

Data Normalization & Functional Dependency

A relation in a relational database is based on a relation schema, which consists of a no. of attributes. A relational database is made up of a no. of relational schemas. In this unit, the major focus is on the issues involved in the design of a database schema using the relational model.

Main objective of this chapter is to design a good relational schema. Normally, the relational schemas face some undesirable problems such as redundancy, redundancy & various types of anomalies (Insertion, Update & Deletion) which <sup>inconsistencies</sup> constitute a bad database design.

Consider the relational schema Supplier (Name, address, Item, Price)

Supplier		Item	Price
Name	address		
Harell's India Ltd	Civil Lines	Cable	400
Harell's India Ltd	Civil Lines	Switch	300
Bajaj Hd.	South Ext, Part I	Fan	450
Bajaj Ltd.	South Ext, Part I	Switchgear	500
:	:	:	:

Here, we can see several problems associated with this relation schema which prove inadequate:

- i) Redundancy The aim of the database system is to reduce redundancy meaning that information is to be stored only once, storing information several times leads to wastage of storage space and an increase in the total size of the data stored. In the above table

Relationships schema, the address of the supplies are repeated several times for each data item supplied.

ii) Update anomalies - Multiple copies of the same fact may lead to update anomalies or inconsistencies, when an update is made and only some of the multiple copies are updated. Thus, a change in the address of supplier must be made for consistency in all the tuples pertaining to the Shame. If one of them is not changed, there will be an inconsistency in the data.

iii) Insertion anomalies -

If this is the only relation in the database showing the association between a Sname and items supplied the fact that we cannot insert an item in the database unless a supplier name is entered in the database. relation.

iv) Deletion anomalies

If this is the only relation showing the association between Sname & Item supplied, the fact that which item is supplied by which supplier will be lost.

In the relational model, the above problem can be overcome by the process called Decomposition.

Decomposition (Definition) - The decomposition of a relational schema  $R = (A_1, A_2, \dots, A_n)$  is its replacement by a set of relational schemas

$\{R_1, R_2, \dots, R_m\}$  such that  $R_i \subseteq R$  for  $i = 1 \text{ to } m$  &  $R_1 \cup R_2 \cup \dots \cup R_m = R$

Hence, the problem in the Supplier schema can be resolved only if we can replace it with the following relational schemas:-

SA (Sname, address) -  $R_1$   
 SIP (Sname, item, price) -  $R_2$

The first schema SA gives the address of each supplier exactly once, hence there is no redundancy.

We can insert an address for a supplier even if it currently supplies no items.

Similar concepts will be applied to SIP.

where, Supplier = SA  $\cup$  SIP &  
 $SA \subseteq \text{Supplier}$ ,  $SIP \subseteq \text{Supplier}$ .

### Functional Dependency

- It is a constraint between two sets of attributes from the database.
- A functional dependency (FD), denoted by  $X \rightarrow Y$ , between two sets of attributes X & Y that are subsets of R specifies a constraint on the possible tuples that can form a relation State  $\tau$  of R. The constraint is that, for any two tuples  $t_1$  &  $t_2$  in  $\tau$  that have

$$t_1[X] = t_2[X], \text{ then we must have } t_1[Y] = t_2[Y]$$

This means that the values of the Y component of a tuple depend on the values of the X component.

or alternatively, the values of X-Component of a tuple uniquely or functionally determines the value of its other Y-Component.

The FD  $X \rightarrow Y$

$Y$  is said to be functionally dependent on  $X$  or  $X$  functionally determines  $Y$ .

$X$  is called left hand side attribute also called determinant &  $Y$  is called right hand side attribute.

From the above relation Supplier, the following FDs occurs -

Sname  $\rightarrow$  address

Sname, Item  $\rightarrow$  Price

In other word, whenever 2 tuples of R agree in their X-values, they also agree on their Y-values.

Consider another example ..

The following FDs holds:

S	S#	Sname	Status	City
	S1	Smith	20	London
	S2	Jones	10	Paris
	S3	Joe	30	Paris

$S\# \rightarrow Sname$   
 $S\# \rightarrow Status$   
 $S\# \rightarrow City$

but  
 $City \nrightarrow S\#$ .

e.g 2

(a)  $SSN \rightarrow Ename$ , specify that the value of an employee's social security no uniquely determines e name.

(b) Pnumber  $\rightarrow \{Pname, Plocation\}$

The value of Pnumber uniquely determines the pname & location

(c)  $SSN, Pnumber \rightarrow hours$

a combination of SSN & Pnumber values uniquely determines the no. of hrs.

Prime Attribute - An attribute is said to be a prime attribute if it is the part of the candidate key. Otherwise, it is called non-prime attribute.

8. Student (Rollno, name, address, subject).
- here, Rollno. is the prime attribute & name, address & subjects are called non-prime attribute; because Rollno is the P. key. & FDI's are  $\text{Rollno} \rightarrow \text{name}$   
 $\text{Rollno} \rightarrow \text{address}$  &  
 $\text{Rollno} \rightarrow \text{subject}$ .

SPJ (Sno, Pro, Jno, Qty)

here Sno, Pro, Jno are the prime attributes.  
Qty is non-prime attribute as the FD for this is  
 $\text{Sno}, \text{Pro}, \text{Jno} \rightarrow \text{Qty}$

Different types of functional dependency:

i) Trivial FD: A FD in the form  $X \rightarrow Y$  is trivial if  $Y \subseteq X$

ii) fully functional dependency

A functional dependency  $X \rightarrow Y$  is a full functional dependency if removal of any attribute A from X means that the dependency doesn't hold any more.

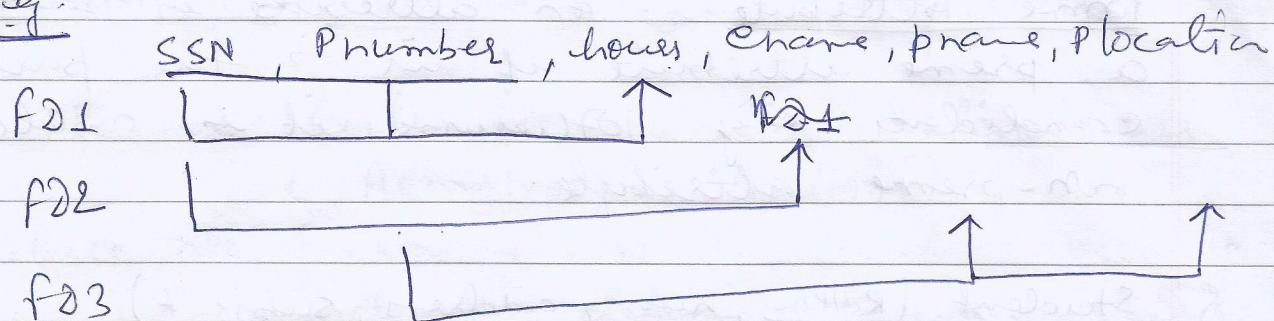
i.e. for any  $A \in X$   
 $(X - \{A\})$  doesn't functionally determine Y.

e.g. Student (Rollno, name, address, subject) & an FD  $\text{Rollno} \rightarrow \text{name}$  is fully FD. But in SPJ (Sno, Pro, Jno, Qty) & an FD

$\text{Sno}, \text{Pro} \rightarrow \text{Qty}$  doesn't hold as Qty is not fully dependent on the primary key which is  $\text{Sno}, \text{Pro} \text{ & } \text{Jno}$ .

Note

e.g. EMP-PROJ



F1: SSN, Pnumber  $\rightarrow$  hours  
is fully FD.

F2: SSN  $\rightarrow$  ename

not fully FD because the P key is the combination of SSN & Pnumber.

F3: Pnumber  $\rightarrow$  Pname, plocation

not fully FD. - same reason.

### iii) Partial Functional Dependency

A FD  $X \rightarrow Y$  is a partial dependency if some attribute  $A \in X$ , can be removed from  $X$  and the dependency still holds.

e.g. In the above example.

F2: SSN  $\rightarrow$  ename is partial dependency as SSN is part of the P-key.

F4: SSN, Pnumber  $\rightarrow$  ename is also partial dependency as SSN  $\rightarrow$  ename holds.

### iv) Transitive FD -

A FD  $X \rightarrow Y$  in a relational schema R is transitive dependency if there is a set of attributes Z that is neither a candidate

Key not a subset of any key of R & both  
 $X \rightarrow Z$  &  $Z \rightarrow Y$  hold.

g.

Given R(A, B, C, D, E) & the FDs

$$F = \{ A \rightarrow C, B \rightarrow D, C \rightarrow E, B \rightarrow C \} \text{ then}$$

from the above FDs, of

$$A \rightarrow C \text{ & } C \rightarrow E \text{ then}$$

$A \rightarrow E$  i.e. E is transitively dependent on A.

2. Prof(

(Prof-name, Deptname, chairperson) & the FDs are

$$\text{Prof-name} \rightarrow \text{Deptname} \text{ &}$$

$$\text{Deptname} \rightarrow \text{chairperson} \text{ then}$$

$$\text{Prof-name} \rightarrow \text{chairperson} \text{ holds.}$$

## Normalization

### What is normalization?

Normalization is a process of decomposing the given relation schema into a set of relational schemas based on their FDs and primarily keeps to achieve the desirable properties of

- minimizing data redundancy
- minimizing the insertion, deletion & update anomalies
- Reduces the memory usage

It also provides a set of normal form tests that can be carried out on individual relation schemas so that the relational

database can be normalized to any desired degree.

### Inference Rules for Functional Dependencies ( Armstrong's Axioms)

- There are following rules that logically implied functionally dependency with a set of FD's  $F$  & suppose  $X, Y \& Z$  denotes sets of attributes over a relational schema  $R$ .
  - also called Armstrong's axioms
- We use the notation  $F \models X \rightarrow Y$  i.e.  $X \rightarrow Y$  is inferred from  $F$ .

#### R1. Reflexivity rule

If  $X$  is a set of attributes &  $Y \subseteq X$ , then  
 $X \rightarrow Y$  holds if  $F \models X \rightarrow Y$

#### R2. Augmentation rule

If  $X \rightarrow Y$  holds &  $Z$  is a set of attributes,  
then  $XZ \rightarrow YZ$  or  $F \models XZ \rightarrow YZ$

#### R3. Transitivity Rule

If  $X \rightarrow Y$  &  $Y \rightarrow Z$  then  
 $X \rightarrow Z$  holds. i.e.  $F \models X \rightarrow Y$

#### R4. Union or Additive rule

If  $X \rightarrow Z$  &  $X \rightarrow Y$ , then  
 $X \rightarrow YZ$  holds. or  $F \models X \rightarrow YZ$

#### R5. Decomposition or Projective Rule

If  $X \rightarrow YZ$  holds then  
 $X \rightarrow Y$  holds &  
 $X \rightarrow Z$  holds.  
or  $F \models (X \rightarrow Y, X \rightarrow Z)$

~~Contd.~~

## Inference rules for FDs (Functional Dependencies)

R6 Pseudo-transitivity rule

If  $X \rightarrow Y$  holds &  $WY \rightarrow Z$  holds  
then  $WX \rightarrow Z$  holds.

Example: Prove the following:

$$1) X \rightarrow Y \text{ & } Z \subseteq Y \models \{X \rightarrow Z\}$$

Proof: i)  $X \rightarrow Y$  given

ii) As  $Z \subseteq Y \Rightarrow Y \rightarrow Z$  (R1, Reflexivity)

iii)  $X \rightarrow Y, Y \rightarrow Z \models X \rightarrow Z$  (Transitivity, R3)

$$2) \{W \rightarrow Y, X \rightarrow Z\} \models \{WX \rightarrow Y\}$$

i)  $W \rightarrow Y$  given

$\Rightarrow WX \rightarrow YX$  (R2, Augmentation)

$YX \rightarrow Y$  as  $Y \subseteq YX$  (R1, Reflexivity)

$\therefore WX \rightarrow Y$  (Transitivity R3).

## Normal Forms

The normal forms are used to ensure that the various types of anomalies & inconsistencies are not introduced into the database. For determining whether a particular relation is in normal forms or not, the FDs between the attributes in the relation are examined & not the current contents of the relation. Various types of normal forms are used in relational database are as follows:

1. First normal form (1NF)
2. Second " " (2NF)
3. Third " " (3NF)
4. Boyce-Codd " " (BCNF)
5. Fourth " " (4NF)
6. Fifth " " (5NF).

A relational schema is said to be in a particular ~~relax~~ normal forms if it satisfies a certain prescribed set of conditions as in discuss different normal form. To explain the various normal forms, first consider an example of Innormalized relation.

P

Pno	Pname	Eid	Ename	Category	Rate
P001	FIS	E001	Smith	A	700
○	○	E002	Robin	B	500
P002	EIS	E006	John	A	700
○	○	E007	Joe	B	300

fig(a)

Here, this unnormalized rel.<sup>n</sup>P contains non-atomic values.

i) first normal form (1NF) - A relation is said to be in 1NF, if the values in the domain of each attribute of the relation are atomic i.e. only one value is associated with each attribute and the value is not a set of values or list of values.

Hence, the above relation P can be normalized into a relation (1NF)

ex-1 P

Projno	Projname	Eno	Ename	Category	Rate
P001	FIS	E001	Smith	A	700
P001	FIS	E002	Robin	B	500
P002	EIS	E006	John	A	700
P002	EIS	E007	Joe	B	300

info Emp-proj(Projno, Projname)  $\Sigma$  Emp(eno, Projno, ename, Category, Rate)  
PK PK

A table is in 1NF if

- there are no repeating groups.
- all the key attributes are defined
- all attributes are dependent on a primary key.

example 2.

(a) Department

	Dname	eno	Dmgr-ssn	Dloc	
	Research	5	11325	(Delhi, Bangalore, Noida)	→ Multivalued

(b)

	Dname	eno	Dmgr-ssn	Dloc	
	Research	5	11325	(Delhi, Bangalore, Noida)	→ Multivalued
	Admin	4	52325	(Delhi, Chennai)	→ Multivalued.
	IT	1	23432	Noida	→

State of Relation Department

Data 1NF

(c) Department (with P. Key Dno + Dloc)

Dname	<u>Dno</u>	Emp-ssn	<u>Dloc</u>
Research	5	11325	delhi
Research	5	11325	Banglore
Research	5	11325	Noida
Admin	4	52325	delhi
Admin	4	52325	chennai
IT	1	23432	Noida

atomic value

atomic value

Q: Normalized the relation (EMP-Proj) into 1NF.

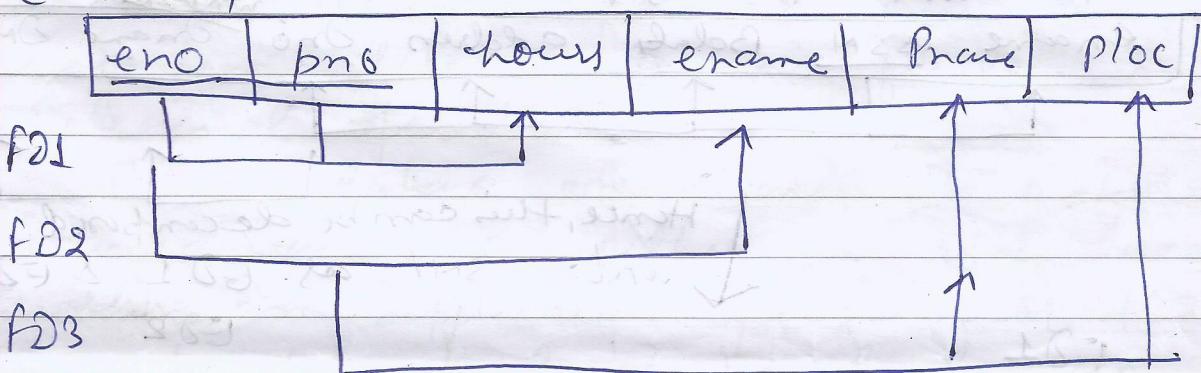
EMP-Proj (ssn, ename, pno, hours)

EMP-Proj 1 (ssn, ename) ↓ PKEMP-Proj 2 (ssn, pno, hours)

## (2) Second Normal Form (2NF)

- The 2NF is based on the concept of full functional dependency.
  - A Relation is in 2NF if it is in 1NF and every non-prime attributes are fully functionally dependent on the primary key of R.
- eg 1

EMP-Proj



Can be decomposed into 3 1st normal relations as:

EP1

eno	pno	hours

FD1: eno → pno

EP2

eno	ename

FD2: eno → ename

EP3

pno	Phone	Ploc

FD3: pno → Phone

FD4: Phone → Ploc

### ③ Third Normal Form (3NF)

- 3NF is based on the concept of transitive dependency

- A relation schema R is in 3NF if it is in 2NF and no non-prime attribute of R is transitively dependent on the primary key.

Ex ① EMP-Dept - is not in 3NF because of the transitive dependency of Dname on SSN via Dno.

ename	<u>ssn</u>	Bdate	address	Dno	<sup>non-prime</sup> Dname	DEPSSN

Hence, this can be decomposed into 3NF as E1 & E2.

E1

E2

ename	<u>ssn</u>	Bdate	address	Dno	Dno	Dname	DEPSSN

Ex 2

Major

non-prime

Stud-name(PK)	Major	Department

as the non-prime attribute Major is transitively dependent on R.

M1

M2

Stud-name	Major	Major	Department

i.e.  $\text{Stud-name} \rightarrow \text{Major}$        $\Rightarrow \text{Stud-name} \rightarrow \text{Department}$   
 $\text{Major} \rightarrow \text{Department}$       shouldn't occurs.