

What is Normalization?

- Database designed based on the E-R model may have some amount of
 - Inconsistency
 - Uncertainty
 - Redundancy

To eliminate these draw backs some **refinement** has to be done on the database.

- **Refinement** process is called **Normalization**
- Defined as a step-by-step process of decomposing a complex relation into a simple and stable data structure.
- The formal process that can be followed to achieve a good database design
- Also used to check that an existing design is of good quality
- The different stages of normalization are known as “normal forms”
- To accomplish normalization we need to understand the concept of Functional Dependencies.

Source: Infosys Campus Connect Study Material

Need for Normalization

Student_Course_Result Table

Student_Details			Course_Details				Result_Details		
101	Davis	11/4/1986	M4	Applied Mathematics	Basic Mathematics	7	11/11/2004	82	A
102	Daniel	11/6/1987	M4	Applied Mathematics	Basic Mathematics	7	11/11/2004	62	C
101	Davis	11/4/1986	H6	American History		4	11/22/2004	79	B
103	Sandra	10/2/1988	C3	Bio Chemistry	Basic Chemistry	11	11/16/2004	65	B
104	Evelyn	2/22/1986	B3	Botany		8	11/26/2004	77	B
102	Daniel	11/6/1987	P3	Nuclear Physics	Basic Physics	13	11/12/2004	68	B
105	Susan	8/31/1985	P3	Nuclear Physics	Basic Physics	13	11/12/2004	89	A
103	Sandra	10/2/1988	B4	Zoology		5	11/27/2004	54	D
105	Susan	8/31/1985	H6	American History		4	11/22/2004	87	A
104	Evelyn	2/22/1986	M4	Applied Mathematics	Basic Mathematics	7	11/11/2004	65	B

- Data Duplication
- Delete Anomaly

- Insert Anomaly
- Update Anomaly

Source: Infosys Campus Connect Study Material

Need for Normalization

- **Duplication of Data** – The same data is listed in multiple lines of the database
- **Insert Anomaly** – A record about an entity cannot be inserted into the table without first inserting information about another entity – Cannot enter a student details without a course details
- **Delete Anomaly** – A record cannot be deleted without deleting a record about a related entity. Cannot delete a course details without deleting all of the students' information.
- **Update Anomaly** – Cannot update information without changing information in many places. To update student information, it must be updated for each course the student has placed

Functional dependency

- In a given relation R, X and Y are attributes. Attribute Y is **functionally dependent** on attribute X if each value of X determines **EXACTLY ONE** value of Y, which is represented as $X \rightarrow Y$ (X can be composite in nature).
- We say here “x determines y” or “y is functionally dependent on x”
 $X \rightarrow Y$ does not imply $Y \rightarrow X$
- If the value of an attribute “Marks” is known then the value of an attribute “Grade” is determined since $\text{Marks} \rightarrow \text{Grade}$
- Types of functional dependencies:
 - Full Functional dependency
 - Partial Functional dependency
 - Transitive dependency

Functional dependency

Consider the following Relation

REPORT (STUDENT#, COURSE#, CourseName, IName, Room#, Marks, Grade)

- **STUDENT#** - Student Number
- **COURSE#** - Course Number
- **CourseName** - Course Name
- **IName** - Name of the Instructor who delivered the course
- **Room#** - Room number which is assigned to respective Instructor
- **Marks** - Scored in Course COURSE# by Student STUDENT#
- **Grade** - obtained by Student STUDENT# in Course COURSE#

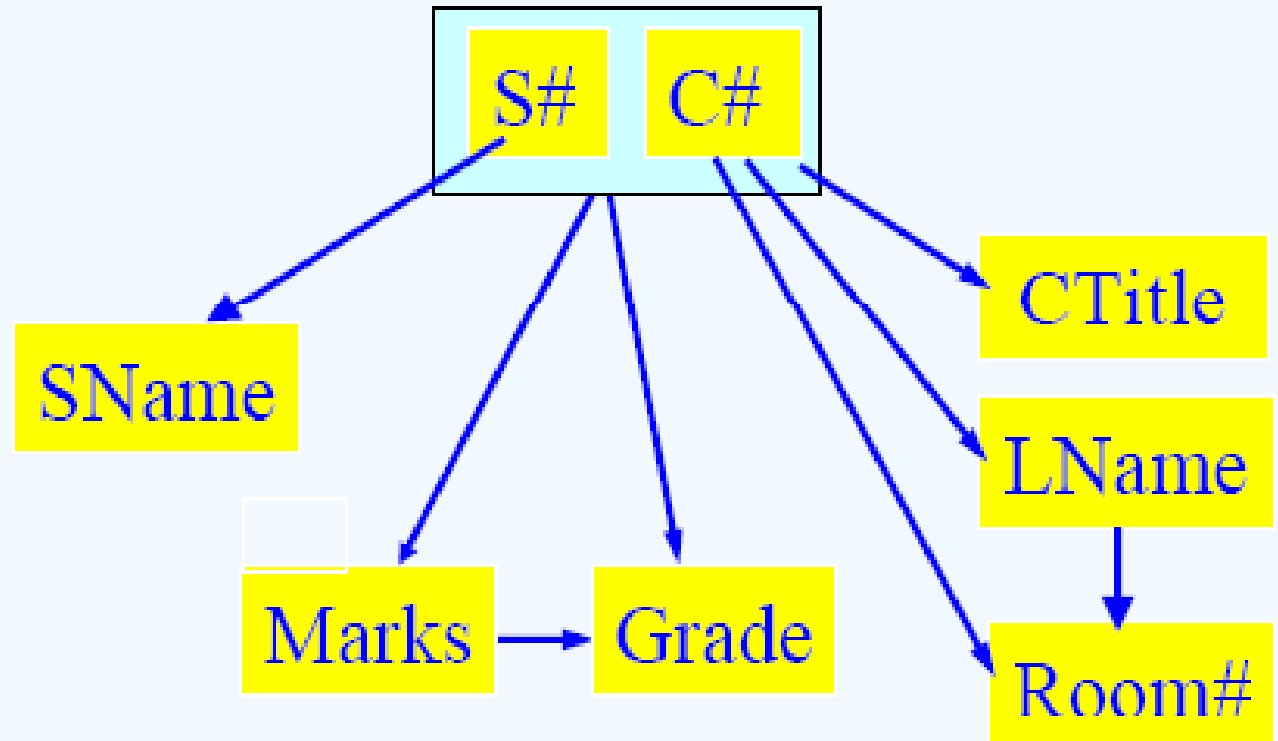
Functional dependency

- **STUDENT# COURSE# → Marks**
- **COURSE# → CourseName,**
- **COURSE# → IName** (Assuming one course is taught by one and only one Instructor)
- **IName → Room#** (Assuming each Instructor has his/her own and non-shared room)
- **Marks → Grade**

Dependency Diagram

Report(S#, C#, SName, CTitle, LName, Room#, Marks, Grade)

- $S\# \rightarrow SName$
- $C\# \rightarrow CTitle,$
- $C\# \rightarrow LName$
- $LName \rightarrow Room\#$
- $C\# \rightarrow Room\#$
- $S\# C\# \rightarrow Marks$
- $Marks \rightarrow Grade$
- $S\# C\# \rightarrow Grade$



Assumptions:

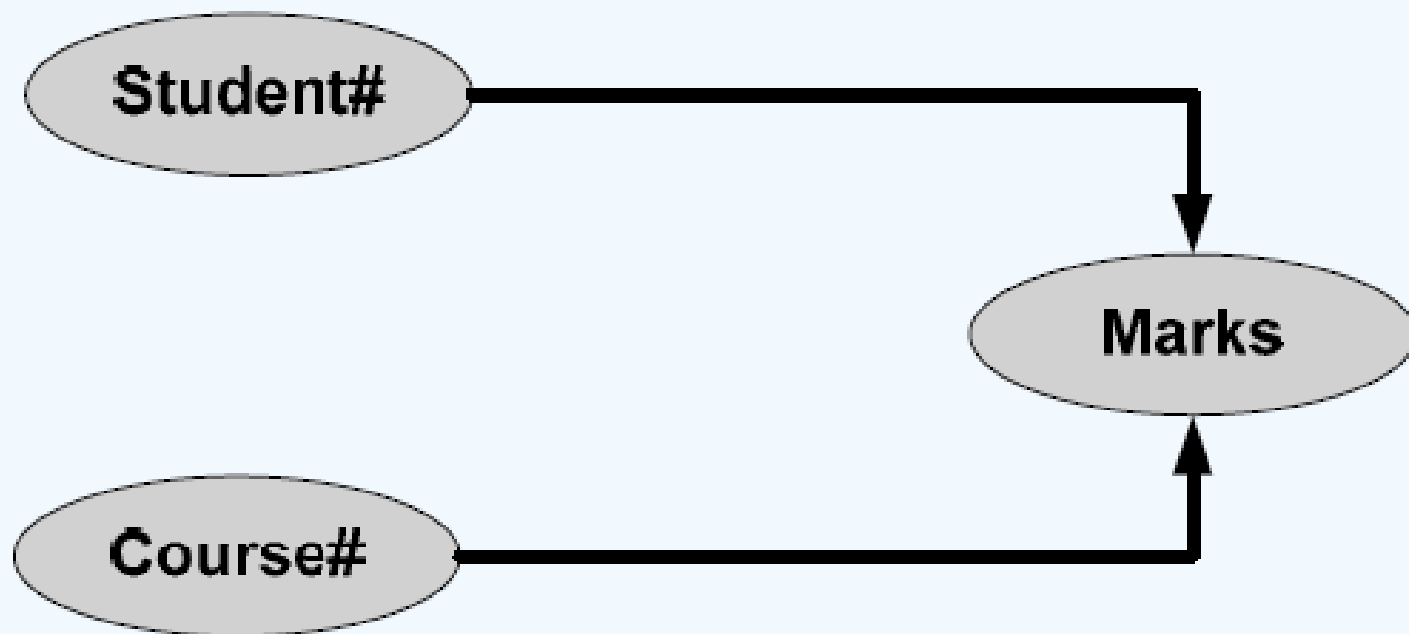
- Each course has only one lecturer and each lecturer has a room.
- Grade is determined from Marks.

Full Dependency

X and Y are attributes.

X Functionally determines Y

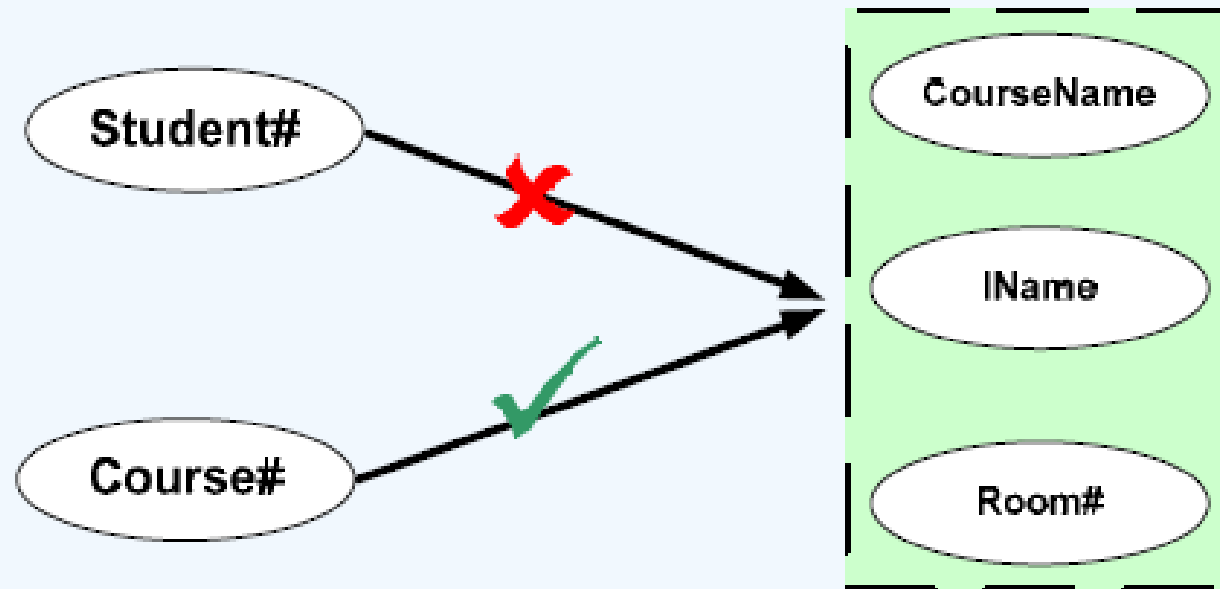
Note: Subset of X should not functionally determine Y



Partial Dependency

X and Y are attributes.

Attribute Y is partially dependent on the attribute X only if it is dependent on a sub-set of attribute X.



Source: Infosys Campus Connect Study Material

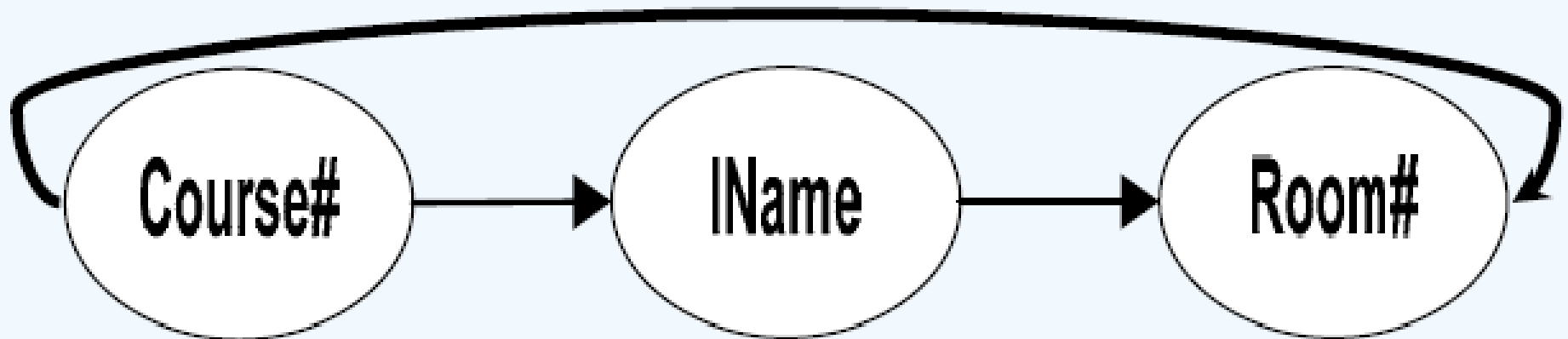
Transitive Dependency

X Y and Z are three attributes.

$X \rightarrow Y$

$Y \rightarrow Z$

$\Rightarrow X \rightarrow Z$



First Normal Form

- Domain is **atomic** if its elements are considered to be **indivisible units**
 - Examples of non-atomic domains:
 - ▶ Set of names, composite attributes
 - ▶ Identification numbers like CS101 that can be broken up into parts
- A relational schema R is in **first normal form** if the domains of all attributes of R are atomic
- Non-atomic values complicate storage and encourage redundant (repeated) storage of data

First Normal Form (Cont'd)

- A relation schema is in 1NF :
 - if and only if all the attributes of the relation R are atomic in nature.
 - **Atomic:** the smallest level to which data may be broken down and remain meaningful

Example ... Without Normalization

Student_Course_Result Table

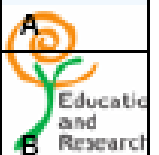
Student_Details			Course_Details				Result_Details		
101	Davis	11/4/1986	M4	Applied Mathematics	Basic Mathematics	7	11/11/2004	82	A
102	Daniel	11/6/1987	M4	Applied Mathematics	Basic Mathematics	7	11/11/2004	62	C
101	Davis	11/4/1986	H6	American History		4	11/22/2004	79	B
103	Sandra	10/2/1988	C3	Bio Chemistry	Basic Chemistry	11	11/16/2004	65	B
104	Evelyn	2/22/1986	B3	Botany		8	11/26/2004	77	B
102	Daniel	11/6/1987	P3	Nuclear Physics	Basic Physics	13	11/12/2004	68	B
105	Susan	8/31/1985	P3	Nuclear Physics	Basic Physics	13	11/12/2004	89	A
103	Sandra	10/2/1988	B4	Zoology		5	11/27/2004	54	D
105	Susan	8/31/1985	H6	American History		4	11/22/2004	87	A
104	Evelyn	2/22/1986	M4	Applied Mathematics	Basic Mathematics	7	11/11/2004	65	B

Source: Infosys Campus Connect Study Material

Table in 1NF

Student_Course_Result Table

Student#	Student Name	Dateof Birth	Cour se #	CourseName	Pre Requisite	Dura tion InDa ys	DateOf Exam	Marks	Grade
101	Davis	04-Nov-1986	M4	Applied Mathematics	Basic Mathematics	7	11-Nov-2004	82	A
102	Daniel	06-Nov-1986	M4	Applied Mathematics	Basic Mathematics	7	11-Nov-2004	62	C
101	Davis	04-Nov-1986	H6	American History		4	22-Nov-2004	79	B
103	Sandra	02-Oct-1988	C3	Bio Chemistry	Basic Chemistry	11	16-Nov-2004	65	B
104	Evelyn	22-Feb-1986	B3	Botany		8	26-Nov-2004	77	B
102	Daniel	06-Nov-1986	P3	Nuclear Physics	Basic Physics	13	12-Nov-2004	68	B
105	Susan	31-Aug-1985	P3	Nuclear Physics	Basic Physics	13	12-Nov-2004	89	A
103	Sandra	02-Oct-1988	B4	Zoology		5	27-Nov-2004	54	D
105	Susan	31-Aug-1985	H6	American History		4	22-Nov-2004	87	A
104	Evelyn	22-Feb-1986	M4	Applied Mathematics	Basic Mathematics	7	11-Nov-2004	65	B



First Normal Form Example

Course_Pref_Table			
Dept	Prof	Course Pref	
		Course	Course_dept
CE	Rajiv	101	CS
		102	CS
		103	EC
	Mahesh	101	CS
		102	CS
		103	EC
		104	EC
CL	Ruchika	101	CS
		103	EC
		106	EE
IT	Rajesh	103	EC
		104	EC
		106	EE
		102	CS
		105	EE

First Normal Form Example

Course_Pref_Table			
Dept	Prof	Course	Course_dept
CE	Rajiv	101	CS
CE	Rajiv	102	CS
CE	Rajiv	103	EC
CE	Mahesh	101	CS
CE	Mahesh	102	CS
CE	Mahesh	103	EC
CE	Mahesh	104	EC
CL	Ruchika	101	CS
CL	Ruchika	103	EC
CL	Ruchika	106	EE
IT	Rajesh	103	EC
IT	Rajesh	104	EC
IT	Rajesh	106	EE
IT	Rajesh	102	CS
IT	Rajesh	105	EE

Second normal form: 2NF

- *A Relation is said to be in Second Normal Form if and only if :*
 - *It is in the First normal form, and*
 - *No partial dependency exists between non-key attributes and key attributes.*
- An attribute of a relation R that belongs to any key of R is said to be a prime attribute and that which doesn't is a **non-prime attribute**

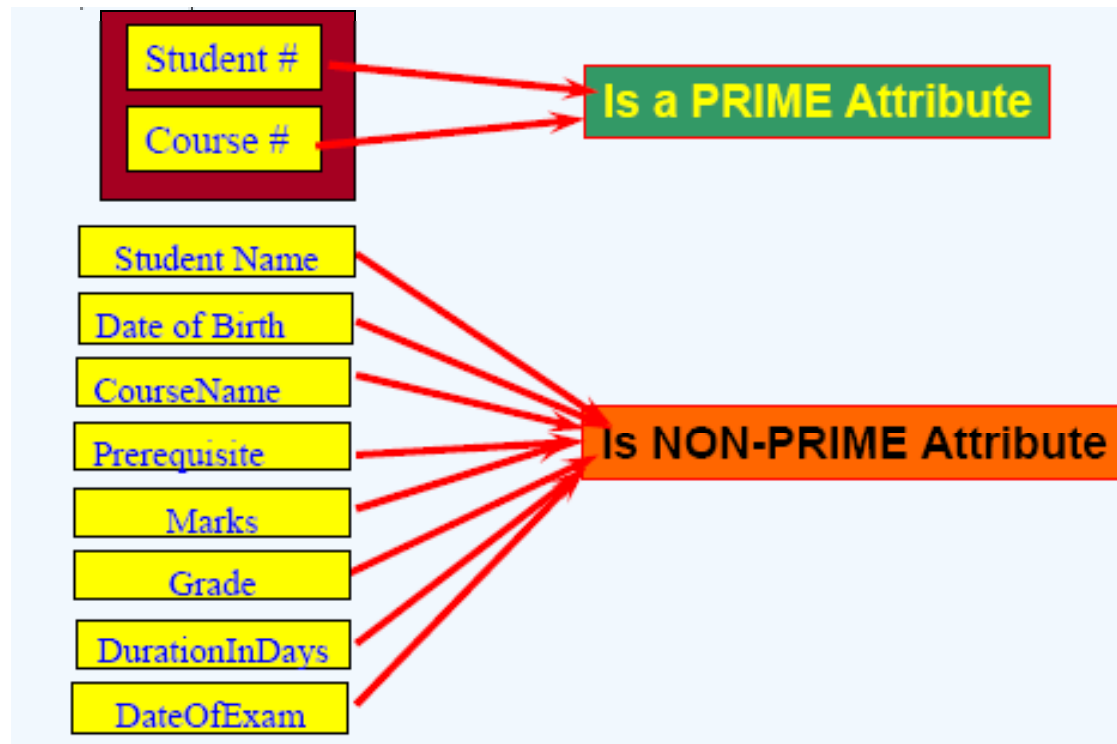
To make a table 2NF compliant, we have to remove all the partial dependencies

Note : - All partial dependencies are eliminated

Prime Vs Non-Prime Attributes

- An attribute of a relation R that belongs to any key of R is said to be a **prime** attribute and that which doesn't is a **non-prime attribute**

Report(S#, C#, StudentName, DateOfBirth, CourseName, PreRequisite, DurationInDays, DateOfExam, Marks, Grade)



Source: Infosys Campus Connect Study Material

Second normal form: 2NF

- STUDENT# is key attribute for Student,
- COURSE# is key attribute for Course
- STUDENT# COURSE# together form the composite key attributes for Results relationship.
- Other attributes like StudentName (Student Name), DateofBirth, CourseName, PreRequisite, DurationInDays, DateofExam, Marks and Grade are non-key attributes.

To make this table 2NF compliant, we have to remove all the partial dependencies.

Student #, Course# -> Marks, Grade

Student# -> StudentName, DOB,

Course# -> CourseName, Prerequisite, DurationInDays

Course# -> Date of Exam

Second normal form: 2NF

S#,C#



Marks

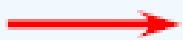
S#,C#



Grade

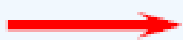
Fully Functionally
dependent on composite
Candidate key

S#



StudentName

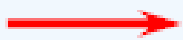
S#



DOB

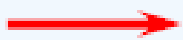
Partial Dependency

C#



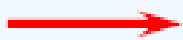
CourseName

C#



Prerequisite

C#



Duration

Partial Dependency

C#



DateOfExam

Partial Dependency



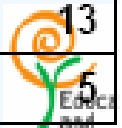
Second normal form: Table in 2NF

STUDENT TABLE

Student#	StudentName	DateofBirth
101	Davis	04-Nov-1986
102	Daniel	06-Nov-1987
103	Sandra	02-Oct-1988
104	Evelyn	22-Feb-1986
105	Susan	31-Aug-1985
106	Mike	04-Feb-1987
107	Juliet	09-Nov-1986
108	Tom	07-Oct-1986
109	Catherine	06-Jun-1984

COURSE TABLE

Course#	Course Name	Pre Requisite	Duration InDays
M1	Basic Mathematics		11
M4	Applied Mathematics	M1	7
H6	American History		4
C1	Basic Chemistry		5
C3	Bio Chemistry	C1	11
B3	Botany		8
P1	Basic Physics		8
P3	Nuclear Physics	P1	13
B4	Zoology		5



Source: Infosys Campus Connect Study Material

Second normal form: Table in 2NF

Student#	Course#	Marks	Grade
101	M4	82	A
102	M4	62	C
101	H6	79	B
103	C3	65	B
104	B3	77	B
102	P3	68	B
105	P3	89	A
103	B4	54	D
105	H6	87	A
104	M4	65	B

Exam_Date Table

Course#	DateOfExam
M4	11-Nov-04
H6	22-Nov-04
C3	16-Nov-04
B3	26-Nov-04
P3	12-Nov-04
B4	27-Nov-04

Source: Infosys Campus Connect Study Material

Second normal form ... Example

Example: The following relation is in First Normal Form, but not Second Normal Form:

Cust_Order_table

OrderNo	Customer	ContactPerson	Total
1	Acme Widgets	John Doe	\$134.23
2	ABC Corporation	Fred Flintstone	\$521.24
3	Acme Widgets	John Doe	\$1042.42
4	Acme Widgets	John Doe	\$928.53

OrderNo Customer → Total

Customer → ContactPerson

Second normal form ... Example

Customer table

Customer	ContactPerson
Acme Widgets	John Doe
ABC Corporation	Fred Flintstone

Customer → ContactPerson

Order_table

OrderNo	Customer	Total
1	Acme Widgets	\$134.23
2	ABC Corporation	\$521.24
3	Acme Widgets	\$1042.42
4	Acme Widgets	\$928.53

OrderNo Customer → Total

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 (i.e., $\beta \subseteq \alpha$)
- α is a superkey for R

Example schema *not* in BCNF:

bor_loan = (customer_id, loan_number, amount)

because ***loan_number \rightarrow amount*** holds on ***bor_loan*** but ***loan_number*** is not a superkey

Decomposing a Schema into BCNF

- Suppose we have a schema R and a **non-trivial dependency** $\alpha \rightarrow \beta$ causes a violation of **BCNF**.

We decompose R into:

- $(\alpha \cup \beta)$
- $(R - (\beta - \alpha))$

- In our example,

- $\alpha = \text{loan_number}$
- $\beta = \text{amount}$

and **bor_loan** is replaced by

- $(\alpha \cup \beta) = (\text{loan_number}, \text{amount})$
- $(R - (\beta - \alpha)) = (\text{customer_id}, \text{loan_number})$

Decomposing a Schema into BCNF

- *Lending-schema = (B_name, assets, B_city, L_no, cust_name, amount)*

B_name → assets B_city (not trivial and B_name is not a super key)

L_no → amount B_name (not trivial and L_no is not a super key)

Candidate key for this Schema is { L_no, cust_name}. This Schema is not in BCNF form. So decompose this schema into below given two schemas

Branch-schema = (B_name, B_city, assets)

Loan-info-schema = (B_name, cust_name, L_no, amount)

- *B_name → assets B_city*, the augmentation rule for functional dependencies implies that *B_name → B_name assets B_city*
- *B_name is super key in Branch_schema.*

Decomposing a Schema into BCNF

Loan-info-schema = (B_name, cust_name, L_no, amount)

L_no → amount B_name (not trivial and L_no is not a super key)

- *This Schema is not in BCNF form. So decompose this schema into below given two schemas*

Loan-schema = (B_name, L_no, amount)

Borrow-schema = (cust_name, L_no)

- *Both of these two schemas are in BCNF.*
- *Decomposition of Lending-schema to all these three schema Branch-schema, Loan-schema and Borrow-schema having dependency preservation and lossless decomposition.*

Third Normal Form

- A relation schema R is in third normal form (3NF) if for all:

$$\alpha \rightarrow \beta \text{ in } F^+$$

at least one of the following holds:

- $\alpha \rightarrow \beta$ is trivial (i.e., $\beta \in \alpha$)
- α is a superkey for R
- Each attribute A in $(\beta - \alpha)$ is contained in a candidate key for R .

(NOTE: each attribute may be in a different candidate key)

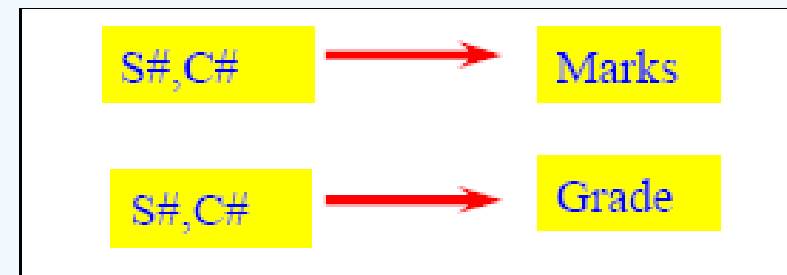
- If a relation is in BCNF it is in 3NF (since in BCNF one of the first two conditions above must hold).
- Third condition is a minimal relaxation of BCNF to ensure dependency preservation.

Third Normal Form

A relation R is said to be in the Third Normal Form (3NF) if and only if

- It is in 2NF and*
- No transitive dependency exists between non-key attributes and key attributes.*

- STUDENT# and COURSE# are the key attributes.
- All other attributes, except grade are non-partially, non-transitively dependent on key attributes.
- **Student#, Course# - > Marks**
- **Marks -> Grade**



Note : - All transitive dependencies are eliminated



Third Normal Form

Note that 3NF is concerned with transitive dependencies which do not involve candidate keys. A 3NF relation with more than one candidate key will clearly have transitive dependencies of the form:

primary_key \rightarrow other_candidate_key \rightarrow any_non-key_column

Third Normal Form

Student#	Course#	Marks	Grade
101	M4	82	A
102	M4	62	C
101	H6	79	B
103	C3	65	B
104	B3	77	B
102	P3	68	B
105	P3	89	A
103	B4	54	D
105	H6	87	A
104	M4	65	B

Source: Infosys Campus Connect Study Material

Third Normal Form

Student#	Course#	Marks
101	M4	82
102	M4	62
101	H6	79
103	C3	65
104	B3	77
102	P3	68
105	P3	89
103	B4	54
105	H6	87
104	M4	65

MARKSGRADE TABLE		
UpperBound	LowerBound	Grade
100	95	A+
94	85	A
84	70	B
69	65	B-
64	55	C
54	45	D
44	0	E

Source: Infosys Campus Connect Study Material

Third Normal Form: Motivation

- There are some situations where
 - BCNF is not dependency preserving, and
 - efficient checking for FD violation on updates is important
- Solution: define a weaker normal form, called Third Normal Form (3NF)
 - Allows some redundancy (with resultant problems; we will see examples later)
 - But functional dependencies can be checked on individual relations without computing a join.
 - There is always a lossless-join, dependency-preserving decomposition into 3NF.

Comparison of BCNF and 3NF

■ Relations in BCNF and 3NF

- Relations in BCNF: no repetition of information
- Relations in 3NF: problem of repetition of information

■ Decomposition in BCNF and in 3NF

- It is always possible to decompose a relation into relations in 3NF and
 - ▶ the decomposition is lossless
 - ▶ **dependencies are preserved**
- It is always possible to decompose a relation into relations in BCNF and
 - ▶ the decomposition is lossless
 - ▶ **May some of the dependencies are not preserved.**

Merits of Normalization

- Normalization is based on a mathematical foundation.
- Removes the redundancy to a greater extent. After 3NF, data redundancy is minimized to the extent of foreign keys.
- Removes the anomalies present in INSERTs, UPDATEs and DELETEs.

Demerits of Normalization

- Data retrieval or SELECT operation performance will be severely affected.
- Normalization might not always represent real world scenarios.

Summary of Normal Forms

Input	Operation	Output
Un-normalized Table	Create separate rows or columns for every combination of multivalued columns	Table in 1 NF
Table in 1 NF	Eliminate Partial dependencies	Tables in 2NF
Tables in 2 NF	Eliminate Transitive dependencies	Tables in 3 NF
Tables in 3 NF	Eliminate Overlapping candidate key columns	Tables in BCNF



Source: Infosys Campus Connect Study Material

Points to Remember:

Normal Form	Test	Remedy (Normalization)
1NF	Relation should have atomic attributes. The domain of an attribute must include only atomic (simple, indivisible) values.	Form new relations for each non-atomic attribute
2NF	For relations where primary key contains multiple attributes (composite primary key), non-key attribute should not be functionally dependent on a part of the primary key.	Decompose and form a new relation for each partial key with its dependent attribute(s). Retain the relation with the original primary key and any attributes that are fully functionally dependent on it.
3NF	Relation should not have a non-key attribute functionally determined by another non-key attribute (or by a set of non-key attributes). In other words there should be no transitive dependency of a non-key attribute on the primary key.	Decompose and form a relation that includes the non-key attribute(s) that functionally determine(s) other non-key attribute(s).



Source: Infosys Campus Connect Study Material

Database System Concepts, 5th Ed.

©Silberschatz, Korth and Sudarshan