

UNIT-IV

Schema Refinement AND NORMAL FORMS

TOPICS:-

- Purpose of Normalization or Schema Refinement.
- Concept of Functional Dependency.
- Normal forms based on Functional Dependency (1NF, 2NF and 3NF).
- Concept of Surrogate Key.
- Boyce-Codd Normal form (BCNF).
- Lossy Join and Dependency preserving Decomposition.
- Fourth Normal form (4NF).

Database Schema (Definition) :-

A database schema is the skeleton structure that represents the logical view of the entire database. It defines how the data is organized and how the relations

among them are associated.

It formulates all the constraints that are to be applied on the data.

A database schema defines its entities and the relationship among them.

It contains a descriptive detail of the database, which can be depicted by means of schema diagrams. The database designers who design the schema to help programming understand the database and make it useful.

A database schema can be divided broadly into two categories:-

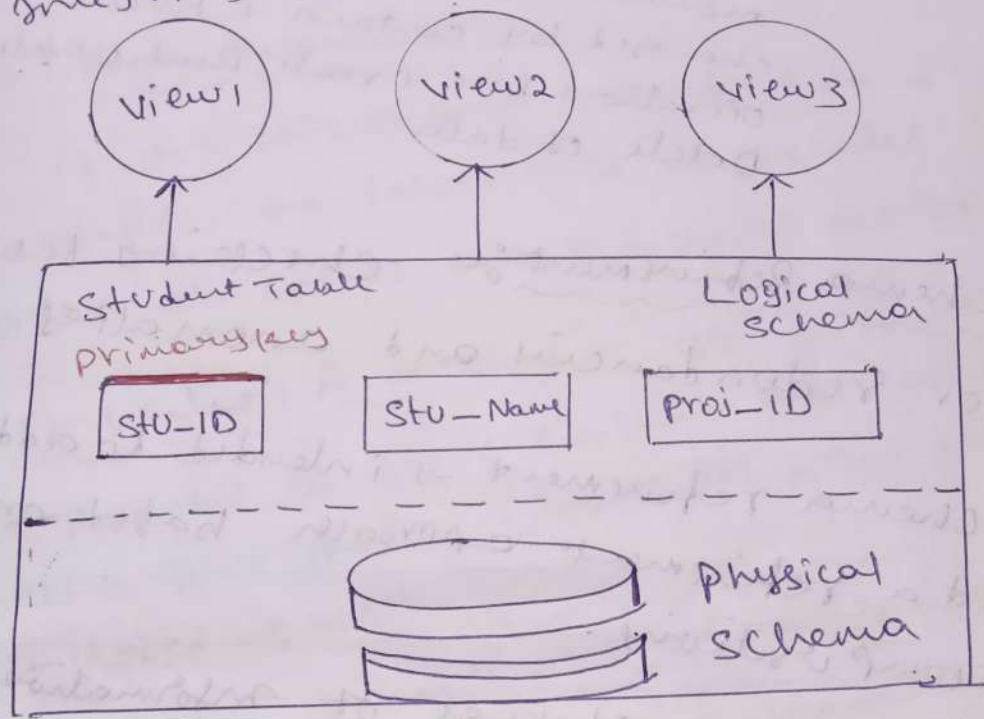
① Physical Database Schema:-

This schema pertains to the actual storage of data and its form of storage like files, indices etc. It ~~can~~ defines how the data will be stored in a secondary storage.

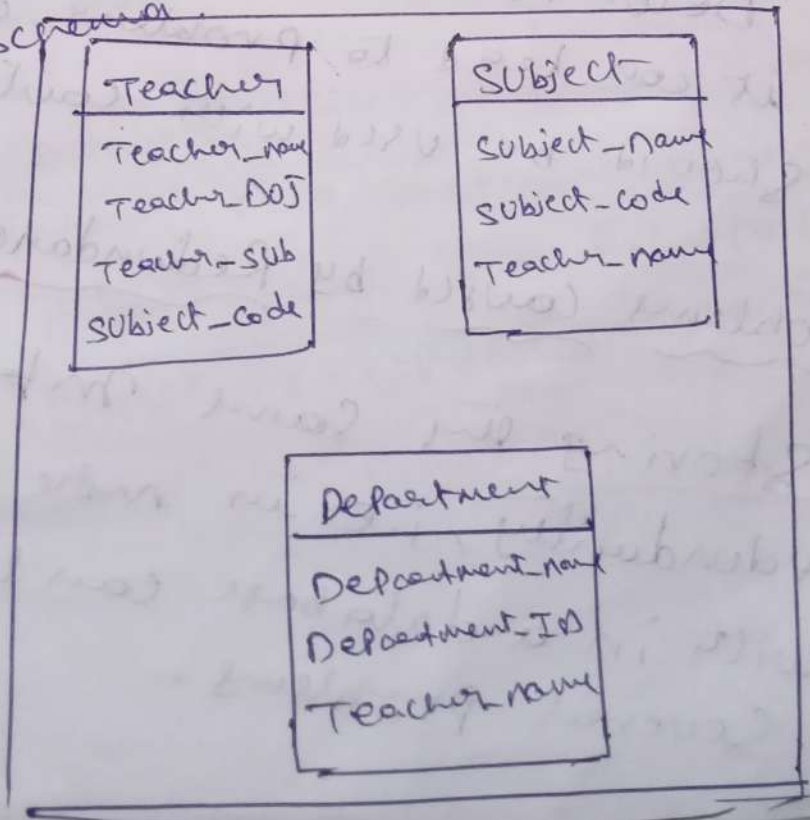
② Logical Database Schema:-

This schema all the logical constraints

that need to be applied on the data stored. It defines tables, views and integrity constraints.



* Schema is an overall description of the database. The basic structure of how the data will be stored in the database is called Schema.



Instance :- Instances are the collections of information stored at a particular moment. The instances can be changed by certain CRUD operations like Create, Read, Update, Delete of data.

Schema Refinement :- Checking tables for redundancies and anomalies.

→ Schema refinement is intended to address and a refinement approach based on decomposition.

→ Redundant storage of information is the root cause of such problems.

→ Decomposition can eliminate redundancy it can lead to problems of its own and should be used with caution.

Problems caused by Redundancy :-

Storing the same information redundantly, i.e. in more than one place within a database can lead to several problems.

→ Redundant Storage:-

Some information is stored repeatedly.

→ Update Anomaly:-

if one copy of such repeated data is updated, an inconsistency is created unless all copies are similarly updated.

→ Insertion Anomaly:- It may not be possible to store certain information unless some other, unrelated, information is stored as well.

→ Deletion Anomaly:- It may not be possible to delete certain information without losing some other, unrelated information as well.

Anomaly indicates there is some difficulty or a some problem or a constraint in any relational database need to be perform insertion, Deletion, Update is quite common.

→ Redundant Storage:-

Some information is stored repeatedly.

→ Update Anomalies:-

if one copy of such repeated data is updated, an inconsistency is created unless all copies are similarly updated.

→ Insertion Anomalies:- It may not be possible to store certain information unless some other, unrelated, information is stored as well.

→ Deletion Anomalies:- It may not be possible to delete certain information without losing some other, unrelated information as well.

Anomaly indicates there is some difficulty or a some problem or a constraint in any relational database need to be perform insertion, Deletion, Update is quite common.

Example:-

SUPPOSE a manufacturing company stores the employee details in a table named Employee that has four attributes: EMP_id storing employee's id, EMP_name storing employee's name, EMP_address storing employee's address and EMP_dept storing the department details in which the employee works. At some point of time the table looks like this:

EMP_ID	EMP_NAME	EMP ADDRESS	EMP DEPT
101	KANTH	VIZAG	D001
101	KANTH	VIZAG	D002
123	MADHU	HYD	D890
166	SATYA	chennai	D900
166	SATYA	chennai	D004

The above table is not normalized. we will see the problems that we face when a table is not normalized.

Update Anomaly:- If we want to update the address of KANTH then we have to update the same in two rows or the data will become inconsistent. The correct address gets updated in one department but not in other as per the database.

KANTH would be having two different addresses, which is not correct and would lead to inconsistent data.

Insert Anomaly:- Suppose a new employee joins the company, who is under training and currently not assigned to any Department then we could not be able to insert the data into the table if emp_dept field doesn't allow nulls.

Delete Anomaly:- Suppose, is at a point of time the company closes the department D890 then deleting the rows that are having emp_dept as D890 would also delete the information of employee Madhu since she is assigned only to this department.

TO overcome these anomalies we need to normalize the data.

- 1) Dependency preservation:- It ensures that each FD is represented in some individual resulting relation after decomposition.
- 2) loss less join property:- It guarantees that the spurious tuple generation problem does not occur with respect to the relation schema created after decomposition.

Concept of Functional Dependency

It is a method that describes the relationship between attributes.

Example 1:-

$X \rightarrow Y$ means

X determines Y

or

Y is determined by X

Determinant



(X)

→ (Y)

→ Dependent attribute

Example 2:-

$Sid \rightarrow Sname$

101 — KANTHA (Laxmi Kant)

102 — KANTH (Rajini Kant)

(or)

$Sid \rightarrow Sname$

1 — KANTH

1 — KANTH

Example 2:-

①

Valid (one student)

$Sid \rightarrow Sname$	
1	KANTH
1	KANTH

$X \rightarrow Y$ Valid

②

2 students

$Sid \rightarrow Sname$	
1	KANTH
2	KANTH

③

Valid

$Sid \rightarrow Sname$	
1	KANTH
2	MADHU

2 different students

④

$Sid \rightarrow Sname$	
1	KANTH
1	MADHU

Invalid case
one sid - two different names

Functional Dependency are of two types :-

- ① Trivial ② Non-Trivial.

① Trivial FD :- if X determines Y then Y is subset of X . Trivial FD

$X \rightarrow Y$ (Reflexive property)

↑
subset

always true.

method \Rightarrow II

$$X \rightarrow Y$$

$$\text{L.H.S } A \text{ R.H.S } \neq \emptyset$$

$$\text{sid, sname} \rightarrow \text{sid} \Rightarrow \text{sid}$$

② Non Trivial FD :- when $X \rightarrow Y$ holds true where Y is not a subset of X .
In a relationship, if attribute Y is not a subset of attribute A , then it is considered as a non-trivial dependency.

$$X \rightarrow Y \text{ then } X \cap Y = \emptyset$$

(complete nontrivial)

eg :- $\text{sid} \rightarrow \text{sname}$

$$\text{sid} \rightarrow$$

If an FD $X \rightarrow Y$ holds, where Y is not a subset of X , then it is called a non-trivial FD.

Properties of Functional Dependency:-

(1) Reflexivity:- if Y is subset of X then

$$X \rightarrow Y.$$

if X is set of attributes and Y is subset of X , then X holds Y .

(2) Augmentation:- if $X \rightarrow Y$ then $XZ \rightarrow YZ$

$$\text{Eg:- Sid} \rightarrow \text{Sname}$$

$\text{Sid, phone} \rightarrow \text{Sname, phone}$
i.e. adding dependencies, does not change its basic dependencies.

(3) Transitive:- if $X \rightarrow Y$ and $Y \rightarrow Z$ then

$$X \rightarrow Z$$

$$\text{Eg:- Sid} \rightarrow \text{Sname and Sname} \rightarrow \text{city}$$

$$\text{Sid} \rightarrow \text{city}$$

$X \rightarrow Y$ is functionally determined by X .

(4) Union:- if $X \rightarrow Y$ and $X \rightarrow Z$ then

$$X \rightarrow YZ \quad (YZ \text{ making union})$$

(5) Decomposition:- if $X \rightarrow YZ$ then

$$X \rightarrow Y \text{ and } X \rightarrow Z$$

check condition:-

$$XY \rightarrow Z, X \rightarrow Z, Y \rightarrow Z$$

we cannot decompose left hand side.

⑥ Pseudo Transitivity:-

if $X \rightarrow Y$ and $WY \rightarrow Z$ then
 $WX \rightarrow Z$

⑦ Composition:-

if $X \rightarrow Y$ and $Z \rightarrow W$ then
 $XZ \rightarrow YW$

Functional Dependency:- The value of one attribute determining the value of another attribute.

Candidate Key:-

Attribute that uniquely identifies a row in a relation could be a combination of (non-redundant) attributes.

- Each non-key field is functionally dependent on every candidate key.

Normalization:- To Remove duplicate values. It is the process of organising the attributes and labels of a relational database to minimise data redundancy with the help of functional dependencies and lossy join property.

NORMAL FORMS BASED ON FUNCTIONAL Dependency

~~As per~~

TO overcome the above Anomalies we need to normalize the data.

1) Normalization :- The most commonly used Normal forms are given below:

- 1) FIRST NORMAL FORM (1NF)
- 2) SECOND NORMAL FORM (2NF)
- 3) Third Normal form (3NF)

1) FIRST NORMAL FORM (1NF) :-

As per the rule of first normal form, an attribute (column) of a table cannot hold multiple values. It should hold only atomic values.

Example :-

SUPPOSE a company wants to store the names and contact details of its employees. It creates a table that looks like this :-

EMP-ID	EMP-Name	EMP-Address	EMP-Mobile
101	KANTH	VIZAG	8912312390
102	madhu	HYD	8812121212
103	KIRAN	CHENNAI	9900012222
104	RAM	Bangalore	7778881212
			9990000123
			8123450987

Two employees (madhu & RAM) are having two mobile numbers so the company stored them in the same field as you can see in the above table.

This table is not in 1NF as the rule says "each attribute of a table must have atomic (single) values", the emp-mobile values for employees madhu & RAM violates that rule.

To make the table complies with 1NF we should have the data like this.

EMP-ID	EMP-Name	EMP-Address	EMP-Mobile
101	KANTH	VIZAG	8912312390
102	madhu	HYD	8812121212
102	madhu	HYD	9900012222
103	KIRAN	CHENNAI	7778881212
104	RAM	Bangalore	9990000123
104	RAM	Bangalore	8123450987

Second Normal Form

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

Rule 1:

The Table of Relation must be in 1st Normal Form.

Rule 2: Every non-key attribute is fully functionally dependent on the ENTIRE primary key. (8)

All the non-prime attributes should be fully functional dependent on candidate

or

No non-prime attribute is dependent on a proper subset of any candidate key of table.

or

- Every non-key attribute must be defined by the entire key, not by only part of the key.

→ NO partial functional dependencies

No non-prime attribute Definition: if an attribute is not the part of a key, then that attribute is called non-primary attribute.

(or)

An attribute that is not part of any candidate key is known as non-prime attribute.

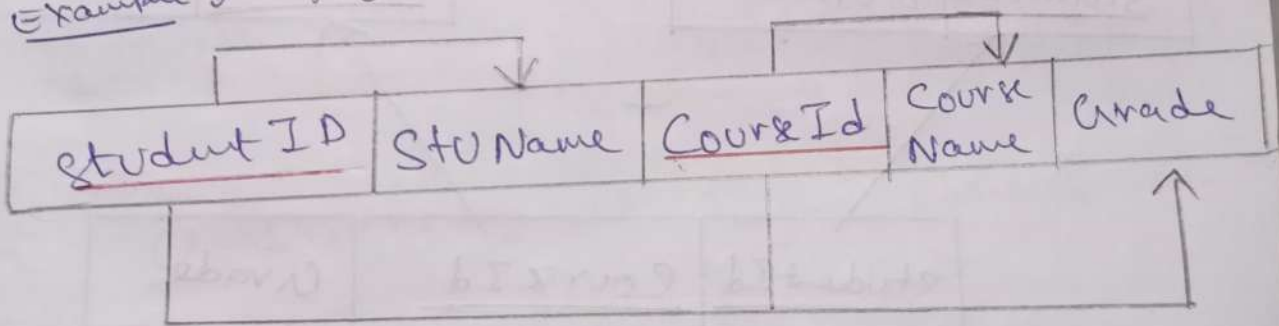
primary attribute Definition:- If an attribute is part of a key, then that attribute is called primary attribute.

Proper Subset:-

A proper subset of a set A is subset of A that is not equal to A. In other words, if B is a proper subset of A, then all elements of B are in A but A contains at least one element that is not in B.

For example:- if $A = \{1, 3, 5\}$ then $B = \{1, 5\}$ is a proper subset of A.

Example:- Functional Dependencies in student



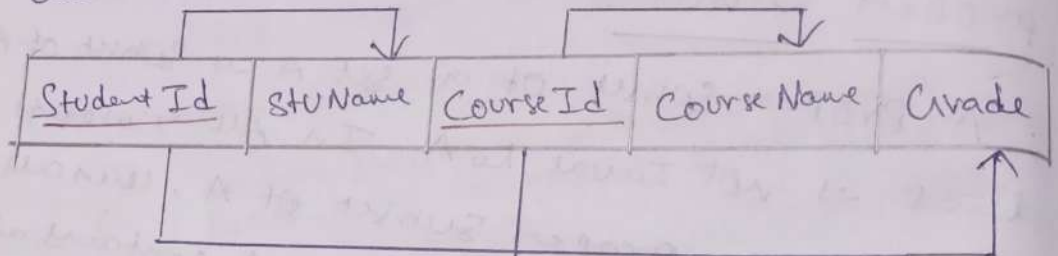
Can represent FD's with arrows as above, or

- Student ID \rightarrow Stu Name.
- Course ID \rightarrow Course Name.
- Student ID, Course ID \rightarrow Grade

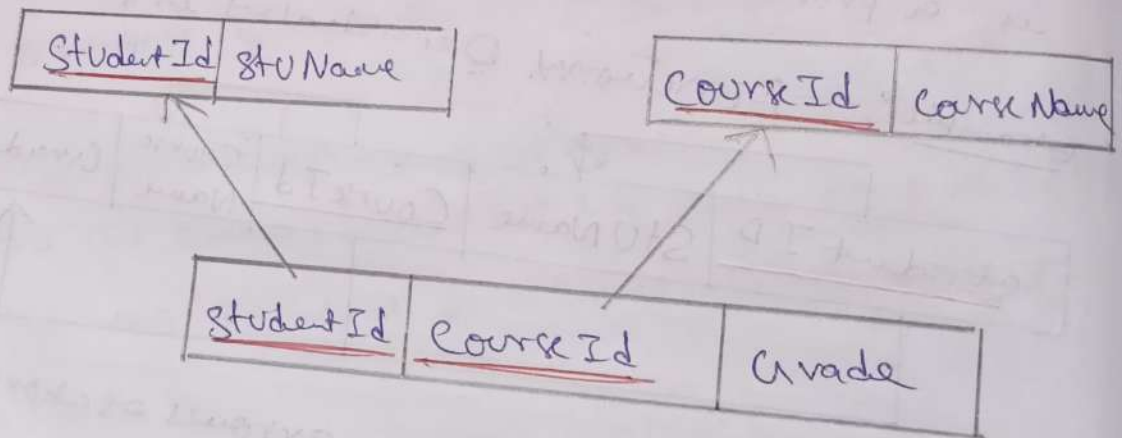
Therefore, NOT in 2nd Normal Form!!

2 Normal form: Normalizing

How do we convert the Partial dependencies into normal ones? By decomposing into more tables.



The above arrows represents functional Dependencies, below they mean Foreignkey Constraints.



THIRD NORMAL FORM

The Table must be in the second Normal form, and there should be no transitive dependency in table.

A transitive dependency is when a non-key attribute depends on another non-key attribute. This is called transitive because the primary key is a determinant for another attribute, which in turn is a determinant for a third attribute.

<u>COURSE</u>	<u>SECNUM</u>	classroom	capacity
---------------	---------------	-----------	----------

classroom \rightarrow capacity — Transitive

Any partial FDs? — NO

Any Transitive FDs? \rightarrow YES!

To eliminate the Transitive Dependency by decomposing the table into two.

<u>COURSE</u>	<u>SECNUM</u>	classroom
		<u>classroom</u> capacity

<u>Roll</u>	State	City
-------------	-------	------

Candidate key = { RollNO }

FD \rightarrow RollNO \rightarrow State

State \rightarrow City

According to transitivity

which is not in FD.

R (ABCD)

FD: $AB \rightarrow C$,

$C \rightarrow D$ ('D' Non prime attribute determined by 'C' Non prime attribute)

Prime attribute: A, B

Non prime attribute: C, D

Closure of $AB^+ = ABCD$

~~Conclusion~~

For example if 'C' & 'D' is Prime attribute then there is no problem.

$AB \rightarrow C, C \rightarrow D$

Example-II
Condition-II

note

R (ABCD)

Closure of A

FD: $AB \rightarrow CD$ ✓

$AB^+ = ABCD$

$D \rightarrow A$ ✓

CK {AB} {DB}

$DB^+ = DBAC$

PA = {A, B, D} {To make Candidate key}

N-PA = {C}

For each FD

LHS must be candidate key or SK or R.H.S must be prime attribute

BOYCE - Codd Normal Form (BCNF)

Boyce - Codd normal form (BCNF) is a normal form used in database normalization. It is a slightly stronger version of 3rd normal form.

A relational schema R is in Boyce - Codd normal form if and only if for every one of its dependencies $X \rightarrow Y$, at least one of the following conditions hold.

- ① $X \rightarrow Y$ is a trivial functional dependency (G.S.N)
- ② X is a superkey for schema R .

if a relational schema is in BCNF then all redundancies based on FD, has been removed, although other types of redundancies may still exist.

(or)

A database table is in BCNF if and only if there are no non-trivial FD of attributes on any attribute other than a super set of a candidate key.

BCNF can also be called as BCNF special case of 3rd Normal Form.

In third Normal form the L.H.S may be candidate key or super key.

ii) R.H.S ~~may~~ may be a prime attribute

In BCNF for all the F.D's, all attributes must be dependent on candidate key.
Restricted to L.H.S.

Example:-

<u>Rollno</u>	Name	<u>Voterid</u>	age
1	KANTH	P123	35
2	MOHAN	M123	30
3	Ravi	R123	25
4	Krishna	R123	20

Candidate key = { Rollno, Voterid }

Functional Dependency = {
 $\text{Rollno} \rightarrow \text{Name}$
 $\text{Rollno} \rightarrow \text{Voterid}$
 $\text{Voterid} \rightarrow \text{age}$
 $\text{Voterid} \rightarrow \text{Rollno}$

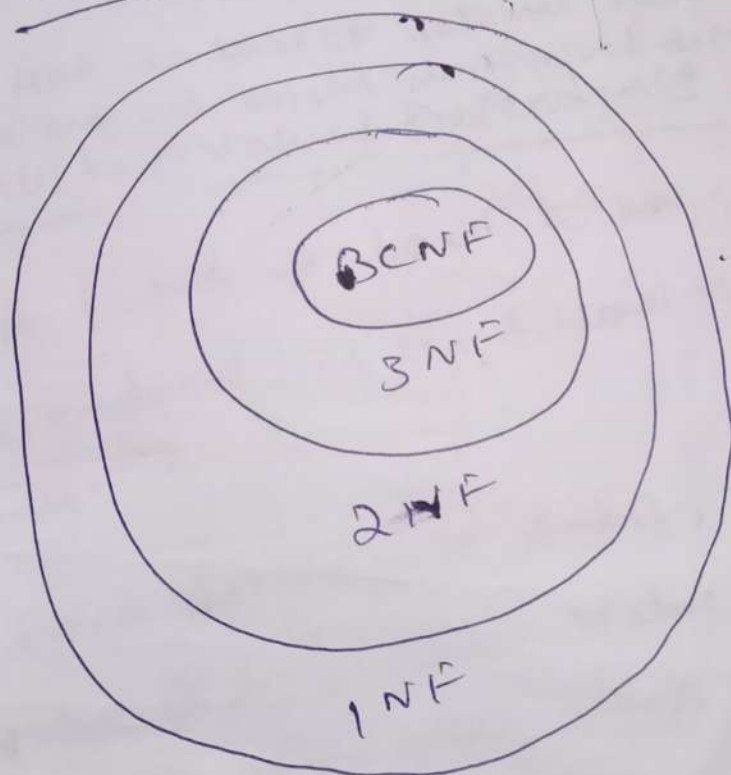
1. FD i.e $\text{Rollno} \rightarrow \text{Name}$ (FD is valid)
 \downarrow
 L.H.S is candidate key or super key

2. FD. i.e. Rollno \rightarrow Votrid (valid)
L.H.S is candidate key

3. Votrid \rightarrow age (valid.)

4. Votrid \rightarrow Rollno (valid)

L.H.S of each FD should be
candidate key or superkey



FOURTH NORMAL FORM (4NF)

For a table to satisfy the Fourth Normal Form, it should satisfy the following two conditions:

1) It should be in the Boyce-Codd Normal form

2) And, the table should not have any multi-valued Dependency.

→ 4NF is a level of database normalization where there are non-trivial multivalued dependencies other than C.V.

MULTI-Valued Dependency

A table is said to have multi-valued dependencies, if the following conditions are true:

① For a dependency $A \twoheadrightarrow B$, if for a single value of A, multiple value of B exists, then the table may have multi-valued dependencies.

2. Also, a table should have at least 3 columns for it to have a multi-valued Dependency.

3. And, for a relation $R(A, B, C)$,

if there is a multi-valued dependency between A and B, then B and C should be independent of each other.

if all these conditions are true for any relation (table), it is said to have multi-valued dependency,

A	B	C
A1	B1	C1
	B2	C2

multi-valued dependencies between

$A \rightarrow B$ and $A \rightarrow C$.

Course	Instructor	Text Book
CSE	RAM SAM KIM	C C++
CST	JIM	DBMS JAVA

As in this table multivalued Dependency exists so this table is suffering from all anomalies so to solve this splitting of table is required.

course → Instructor
course → Text book

course	Instructor
CSE	RAM
CSE	SAM
CSE	KIM
CST	RIM

course	Instructor Text Book
course	C
CSE	C++
CSE	DBMS
CST	JAVA
CST	

Lossless / Lossy Join Decomposition

When we want to normalize, we need to decompose the table. When we do the decomposition we need to follow two rules

- ① Lossless vs lossy Decomposition
- ② Dependency preserving Decomposition

① Lossless versus Lossy Decomposition :-

R

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

R₁ (AB)

A	B
1	2
2	2
3	3

R₂ (BC)

B	C
2	1
2	2
3	2

Query 1:- find the value of C if the value of A = '1'
select R₂.C from R₂ Natural Join R₁ where R₁.A = "1";

R1		R2	
A	B	B	C
1	2	2	1 ✓
1	2	2	2 ✓
1	2	3	2 ✗
2	2	2	1 ✓
2	2	2	2 ✓
2	2	3	2 ✗
3	3	2	1 ✗
3	3	2	2 ✗
3	3	3	2 ✓

Join means

Cross product + Condition
↓
(Equality)

Spurious tuples.

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

we got 5 Rows
it is call spurious.
we can call
lossy join
Decomposition.
Inconsistent.
This

we keep B as common. But B is
→ even if we keep AC, BC also we get (lossy)
→ common Attribute should be candidate key
or superkey of either R1 or R2 or ~~both~~ Both.
→ so finally $R_1(AB)$
 $R_2(AC)$.

Rules :-

- 1) $R_1 \cup R_2 = R$
 $AB \cup AC$
 $ABC = AB \cup AC$

$$R_1 \cap R_2 = \phi$$

$$AB \cap AC$$

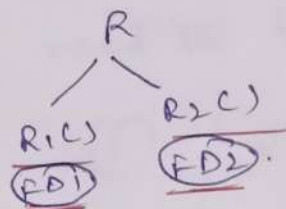
$$A = \phi$$

③ $R_1 \cap R_2 \neq \phi$
of both

Dependency Preserving Decomposition

In a Table there are attributes.
And given Functional Dependencies. We
 $A \rightarrow B, B \rightarrow C$. hidden $F.D^+$
 $A \rightarrow B$
 $B \rightarrow C$
 $A \rightarrow C$

$$FD_1 \cup FD_2 = FD^+$$



The Union of attributes of R_1 and R_2 should be equal to R .

We need to check the dependency preserve.
Let $R(ABCD)$ with functional Dependency

$$\{A \rightarrow B, B \rightarrow C, C \rightarrow D, D \rightarrow B\}$$

if R is decomposed into $R_1(AB), R_2(BC), R_3(BD)$

$$\text{Trivial} \rightarrow A \rightarrow A = A$$

Dependencies

$R_1(AB)$	$R_2(BC)$	$R_3(BD)$	$FD^+ = BCD$
A determine A $A \rightarrow A$ ✓ $B \rightarrow B$ x $A \rightarrow B$ ✓ $B \rightarrow A$ x $A \rightarrow A = A$ But not of leave trivial. Take non-trivial $FD^+ = BCD$ In this table Non-trivial & closure	$B \rightarrow C$ ✓ $C \rightarrow B$ ✓	$R(BD)$ $B \rightarrow D$ ✓ $D \rightarrow B$ ✓	

The possible preserved

$A \rightarrow B$ ✓

$B \rightarrow C$ ✓

$C \rightarrow B$ ✓

$B \rightarrow D$ ✓

$D \rightarrow B$ ✓

$CT = CDD$

Unique key: The values of the attribute in table is unique, Not updatable and may be null.

Surrogate key: It is unique, Not null and updatable.

An Artificial key which aims to uniquely identify each record is called a surrogate key. These kind of keys are unique because they are created when you won't have any natural primary key.

The usual surrogate key used is a unique sequential number.

A surrogate key in a database is a Unique identifier for either an Entity in the modelled world or an object in the database.
The surrogate key is not derived from application data, unlike a natural key which is derived from application data.

or

A surrogate represents an Entity in the outside world. The surrogate is internally generated by the system but is nevertheless visible to the user of application.

Approaches to generating surrogates include

Example :-

Oracle Sequence, or generated as identifier

Dependency-preserving Decomposition:-

The dependency preservation decomposition is another property of decomposed relations. database schema D in which each functional dependency $X \rightarrow Y$ specified in F either appeared directly in one of the Relation Schemas R_i in the decomposed D or could be inferred from the dependencies that appeared in some R_i .

Decomposition $D = \{R_1, R_2, R_3, \dots, R_m\}$ of R is said to be dependency preserving with respect to F if the union of the projections of F on each R_i , in D is equivalent to F .

In other words, RC Join of R_1, R_2 over X .

The dependencies are preserved because each dependency in F represents a constraint on the database. If decomposition is not dependency-preserving, some dependency is lost in the decomposition.