

UNIT - III

NITISH SHARMA

1

Assignment

- Overview of Physical Storage Media,
- File Organization,
- Indexing and Hashing,
- B+ tree Index Files,
- Query Processing Overview,
- Materialized views,
- Database Tuning.

2

Database

A database is an organized collection of data whose content must be quickly and easily

- Accessed
- Managed
- Updated

A relational database is one whose data are split up into tables, sometimes called **relations**.

3

Normalization

Normalization is a process that improves storage efficiency, data integrity, and scalability of a database design by generating relations that are of higher normal forms.

The **main goal** of Database Normalization is to restructure the logical data model of a database to:

- Reduction of redundant data.
- Ensure data dependencies make sense.
- Reduce the potential for data anomalies.

4

History

Edgar F. Codd first proposed the process of normalization and what came to be known as the **1st normal form** in his paper

“A Relational Model of Data for Large Shared Data Banks”

Edgar F. Codd stated:

“There is, in fact, a very simple elimination procedure which we shall call normalization. Through decomposition non-simple domains are replaced by ‘domains whose elements are atomic (non-decomposable) values.’”

- Edgar F. Codd originally established three normal forms: 1NF, 2NF and 3NF.
- 3NF is widely considered to be sufficient for most applications.

5

Why we need Normalization?

- To reduce the redundant data from database.
- To minimize data loss.
- To reduce the potential of data anomalies.
- To improve storage efficiency, data integrity, and scalability of a database design.

6

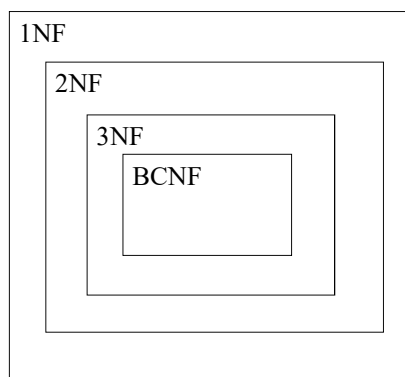
Normalization

Normal Forms are progressive in nature:

1NF is considered the weakest,
2NF is stronger than 1NF,
3NF is stronger than 2NF, and
BCNF is considered the strongest

7

Normalization



a relation in 2NF is also in 1NF

a relation in 3NF is also in 2NF

a relation in BCNF, is also in 3NF

8

Functional Dependency

Consider a relation schema R , and let $\alpha \subseteq R$ and $\beta \subseteq R$.

The functional dependency

$$\alpha \rightarrow \beta$$

holds on schema R if, in any legal relation $r(R)$, for all pairs of tuples t_1 and t_2 in r such that $t_1[\alpha] = t_2[\alpha]$, it is also the case that $t_1[\beta] = t_2[\beta]$.

- **Functional dependencies** sometimes are referred to as “equality-generating dependencies”.

Example:

Suppose each employee is identified by their unique employee number. We say there is a functional dependency of email address on employee number:

$$\text{employee number} \rightarrow \text{email address}$$

9

Determinant

Functional Dependency

$$\text{EmpNum} \rightarrow \text{EmpEmail}$$

Attribute on the LHS is known as the *determinant*

Attribute on the RHS is known as the *determiner*

Here, “EmpNum” is a determinant of “EmpEmail”

10

Functional Dependency

EmpNum	EmpEmail	EmpFname	EmpLname
123	al@npiu.com	Alan	Lee
456	ps@npiu.com	Peter	Smith
555	jd@npiu.com	John	Doe
633	zl@npiu.com	Zhnag	Li
787	xf@npiu.com	Xu	Fing

If EmpNum is the PK then the FDs:

EmpNum \rightarrow EmpEmail

EmpNum \rightarrow EmpFname

EmpNum \rightarrow EmpLname

must exist.

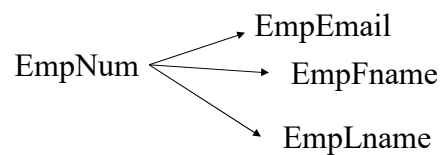
11

Functional Dependency

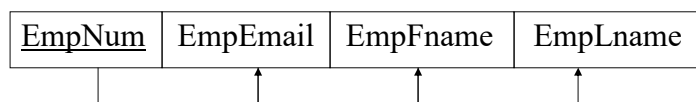
EmpNum \rightarrow EmpEmail

EmpNum \rightarrow EmpFname

EmpNum \rightarrow EmpLname



FDs are depicted in three different ways



12

Introduction to Axioms Rules

- Armstrong's Axioms is a set of rules.
- It provides a simple technique for reasoning about functional dependencies.
- It was developed by William W. Armstrong in 1974.
- It is used to infer all the functional dependencies on a relational database.

13

Various Axioms Rules

A. Primary Rules

Rule 1	Reflexivity If A is a set of attributes and B is a subset of A, then A holds B. $\{A \rightarrow B\}$
Rule 2	Augmentation If A hold B and C is a set of attributes, then AC holds BC. $\{AC \rightarrow BC\}$ It means that attribute in dependencies does not change the basic dependencies.
Rule 3	Transitivity If A holds B and B holds C, then A holds C. If $\{A \rightarrow B\}$ and $\{B \rightarrow C\}$, then $\{A \rightarrow C\}$ A holds B $\{A \rightarrow B\}$ means that A functionally determines B.

14

B. Secondary Rules

Rule 1	Union If A holds B and A holds C, then A holds BC. If $\{A \rightarrow B\}$ and $\{A \rightarrow C\}$, then $\{A \rightarrow BC\}$
Rule 2	Decomposition If A holds BC and A holds B, then A holds C. If $\{A \rightarrow BC\}$ and $\{A \rightarrow B\}$, then $\{A \rightarrow C\}$
Rule 3	Pseudo Transitivity If A holds B and BC holds D, then AC holds D. If $\{A \rightarrow B\}$ and $\{BC \rightarrow D\}$, then $\{AC \rightarrow D\}$

15

Example:

Consider relation $E = (P, Q, R, S, T, U)$ having set of Functional Dependencies (FD).

$P \rightarrow Q$ $P \rightarrow R$
 $QR \rightarrow S$ $Q \rightarrow T$
 $QR \rightarrow U$ $PR \rightarrow U$

Calculate some members of Axioms are as follows,

1. $P \rightarrow T$
2. $PR \rightarrow S$
3. $QR \rightarrow SU$
4. $PR \rightarrow SU$

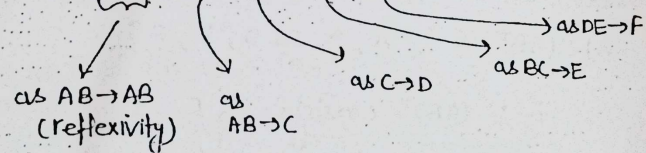
16

➔ Attribute Closure (X^+):-

$X^+ = \{ \text{set of all attributes determined by } X \}$

ex $\{ AB \rightarrow C, C \rightarrow D, DE \rightarrow F, BC \rightarrow E, EG \rightarrow H \}$

$(AB)^+ = \{ A, B, C, D, E, F \}$



17

(ii) $R(ABCDE)$

$\{ A \rightarrow BC, CD \rightarrow E, B \rightarrow D, D \rightarrow A \}$

By Hit and Trial,

$A^+ = \{ A, B, C, D, E \}$ So, A is cand key.

Now

$D \rightarrow A$

So, $D^+ = \{ D, A, B, C, E \}$ So D is cand key

Now,

$B \rightarrow D$

So, $B^+ = \{ B, D, A, C, E \}$ So B is cand key

$\{ A, D, B \}$ are cand key.

18

Equality of FD sets:-

Let F and G be two FD sets:-

then F & G FD sets are equal iff $F^+ \equiv G^+$

But it will become very hard to calculate all members of F^+ and G^+ , so our problem is not solved.

Another method:-

F & G FD sets are equal iff

a) F covers G :-

Every FD of G set must be implied in F .

$$[F \supseteq G]$$

b) G covers F :-

Every FD of F set must be implied in G .

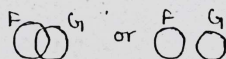
$$[F \subseteq G]$$

19

④

F covers G	G covers F	
✓	✓	$F \equiv G$
✓	✗	$F \supset G$
✗	✓	$F \subset G$
✗	✗	F & G Not Comparable

means



but no one is subset of the other

20

10/

$F = \{ A \rightarrow BCDEF, BC \rightarrow ADEF, B \rightarrow F, D \rightarrow E \}$

$G = \{ A \rightarrow BC, BC \rightarrow AD, B \rightarrow F, D \rightarrow E \}$

then which is true a) $F \subset G$ b) $F \supset G$ c) $F \equiv G$ d) None

In G ,

$A^+ = \{ ABCDEF \}$ so, $A \rightarrow BCDEF$ is covered

$BC^+ = \{ BCADDEF \}$ so $BC \rightarrow ADEF$ is covered

$B \rightarrow F$ & $D \rightarrow E$ is covered.

so, G covers F $F \subset G$

In F , $A \rightarrow BC$ is covered.

$BC \rightarrow AD$ is covered

$B \rightarrow F$ are also covered

$D \rightarrow E$

so, F covers G .

$F \supset G$

$\therefore F \subset G$ & $F \supset G$ so, $F \equiv G$.

(c) is correct.

21

Canonical Cover [Minimal Cover]

Minimal set of FDs (F_m) which are logically equal to given FD set F .

$$F_m \equiv F$$

22

Ex $F = \{ A \rightarrow B, B \rightarrow C, A \rightarrow C, AB \rightarrow C, A \rightarrow A \}$

Find minimal cover of F .

Sol:-

i) $A \rightarrow A$ (Trivial FD)

It should not be included in F .

as if we not have $A \rightarrow A$ then also it can't influence FD set F .

or $A \rightarrow A$ is obviously derivable from F .

So, remove $A \rightarrow A$

ii) $AB \rightarrow C$

$B \rightarrow C \xRightarrow{\text{Augm.}} AB \rightarrow AC \xRightarrow{\text{split}} AB \rightarrow C$

$\therefore AB \rightarrow C$ can be determined by $B \rightarrow C$

So, remove $AB \rightarrow C$

iii) $A \rightarrow C$

$A \rightarrow B, B \rightarrow C \xRightarrow{\text{trans.}} A \rightarrow C$

So remove $A \rightarrow C$.

$\therefore F_{\text{MC}} = \{ A \rightarrow B, B \rightarrow C \}$

23

NOTE:-

Minimal cover of FD set F may not be unique but all minimal covers are logically equal or equal to F .

i.e.

$$(F_{M1} \equiv F_{M2} \equiv F)$$

or Expressive Power of $F_{M1} = \text{that of } F_{M2} = \text{that of } F$.

24

Functional Dependency

CASE I:

The functional dependency $\alpha \rightarrow \beta$ forms redundancy in R iff

- It is not a trivial F.D. & α is not a superkey.

CASE II:

The functional dependency $\alpha \rightarrow \beta$ doesn't forms redundancy in R iff

- It is a trivial F.D. or α is a superkey.

25

What is decomposition?

- Decomposition is the process of breaking down in parts or elements.
- It replaces a relation with a collection of smaller relations.
- It breaks the table into multiple tables in a database.
- It should always be lossless, because it confirms that the information in the original relation can be accurately reconstructed based on the decomposed relations.
- If there is no proper decomposition of the relation, then it may lead to problems like loss of information.

26

Properties of Decomposition

Following are the properties of Decomposition,

1. Lossless Decomposition
2. Dependency Preservation

27

1. Lossless Decomposition

- Decomposition must be lossless. It means that the information should not get lost from the relation that is decomposed.
- It gives a guarantee that the join will result in the same relation as it was decomposed.

Lossless Join Decomposition

Relation R decomposed into $R_1 R_2 \dots R_n$

In General, $(R_1 \bowtie R_2 \bowtie R_3 \bowtie \dots \bowtie R_n) \supseteq R$

If $(R_1 \bowtie R_2 \bowtie \dots \bowtie R_n) \equiv R$
Then Lossless Join Decomposition

If $(R_1 \bowtie R_2 \bowtie \dots \bowtie R_n) \supset R$
Then lossy join decomposition

28

Example: <Employee_Department> Table

Eid	Ename	Age	City	Salary	Deptid	DeptName
E001	ABC	29	Pune	20000	D001	Finance
E002	PQR	30	Pune	30000	D002	Production
E003	LMN	25	Mumbai	5000	D003	Sales
E004	XYZ	24	Mumbai	4000	D004	Marketing
E005	STU	32	Bangalore	25000	D005	Human Resource

- Decompose the above relation into two relations to check whether a decomposition is lossless or lossy.
- Now, we have decomposed the relation that is Employee and Department.

Relation 1 : <Employee> Table

Eid	Ename	Age	City	Salary
E001	ABC	29	Pune	20000
E002	PQR	30	Pune	30000
E003	LMN	25	Mumbai	5000
E004	XYZ	24	Mumbai	4000
E005	STU	32	Bangalore	25000

- Employee Schema contains (Eid, Ename, Age, City, Salary).

29

Relation 2 : <Department> Table

Deptid	Eid	DeptName
D001	E001	Finance
D002	E002	Production
D003	E003	Sales
D004	E004	Marketing
D005	E005	Human Resource

- Department Schema contains (Deptid, Eid, DeptName).
- So, the above decomposition is a Lossless Join Decomposition, because the two relations contains one common field that is 'Eid' and therefore join is possible.
- Now apply natural join on the decomposed relations.

30

Employee » Department

Eid	Ename	Age	City	Salary	Deptid	DeptName
E001	ABC	29	Pune	20000	D001	Finance
E002	PQR	30	Pune	30000	D002	Production
E003	LMN	25	Mumbai	5000	D003	Sales
E004	XYZ	24	Mumbai	4000	D004	Marketing
E005	STU	32	Bangalore	25000	D005	Human Resource

Hence, the decomposition is Lossless Join Decomposition.

- If the <Employee> table contains (Eid, Ename, Age, City, Salary) and <Department> table contains (Deptid and DeptName), then it is not possible to join the two tables or relations, because there is no common column between them. And it becomes **Lossy Join Decomposition**.

31

2. Dependency Preservation

- Dependency is an important constraint on the database.
- Every dependency must be satisfied by at least one decomposed table.
- If $\{A \rightarrow B\}$ holds, then two sets are functional dependent. And, it becomes more useful for checking the dependency easily if both sets in a same relation.
- This decomposition property can only be done by maintaining the functional dependency.
- In this property, it allows to check the updates without computing the natural join of the database structure.

... Relational Schema (R) with FDset F decomposed into R_1, R_2, \dots, R_n subRelations with FDsets F_1, F_2, \dots, F_n respectively.

In General

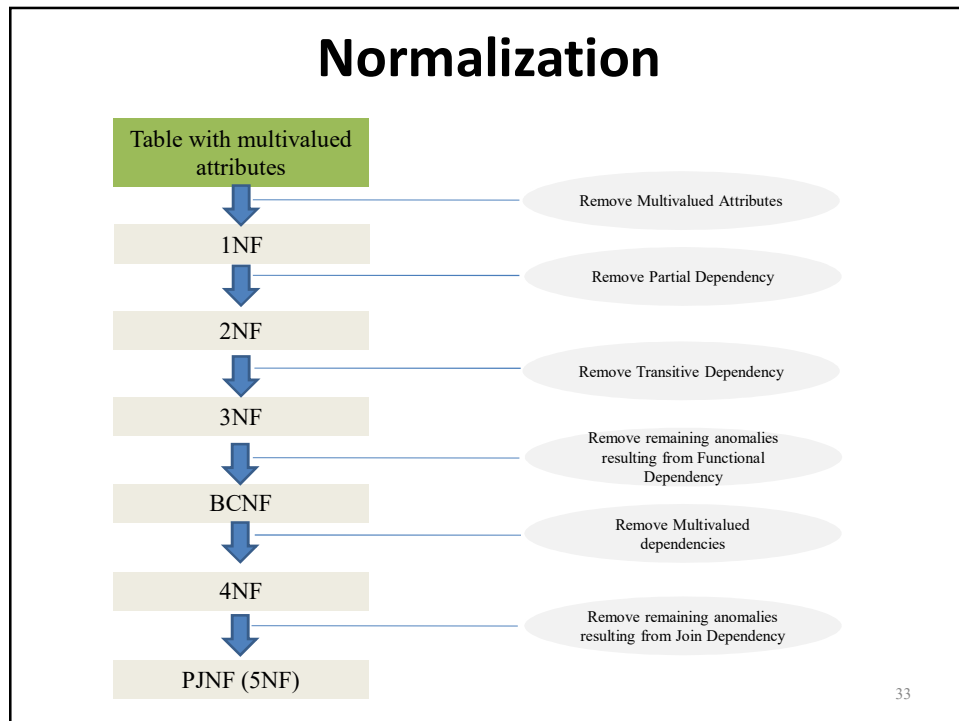
$$[F_1 \cup F_2 \cup \dots \cup F_n] \subseteq F$$

If $[F_1 \cup F_2 \cup \dots \cup F_n] \equiv F$
 then decomposition is Dependency Preserving Decomposition.

If $[F_1 \cup F_2 \cup \dots \cup F_n] \subset F$
 then it is not Dependency Preserving Decomposition.

} No FD is lost becoz of Decomposition

32



First Normal Form

A relation schema R is in **1NF** if the domains of all attributes of R are atomic.

“1NF places restrictions on the structure of relation”

- It does not require additional information such as functional dependencies.

A domain is **atomic** if elements of the domain are considered to be indivisible units.

First Normal Form

The following relational table is **not** in 1NF:

EmpNum	EmpPhone	EmpDegrees
123	233-9876	
333	233-1231	BA, BSc, PhD
679	233-1279	BSc, MSc

EmpDegrees is a multi-valued field:

Employee with EmpNum - 333 has three degrees: *BA, BSc, PhD*

Employee with EmpNum - 679 has two degrees: *BSc* and *MSc*

35

First Normal Form

Employee

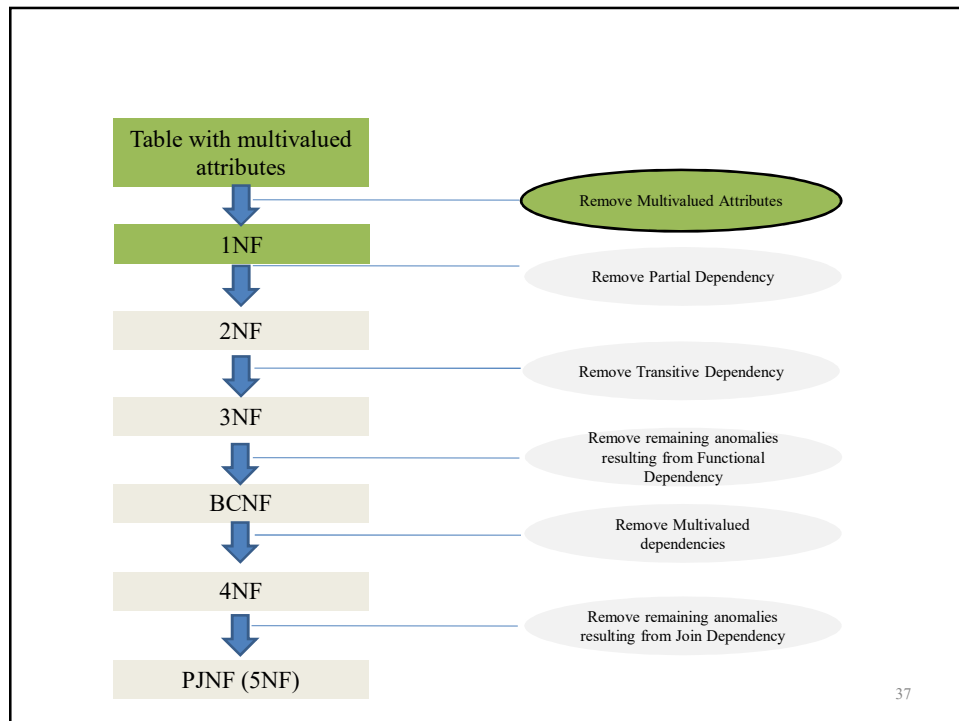
EmpNum	EmpPhone
123	233-9876
333	233-1231
679	233-1231

EmployeeDegree

EmpNum	EmpDegree
333	BA
333	BSc
333	PhD
679	BSc
679	MSc

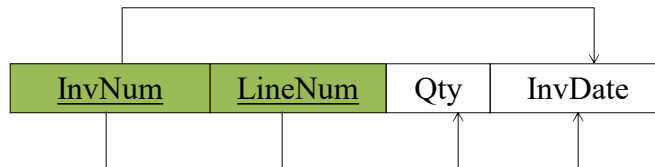
We can obtain original relation by using Join operation on these relations.

36



Partial dependency

A **partial dependency** exists when an attribute B is functionally dependent on an attribute A, and A is a component of a multipart candidate key.



Candidate keys: {InvNum, LineNum}

InvDate is *partially dependent* on {InvNum, LineNum} as **InvNum is a determinant of InvDate and InvNum is part of a candidate key**

Second Normal Form

A relation is in **2NF** if it is in 1NF, and every non-key attribute is fully functional dependent on candidate key.

- It is based on Full Functional dependency.
- 2NF (and 3NF) both involve the concepts of key and non-key attributes.

39

Second Normal Form

Consider this **InvLine** relation (in 1NF):

<u>InvNum</u>	<u>LineNum</u>	ProdNum	Qty	InvDate
---------------	----------------	---------	-----	---------

$\text{InvNum, LineNum} \longrightarrow \text{ProdNum, Qty}$

$\text{InvNum} \longrightarrow \text{InvDate}$

Key Attributes: InvNum, LineNum

InvLine is **not in 2NF** due to the presence of partial dependency of InvDate on InvNum

40

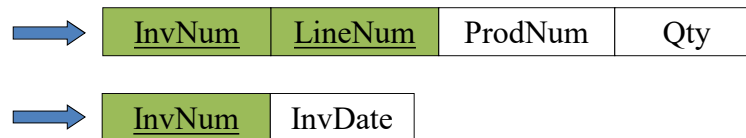
Second Normal Form

InvLine

<u>InvNum</u>	<u>LineNum</u>	ProdNum	Qty	InvDate
---------------	----------------	---------	-----	---------

The above relation has redundancies: the invoice date is repeated on each invoice line.

We can *improve* the database by decomposing the relation into two relations:



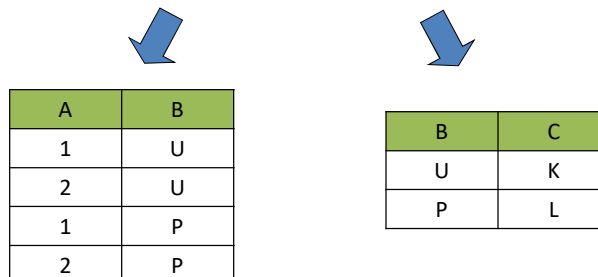
41

R(A,B,C)

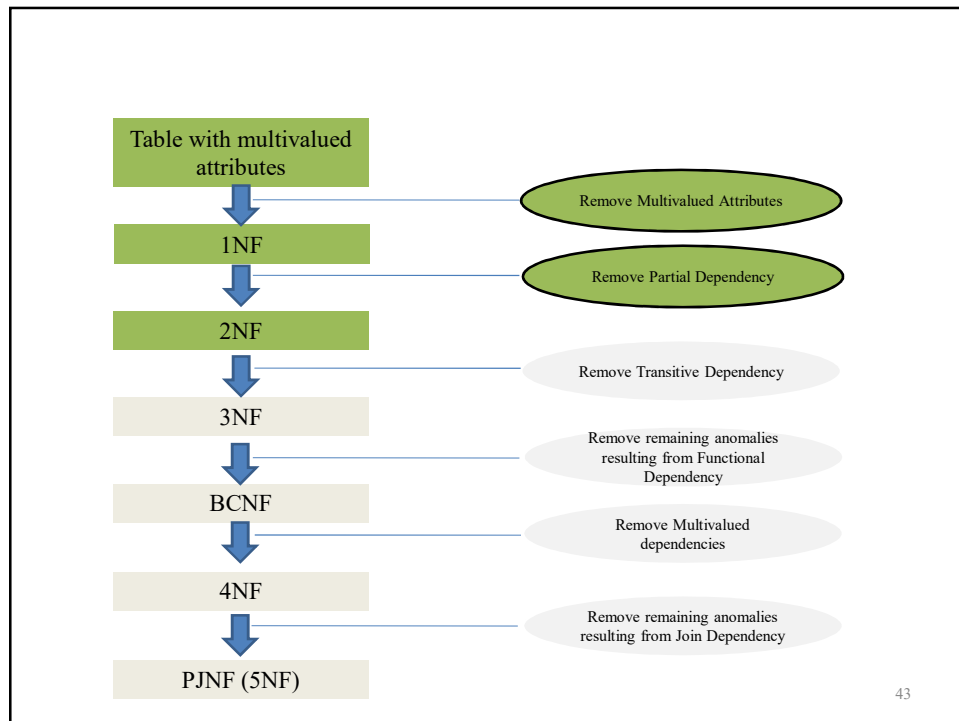
$AB \rightarrow C$
 $B \rightarrow C$

A	B	C
1	U	K
2	U	K
1	P	L
2	P	L

Here, repetition of values for C causes redundancy



42



Transitive dependency

Consider attributes A, B, and C, and where

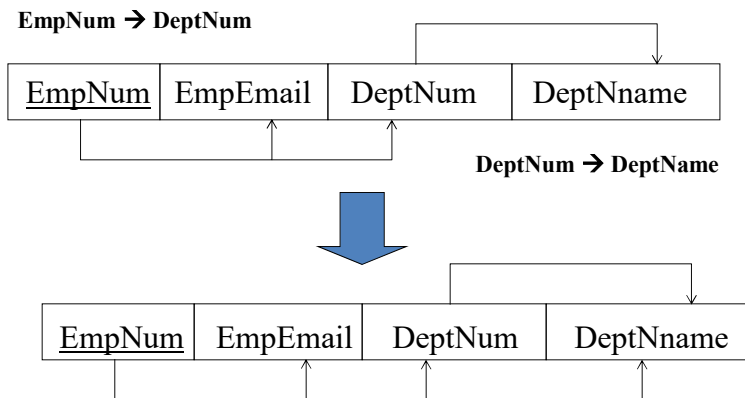
$$A \rightarrow B \text{ and } B \rightarrow C.$$

Functional dependencies are transitive, which means that we also have the functional dependency

$$A \rightarrow C$$

We say that C is transitively dependent on A through B.

Transitive dependency



DeptName is *transitively dependent* on EmpNum via DeptNum
 $\text{EmpNum} \rightarrow \text{DeptName}$

45

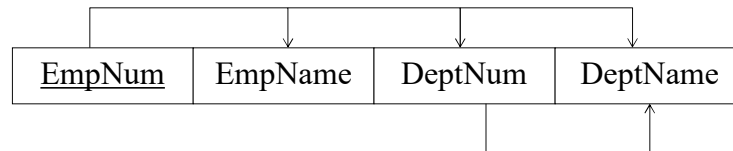
Third Normal Form

A relation schema R is in **3NF** with respect to a set F of functional dependencies if, for all functional dependencies in F^+ of the form $\alpha \rightarrow \beta$, where $\alpha \subseteq R$ and $\beta \subseteq R$, at least one of the following holds:

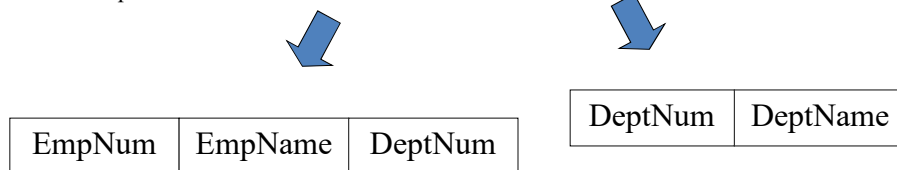
- $\alpha \rightarrow \beta$ is trivial functional dependency (that is, $\beta \subseteq \alpha$).
- α is a superkey for schema R .
- Each attribute A in $\beta - \alpha$ is contained in a candidate key for R .

46

Third Normal Form



We correct the situation by decomposing the original relation into two 3NF relations. Note the decomposition is *lossless*.



47

$R(A,B,C)$

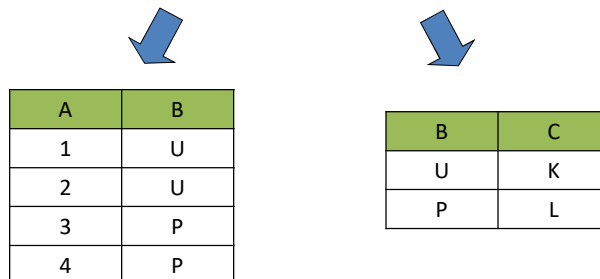
$A \rightarrow B$

$B \rightarrow C$

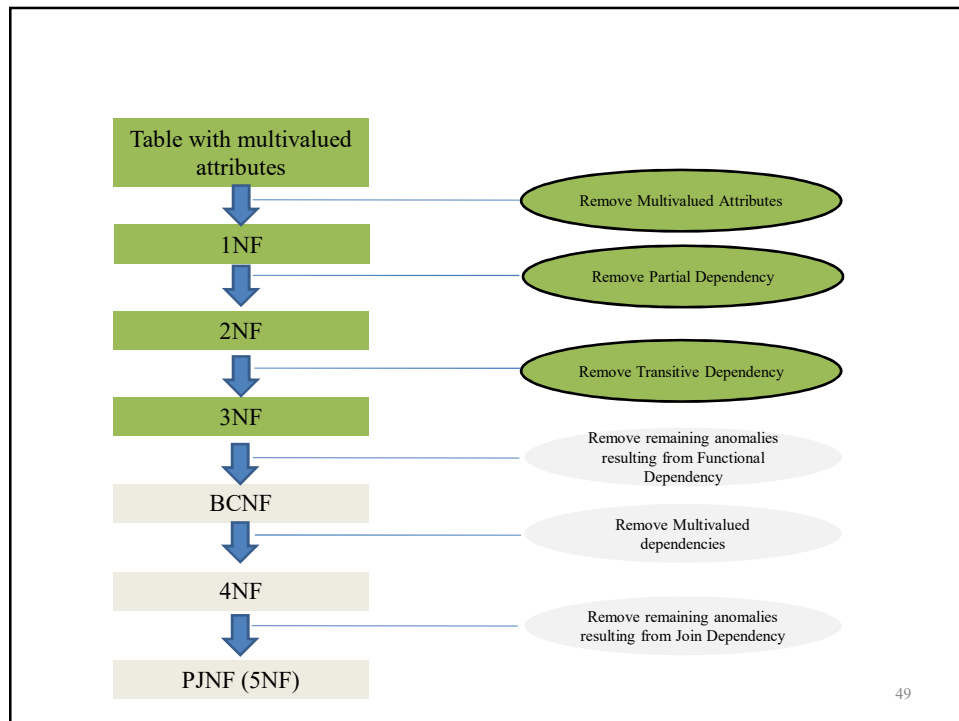
$A \rightarrow C$

A	B	C
1	U	K
2	U	K
3	P	L
4	P	L

Here, repetition of values for C causes redundancy



48

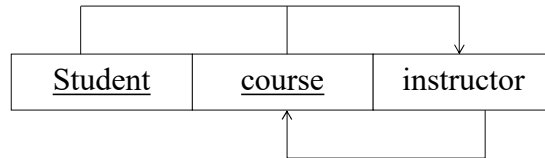


Boyce-Codd Normal Form

A relation schema R is in BCNF with respect to a set F of functional dependencies if, for all functional dependencies in F^+ of the form $\alpha \rightarrow \beta$, where $\alpha \subseteq R$ and $\beta \subseteq R$, at least one of the following holds:

- $\alpha \rightarrow \beta$ is trivial functional dependency (that is, $\beta \subseteq \alpha$).
- α is a superkey for schema R .

Boyce-Codd Normal Form



$\{Student, course\} \rightarrow Instructor$
 $Instructor \rightarrow Course$

Decomposing into 2 schemas

- $\{Student, Instructor\} \{Student, Course\}$
- $\{Course, Instructor\} \{Student, Course\}$
- $\{Course, Instructor\} \{Instructor, Student\}$

51

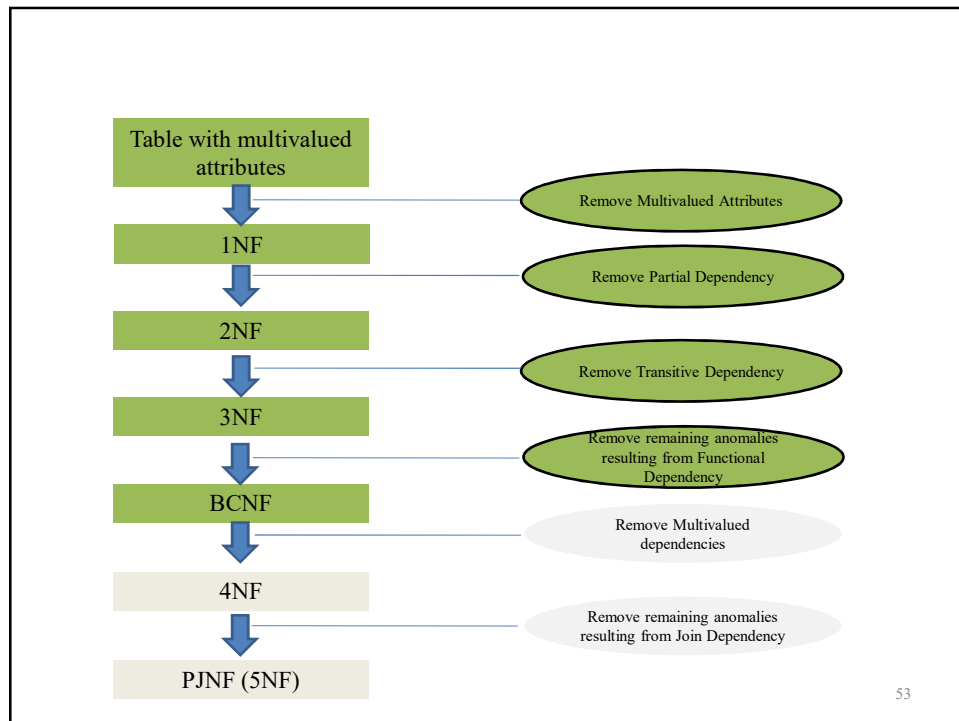
Comparison of BCNF and 3NF

BCNF requires that all non-trivial dependencies be of the form $\alpha \rightarrow \beta$, where α is a superkey.

3NF relaxes this constraint slightly by allowing non-trivial functional dependencies whose determinant is not a superkey.

- It is possible to obtain a 3NF design without sacrificing a lossless join or dependency preservation.
- 3NF allows certain functional dependencies that are not allowed in BCNF.
- Unlike BCNF, 3NF decompositions may contain some redundancy in the decomposed schema.

52



Fourth Normal Form (4NF)

A relation schema R is in **4NF** with respect to a set D of functional and multivalued dependencies if, for all multivalued dependencies in D^+ of the form $\alpha \twoheadrightarrow \beta$, where $\alpha \subseteq R$ and $\beta \subseteq R$, at least one of the following holds:

- $\alpha \twoheadrightarrow \beta$ is trivial multivalued dependency (if $\beta \subseteq \alpha$ or $\beta \cup \alpha = R$).
 - α is a superkey for schema R .
- Every 4NF schema is in BCNF.
 - Multivalued dependencies sometimes are referred to as “tuple-generating dependencies”.

Fourth Normal Form (4NF)

Consider this **Movie** relation:

MovieName	ScreeningCity	Genre
-----------	---------------	-------

Candidate Key: {MovieName, ScreeningCity, Genre}

1. All columns are a part of the only candidate key, hence BCNF
2. Many Movies can have the same Genre
3. Many Cities can have the same movie
4. Violates 4NF

MovieName	ScreeningCity	Genre
Hard Code	Los Angeles	Comedy
Hard Code	New York	Comedy
Bill Durham	Santa Cruz	Drama
Bill Durham	Durham	Drama
The Code Warrior	New York	Horror

55

Fourth Normal Form (4NF)

We correct the situation by decomposing the original relation into two relations.

1. Move the two multi-valued relations to separate tables
2. Identify a primary key for each of the new entity.

MovieName	ScreeningCity	Genre
-----------	---------------	-------



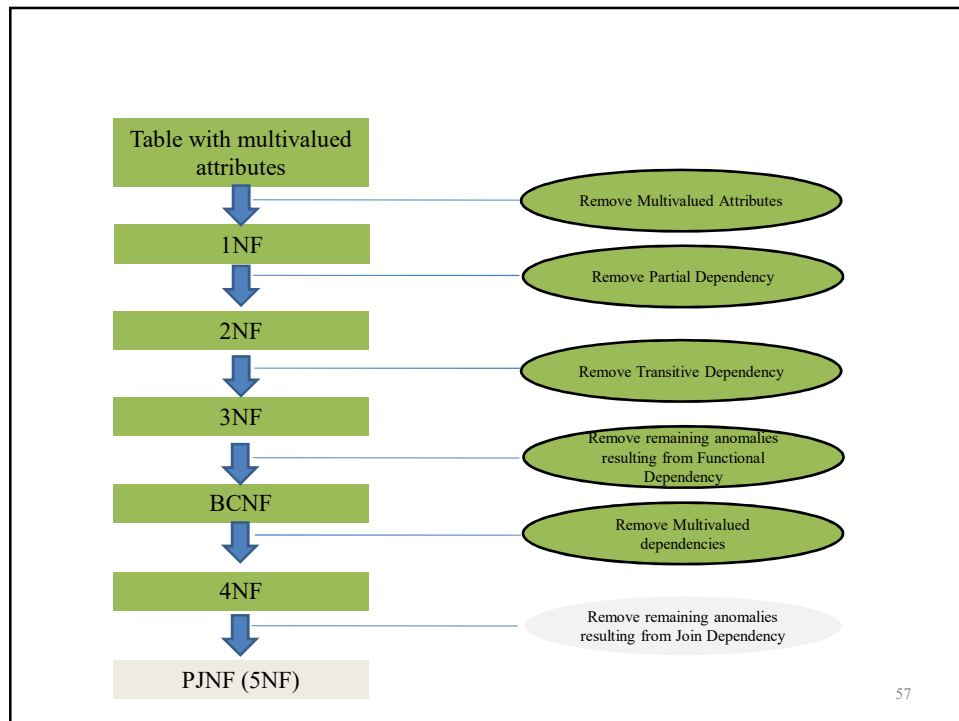
MovieName	ScreeningCity
-----------	---------------

Movie	ScreeningCity
Hard Code	Los Angeles
Hard Code	New York
Bill Durham	Santa Cruz
Bill Durham	Durham
The Code Warrior	New York

MovieName	Genre
-----------	-------

Movie	Genre
Hard Code	Comedy
Bill Durham	Drama
The Code Warrior	Horror

56



More Normal Forms

The **4NF** is “ultimate” normal form.

because, multivalued dependency helps us to tackle some forms of repetition of information that cannot be understood in terms of functional dependencies.

There are types of constraints called join dependencies that generalize multivalued dependencies, and lead to another normal form called “**Project-Join Normal Form**”

There is a class of even more general constraints, which leads to a normal form called “**Domain-Key Normal Form**”.

Analysis of Normal Forms for Redundancy

$X \rightarrow Y$ Non-trivial F.D. & X not a Superkey (forms redundancy)	1NF	2NF	3NF	BCNF
[Proper subset of C.K] \rightarrow [Non prime attribute]	Allowed	Not allowed	Not allowed	Not allowed
[Non prime attribute] \rightarrow [Non prime attribute]	Allowed	Allowed	Not allowed	Not allowed
[Proper subset of C.K & Non prime attribute] \rightarrow [Non prime attribute]	Allowed	Allowed	Not allowed	Not allowed
[Proper subset of C.K] \rightarrow [Proper subset of other C.K]	Allowed	Allowed	Allowed	Not allowed

59

Database Design Goals

Goals	1NF	2NF	3NF	BCNF	4NF
0% Redundancy	No	No	No	Yes: over F.D's No: over M.D's	Yes: over F.D's & M.D's
Lossless join decomposition satisfy	Yes	Yes	Yes	Yes	May not
Dependency Preservation satisfy	Yes	Yes	Yes	May not	May not

60

Q. R(ABCD) $F = \{ AB \rightarrow C, C \rightarrow A, AC \rightarrow D \}$
 which is the Highest NF of R?

Sol:- Candidate keys: $\{ AB, BC \}$

Testing for BCNF (at least one FD with determinant not SK)

SK $(AB) \rightarrow C$ ✓
 Not SK $(C) \rightarrow A$ X
 Not SK $(AC) \rightarrow D$ X } so not BCNF

Testing for 3NF

SK $(AB) \rightarrow C$ Prime ✓
 $C \rightarrow A$ Prime ✓
 Not SK $(AC) \rightarrow D$ Not Prime X } so not 3NF

Testing for 2NF

Here we take all proper subset of candidate keys and find their closure & check whether they contain all prime attributes or not.

Candidate key $\Rightarrow AB \Rightarrow A^+ = (A)$ Prime ✓
 $B^+ = (B)$ Prime ✓
 $AC \Rightarrow C^+ = (CA)$ Not Prime X
 $C \rightarrow D$ will be a partial dependency in R.
 \therefore not in 2NF

Hence highest NF of R: 1NF

61

Q. R(ABCD)

$\{ AB \rightarrow C, BC \rightarrow D \}$

Candidate keys: $\{ AB \}$

BCNF:

$AB \rightarrow C$ ✓
 Not SK $(BC) \rightarrow D$ X } so not in BCNF

3NF:

$BC \rightarrow D$ Not Prime } so not in 3NF

2NF:

$AB \rightarrow A^+ = (A)$ Prime ✓
 $B^+ = (B)$ Prime ✓ } so, 2NF

62

Q R(ABCDE)

$\{A \rightarrow BC, CD \rightarrow E, B \rightarrow D, E \rightarrow A\}$

Cand Key: $\{\underline{A}, \underline{E}, \underline{CD}, \underline{BC}\}$

BCNF:

$A \rightarrow BC$
 $CD \rightarrow E$
~~not~~ $\textcircled{B} \rightarrow D$ so not in BCNF
 $E \rightarrow A$

3NF:

$B \rightarrow \textcircled{D} \text{ Prime}$
 $A \rightarrow \textcircled{BC} \text{ Prime}$
 $CD \rightarrow \textcircled{E} \text{ Prime}$
 $E \rightarrow \textcircled{A}$ so 3NF