODEL | MANUF_YEAR | COLOR |
|---|---|---|
| M2011 | 2008 | White |
| M2001 | 2008 | Black |
| M3001 | 2013 | White |
| M3001 | 2013 | Black |
| M4006 | 2017 | White |
| M4006 | 2017 | Black |

Here columns COLOR and MANUF_YEAR are dependent on BIKE_MODEL and independent of each other.

In this case, these two columns can be called as multivalued dependent on BIKE_MODEL. The representation of these dependencies is shown below:

1. BIKE_MODEL  → →  MANUF_YEAR
2. BIKE_MODEL  → →  COLOR

This can be read as "BIKE_MODEL multidetermined MANUF_YEAR" and "BIKE_MODEL multidetermined COLOR".

### Join Dependency

- o  Join decomposition is a further generalization of Multivalued dependencies.
- o  If the join of R1 and R2 over C is equal to relation R, then we can say that a join dependency (JD) exists.
- o  Where R1 and R2 are the decompositions R1(A, B, C) and R2(C, D) of a given relations R (A, B, C, D).
- o  Alternatively, R1 and R2 are a lossless decomposition of R.

- A JD ⋈ {R1, R2,..., Rn} is said to hold over a relation R if R1, R2,....., Rn is a lossless-join decomposition.
- The *(A, B, C, D), (C, D) will be a JD of R if the join of join's attribute is equal to the relation R.
- Here, *(R1, R2, R3) is used to indicate that relation R1, R2, R3 and so on are a JD of R.

## Inclusion Dependency

- Multivalued dependency and join dependency can be used to guide database design although they both are less common than functional dependencies.
- Inclusion dependencies are quite common. They typically show little influence on designing of the database.
- The inclusion dependency is a statement in which some columns of a relation are contained in other columns.
- The example of inclusion dependency is a foreign key. In one relation, the referring relation is contained in the primary key column(s) of the referenced relation.
- Suppose we have two relations R and S which was obtained by translating two entity sets such that every R entity is also an S entity.
- Inclusion dependency would be happen if projecting R on its key attributes yields a relation that is contained in the relation obtained by projecting S on its key attributes.
- In inclusion dependency, we should not split groups of attributes that participate in an inclusion dependency.
- In practice, most inclusion dependencies are key-based that is involved only keys.

## Fourth normal form (4NF)

- A relation will be in 4NF if it is in Boyce Codd normal form and has no multi-valued dependency.
- For a dependency A → B, if for a single value of A, multiple values of B exists, then the relation will be a multi-valued dependency.

## Example

**STUDENT**

| STU_ID | COURSE | HOBBY |
|--------|--------|-------|
| 21 | Computer | Dancing |
| 21 | Math | Singing |
| 34 | Chemistry | Dancing |
| 74 | Biology | Cricket |
| 59 | Physics | Hockey |

The given STUDENT table is in 3NF, but the COURSE and HOBBY are two independent entity. Hence, there is no relationship between COURSE and HOBBY.

In the STUDENT relation, a student with STU_ID, **21** contains two courses, **Computer** and **Math** and two hobbies, **Dancing** and **Singing**. So there is a Multi-valued dependency on STU_ID, which leads to unnecessary repetition of data.

So to make the above table into 4NF, we can decompose it into two tables:

**STUDENT_COURSE**

| STU_ID | COURSE |
|--------|--------|
| 21 | Computer |
| 21 | Math |
| 34 | Chemistry |
| 74 | Biology |
| 59 | Physics |

**STUDENT_HOBBY**

| STU_ID | HOBBY |
|--------|-------|
| 21 | Dancing |
| 21 | Singing |
| 34 | Dancing |
| 74 | Cricket |
| 59 | Hockey |

## Fifth normal form (5NF)

o   A relation is in 5NF if it is in 4NF and not contains any join dependency and joining should be lossless.
o   5NF is satisfied when all the tables are broken into as many tables as possible in order to avoid redundancy.
o   5NF is also known as Project-join normal form (PJ/NF).

## Example

| SUBJECT | LECTURER | SEMESTER |
|---------|----------|----------|
| Computer | Anshika | Semester 1 |
| Computer | John | Semester 1 |
| Math | John | Semester 1 |
| Math | Akash | Semester 2 |
| Chemistry | Praveen | Semester 1 |

In the above table, John takes both Computer and Math class for Semester 1 but he doesn't take Math class for Semester 2. In this case, combination of all these fields required to identify a valid data.

Suppose we add a new Semester as Semester 3 but do not know about the subject and who will be taking that subject so we leave Lecturer and Subject as NULL. But all three columns together acts as a primary key, so we can't leave other two columns blank.

So to make the above table into 5NF, we can decompose it into three relations P1, P2 & P3:

**P1**

| SEMESTER | SUBJECT |
|---|---|
| Semester 1 | Computer |
| Semester 1 | Math |
| Semester 1 | Chemistry |
| Semester 2 | Math |

**P2**

| SUBJECT | LECTURER |
|---|---|
| Computer | Anshika |
| Computer | John |
| Math | John |
| Math | Akash |
| Chemistry | Praveen |

**P3**

| SEMSTER | LECTURER |
|---|---|
| Semester 1 | Anshika |
| Semester 1 | John |
| Semester 1 | John |
| Semester 2 | Akash |
| Semester 1 | Praveen |