

Unit-3
Database Design Theory: Fundamental Dependencies, Normal forms
Relational DB Design Algorithm.

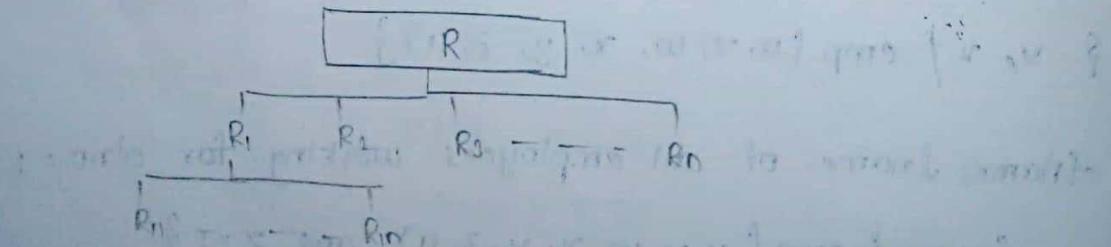
RDB Design approach:

1. ER - to - Relational mapping algorithm
2. Bottom - up Design (RDB design by synthesis)

Based on similarity attributes are grouped into relations if there exists a binary relations between 2 relations they are identified if not separation of attributes takes place

3. top - down Design (RDB design by analysis)

Initially all attributes encountered from requirement are made in single relation called Universal Relation Schema. Again it is decomposed R_1, R_2, \dots, R_n and each relation is decomposed into sub relation



This process continues until leaf node or bottle relations are 3NF or BCNF forms

Design guidelines for RDB Design:

1. Imparting clear semantics to attributes of relation
2. Reducing the redundancy with relations
3. Reducing null valued attributes within a relation

4. This allowing generation of spurious tuples in relations.

① → semantic refers meaningful

meaningful attributes are to be provided in relation

company RDB = { employee, department, project, dependent,
works-on, dept-loc }

Acc to this, in a RDB design construct a relation
with attributes corresponding to ~~not more than~~ single entity type. Or relationship do not combine
attributes belonging to more than one entity type.

② → Acc to this, do not create

Update Anomalies is a problem that is encountered
if redundancy problem is not handled properly with ⁱⁿ relations

1. Insert Anomaly - User faces problems during
insertion of data & updating

2. Deletion Anomaly - problem during deletion

3. Modification Anomaly

③ → Acc to this, form a relation in such a way that
no null value is allowed

④ → Spurious : meaning less

Meaning less tuples even after join.

Allow the join only on relevant attributes

e.g. primary & foreign keys

→ Functional Dependency (FD):

FD is a constraint between two subsets of attributes X and Y in a relation R represented by $X \rightarrow Y$

This constraint holds that all tuples forming a relation state r in a relation R will specify for any 2 tuples t_1 and t_2 if $t_1[x] = t_2[x]$ then $t_1[y] = t_2[y]$ holds means the value for y in t_1 is uniquely determined based on value of x in t_1 then the same must hold in t_2 also. Here x determines a unique value for y .

Relation

emp-proj

	ssn	pno	pname	place	ename	hours
$t_1 \rightarrow$	12345	1	pro-x	newyork	smith	10
$r(R) t_2 \rightarrow$	12346	1	pro-x	newyork	franklin	10.5
	12345	2	pro-y	chicago	smith	12

$$F = \{ ssn \rightarrow ename \}$$

select ename from employee where ssn = 12345 \hookrightarrow smith

select ename from employee where ssn = 12346 \hookrightarrow franklin

select pno, ename from employee where ssn = 12345 \hookrightarrow

so dependency b/w ssn & pno, dependency b/w ssn & ename

$f = \{ \text{ssn} \rightarrow \text{ename},$
 FD's set
 $\text{pno} \rightarrow \{\text{pname}, \text{ploc}\},$
 $\{\text{ssn}, \text{pno}\} \rightarrow \text{hours} \}$

teach	teacher	course	textbook
mpp	dbmu	navatha	
vss	dr	allenweis	
nn	dbmu	navatha	
mpp	broad	benett	

$f = \{ \text{course} \rightarrow \text{textbook} \}$

$\text{teacher} \rightarrow \text{course} X$

to

→ Inference Rules:

IR1: Reflexive Rule: If X is a proper superset $\supseteq Y$ then
always $X \rightarrow Y$ or $X \rightarrow X$ is a valid FD

IR2: Augmentation Rule: $X \rightarrow Y$ is valid FD then \vdash
(intuitively deduce or derive) $XZ \rightarrowYZ$ is also a valid FD.

Ex: $\text{ssn} \rightarrow \text{ename}$

$\{\text{ssn}, \text{ename}\} \rightarrow \{\text{ename}, \text{pno}\}$

Transition Rule: IR3 : If $X \rightarrow Y$ is valid FD & $Y \rightarrow Z$ is
valid FD $\vdash X \rightarrow Z$ is valid FD

These three are called Armstrong

IR4: Decomposition Rule: $x \rightarrow \{A_1, A_2, A_3\} \models x \rightarrow A_1, x \rightarrow A_2, x \rightarrow A_3$

IR5: Union Rule (Reverse of decomposition)

If $x \rightarrow A_1, x \rightarrow A_2, \dots, x \rightarrow A_3 \models x \rightarrow \{A_1, A_2, \dots, A_3\}$

IR6: pseudo transitive rule: If $x \rightarrow y$ is available then $wy \rightarrow z$
 \models (then we can infer) $wx \rightarrow z$

→ Closure Set of FD: If F is a FD set given on R
then closure set of F^+ for F includes all FD's in F
along with additional FD's inferred through IR1 to IR3

→ Algorithm for Obtaining Closure Set: This method
obtains closure set for the given FD set in F
based on each LHS attribute x

$$x^+ = x \quad (x^+ \leftarrow x)$$

Repeat

$$\text{old } x^+ = x^+$$

for each FD $y \rightarrow z$ in F

if $x^+ \supseteq y$ then $x^+ \cup z$

until $(x^+ = \text{old } x^+)$

R \Rightarrow ssn pno pname place ename hours

F = {ssn \rightarrow ename,

pno \rightarrow {pname, place}

{ssn, pno} \rightarrow hours}

for each LHS attribute x of FD in F ,

$$x = \text{ssn}$$

$$\{\text{ssn}\}^+ = \{\text{ssn}, \text{ename}\}$$

$$x = \text{pno}$$

$$\{\text{pno}\}^+ = \{\text{pno}, \text{pname}, \text{ploc}\}$$

$$x = \{\text{ssn}, \text{pno}\}$$

$$\{\text{ssn}, \text{pno}\}^+ = \{\text{ssn}, \text{pno}, \text{ename}, \text{pname}, \text{ploc}, \text{hours}\}$$

then

$$F^+ = \{ \text{ssn} \rightarrow \{\text{ssn}, \text{ename}\},$$

$$\text{pno} \rightarrow \{\text{pno}, \text{pname}, \text{ploc}\}$$

$$\{\text{ssn}, \text{pno}\} \rightarrow \{\text{ssn}, \text{pno}, \text{ename}, \text{pname}, \text{ploc}, \text{hours}\}$$

→ Trivial FD : A FD $x \rightarrow y$ is trivial FD if $x \subseteq y$

Otherwise FD is non trivial FD

→ Equivalence functional Depending sets:

Consider two FD set E & F

i) cover of FD: FD set E is said to be covered if it contains all FDs available in F , in its E^+ . Then we say E covers F

ii) equivalent FD sets - E and F are equivalent FD sets if $E^+ = F^+$ (or)

If E covers F and F covers E Then E and F are two equivalent FD

→ Check $F = \{ A \rightarrow C, AC \rightarrow D, E \rightarrow AD, F \rightarrow H \}$ and
 $G_1 = \{ A \rightarrow CD, E \rightarrow AH \}$ equivalent are not.

$$R = \{ A | B | C | D | E | H \}$$

$$X = A$$

$$\{ A^+ \} = \{ A, C \}$$

$$X = AC$$

$$\{ ACy^+ \} = \{ A, C, D \}$$

$$X = E$$

$$\{ E^+ \} = \{ E, A, D, H \}$$

$$X \neq 1$$

→ Minimal Cover set: Minimal coverset F for the given FD set E :-

algorithm to obtain a minimal set for the given FD set E

1. Initially assign E to F ; $F = E$
2. If F contains FD of the form $X \rightarrow A_1, A_2, \dots, A_n$ then replace it in F by $X \rightarrow A_1, X \rightarrow A_2, \dots, X \rightarrow A_n$
3. for each FD $X \rightarrow A$ do
 - if $B \in X$ and $\{F - \{X \rightarrow A\} \cup \{(X - \{B\}) \rightarrow A\}\}$ is same as F then replace $X \rightarrow A$ in F by $\{(X - \{B\}) \rightarrow A\}$
4. for each FD $X \rightarrow A$ do
 - if $F - \{X \rightarrow A\}$ is same as F means $X \rightarrow A$ is redundant functional dependency may be appeared in F but can be obtained by applying inference rules then remove $X \rightarrow A$ from F

→ Compute minimal set F for $E = \{B \rightarrow A, D \rightarrow A, AB_0 \rightarrow D\}$

Method minimal set F = $\{E\} = \{B \rightarrow A, D \rightarrow A, AB \rightarrow D\}$

Step-2 Not applied

Step-3 FD $X \rightarrow A$ of F is with more than one attribute in X

$AB \rightarrow D$

$\{F - \{AB \rightarrow D\}\} \cup \{A \rightarrow D\}$ is not as same as F

$\{F - \{AB \rightarrow D\}\} \cup \{B \rightarrow D\}$ is also not same as F

Step 3 cannot be applied.

4. Checking for Redundant FD $X \rightarrow A$ in F

① Let $B \rightarrow A$ from F

Then add B to LHS & RHS, according to

Augmentation rule

$$BB \rightarrow AB \Rightarrow B \rightarrow AB \Rightarrow B \rightarrow D$$

so $AB \rightarrow D$ is replaced with $B \rightarrow D$

$$F = \{B \rightarrow A, D \rightarrow A, B \rightarrow D\},$$

($X \rightarrow Y$, $Y \rightarrow Z$ then replace with $X \rightarrow Z$ according to)

By applying $X \rightarrow Y$, $Y \rightarrow Z$ and $X \rightarrow Z$ is available in F then that is redundant so remove that from F

$$B \rightarrow D \quad D \rightarrow A \models B \rightarrow A$$

$$F = \{B \rightarrow D, D \rightarrow A\}$$

algorithm to compute key attribute K from Relation R:
with set of attributes

for each attribute A in R

compute Σ_{K-A}^+ w.r.t. FD set F on. R

if $\{\Sigma_{K-A}\}^+$ contains all the attributes of R

then $K = K - \{A\}$

$$\rightarrow R = \boxed{A | B | D}$$

$$F = \{B \rightarrow A, D \rightarrow A, AB \rightarrow D\}$$

Normal forms based on keys & FD's:

1. First Normal Form (1NF)
2. Second Normal Form (2NF)
3. Third Normal Form (3NF)
4. Boyce - Codd Normal Form (BCNF) } minimal normal form
5. Fourth Normal Form (4NF)
6. Fifth Normal Form (5NF) } 4N Normal Form

- ① → concept based on atomic values (simple, composite, multivalued attributes) and Nested relations
- ② → concept based on fully functional dependency (FFD)
- ③ → concept of transitive dependency
- ④ → concept based on superkey
- ⑤ → multi valued dependency (MVD)
- ⑥ → Join dependency

4. prime attribute : If attribute is member of candidate key, then it is prime attribute
5. Non prime attributes If attribute in a relation is not member of candidate
6. Foreign key : If attribute in a relation is reference an attribute in other relation and should have same domain.

Denormalization: It is a process of storing a join of lower normal form relations in order to satisfy higher normal form condition.

→ 1NF: If a relation R is said to be 1NF if doesn't contain multivalued attributes, nested relations (and sometimes composite attributes).

Eg: dept

dno	dname	dmgrssn	dloc
1	research	12345	{Blore, chennai, hyd}
4	sales	12346	hyd

dept contains multivalued attribute dloc, so it is not in 1NF. To make 1NF for each value store as separate tuple.

then

dno	dname	dmgrssn	dloc
1	research	12345	blore
1	research	12345	chennai

1	research	12345	hyd
4	sales	12346	hyd

But due to this redundancy occurs. To avoid that

- (2) for multi-valued attribute if all distinct values are known create distinct columns for those values

dno	dname	dmgrssn	dloc1	dloc2	dloc3
1	research	12345	blore	chennai	hyd
4	sales	12346	hyd		

But to this null value problem arises. To avoid this

- (3) Remove the multivalued column from that relation and create a new relation along with primary key of old relation. and make that attribute (pk) as foreign key. Combination of those attributes make as key attribute

R ₁	<u>dno</u>	<u>dname</u>	<u>dmgrssn</u>	key attribute
1		research	12345	
4		sales	12346	

R ₂	<u>fk</u> <u>dno</u>	<u>dloc</u>	<u>dno</u> <u>dloc</u>
	1	blore	
	1	chennai	
	1	hyd	
	4	hyd	

Here also redundancy occurs but minimal redundancy is obtained.

This procedure is the preferable procedure.

$R :=$

	emp		project		← Nested Relation
	ssn	ename	pno	hours	
R_1	12345	smith	1	10	
R_2	12345	smith	4	15	
	12345	smith	5	16	
	12346	smith	1	10	
	12346	John	4	10.5	
		enreith			
		John			

(few attribute are taken from emp entity type few from project entity type and some from works-on entity type)

$R_1:$

ssn	ename
12345	smith
12346	John

pno	hours	essn
1	10	12345
4	15	12345

Add key attribute of one relation to another.

→ 2NF: A relation R is said to be in 2NF if every non-prime attribute is fully functionally dependent on the primary key of R if it is 1NF

→ FFD: If $X \rightarrow Y$ is a FD for each attribute B is member of X ($B \in X$), then $X - \{B\} \rightarrow Y$ do not holds, then FD is said to be ffd

otherwise if $X - \{B\} \rightarrow Y$ holds a valid FD then $X \rightarrow Y$ is partial FD

ssn	pno	ename	hours	pname	ploc

Given FD set $F = \{ \text{ssn} \rightarrow \text{ename}, \text{pno} \rightarrow \{\text{pname}, \text{ploc}\}, \{\text{ssn}, \text{pno}\} \rightarrow \text{hours} \}$

key attribute = $\{\text{ssn}, \text{pno}\}$

remove one attribute from key attribute

$(\text{ssn}, \text{pno}) \rightarrow \text{ename}$

$x \rightarrow y$

det $B = \text{ssn}$ (remove B)

$\text{pno} \rightarrow \text{ename}$

We have to check whether it is FD or not

so $\text{pno} \rightarrow \text{ename}$ is invalid

det $B \rightarrow \text{pno}$ (remove B)

$\text{ssn} \rightarrow \text{ename} \rightarrow \text{valid}$

Here both should be invalid only then it is in 2NF i.e. not 1FD (1PFN)

so $(\text{ssn}, \text{pno}) \rightarrow \text{ename}$ is 1FD

So this relation is not in 2NF

NPA $y = \text{hours}$

$X = \{\text{ssn}, \text{pno}\}$

$B = \text{pno}$

$\text{ssn} \rightarrow \text{hours}$ is invalid

$B = \{ssn\}$

$pno \rightarrow \text{hours}$ (invalid or do not holds)

so $\{ssn, pno\} \rightarrow \text{hours}$ is FFD

$Y = \{\text{pname}\}$

$X = \{ssn, pno\}$

$\{ssn, pno\} \rightarrow \{\text{pname}\}$

$B = pno$

$ssn \rightarrow \{\text{pname}\}$ (invalid)

$B = ssn$

$pno \rightarrow \{\text{pname}\}$ (valid)

so $\{ssn, pno\} \rightarrow \{\text{pname}\}$ is pfd.

$Y = \{\text{ploc}\}$

$X = \{ssn, pno\}$

$\{ssn, pno\} \rightarrow \{\text{ploc}\}$

$B = pno$

$ssn \rightarrow \{\text{ploc}\}$ (invalid)

$B = ssn$

$pno \rightarrow \{\text{ploc}\}$ (valid)

$\{ssn, pno\} \rightarrow \{\text{ploc}\}$ is in pfd

R is not in 2NF bcoz ename, pname, ploc are not fully fd.

Take the non-prime attribute which is pfd make a new relation along with key attribute

$$R_1 = \boxed{\text{ssn} | \text{pnol hours}}$$

And non p a which is pfd make new relation with valid ffd.

$$R_2 = \boxed{\text{ssn} | \text{ename}}$$

$$R_3 = \boxed{\text{pnol} | \text{pname} | \text{ploc}}$$

$$R_4 = \boxed{\text{pnol} | \text{ploc}}$$

\rightarrow 3NF: It is based on concept of transitive dependency

FD $X \rightarrow Y$ is transitive dependent if there $X \rightarrow Z, Z \rightarrow Y$ in FD set

Given relation R is in 3NF, if it satisfies 2NF condition, and no non-prime attribute is transitively dependent on key of R, then we say R is in 3NF.

A relation R is said to be in 3NF for each non trivial FD $X \rightarrow Y$ either X must be super key or Y is a prime attribute and also must be in 2NF.

Ex: emp-dept

ssn	ename	bdate	address	sal	dno	dmgrssn	dname

not

$$F = \{ \text{ssn} \rightarrow \{\text{ename}, \text{bdate}, \text{address}, \text{sal}, \text{dno}\} \\ \text{dno} \rightarrow \{\text{dname}, \text{dmgrssn}\} \}$$

Given FD is in 1NF and 2NF also but it is not in 3NF because $\text{ssn} \rightarrow \text{dname}$, $\text{ssn} \rightarrow \text{dmgrssn}$
 All valid transitive dependency ($\text{ssn} \rightarrow \text{dno}$, $\text{dno} \rightarrow \text{dname}, \text{dmgrssn}$)
 i.e. non prime attributes $\text{dname}, \text{dmgrssn}$ are transitively dependent on SSN key of R.

To satisfy 3NF, decompose R into sub relations
 identify the attribute A from the set of non-prime attributes by which valid transitive dependencies are inferred. Then all the attributes in a Relation along with A, A is a foreign key include into relation R_1 . A along with the non-prime attributes which are transitively dependent include in relation R_2 in which A is primary key of R_2 .

R_1 & R_2 are in 3NF condition.

$R_1 =$	<table border="1"> <tr> <td><u>ssn</u></td><td>ename</td><td>bdate</td><td>address</td><td>sal</td><td>FK dno</td></tr> </table>	<u>ssn</u>	ename	bdate	address	sal	FK dno
<u>ssn</u>	ename	bdate	address	sal	FK dno		
$R_2 =$	<table border="1"> <tr> <td>dno</td><td>dname</td><td>dmgrssn</td></tr> </table>	dno	dname	dmgrssn			
dno	dname	dmgrssn					

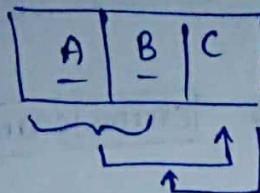
3NF

Non-prime attributes dname, dmgrssn are transitively dependent on primary key ssn.

Boyce Codd Normal Form (BCNF):

A relation in BCNF is also in 3NF but not vice-versa.

If R has FD's of form



Trivial {
x - proper subset of y
y - superset of y

$$F = \{(A, B) \rightarrow C,$$

$$C \rightarrow B\}$$

Then R is always in 3NF

A relation 'R' is said to be in BCNF, for each non-trivial FD $X \rightarrow Y$ in F^+ , X must be a super key.

→ Student

sname	cname	instructor
smith	dbms	mpp
smith	os	nn
adams	dbms	vss
scott	os	kgs
tiger	dbms	mpp

Note: 1. Two students with same name cannot register course

2. Instructor is allowed to teach only one course.

$$F = \{ (sname, cname) \rightarrow instructor, \\ instructor \rightarrow cname \}$$

Relation student is in 3NF but not in BCNF because with non-trivial FD $instructor \rightarrow cname$, instructor is not.

super key. It is in 3NF because of $AB \rightarrow C$ and $C \rightarrow B$

In order to satisfy BCNF condition the above relation can be decomposed into 3 ways.

i, R_1

sname	cname
-------	-------

R_2

sname	instructor
-------	------------

ii, R_1

cname	instructor
-------	------------

R_2

cname	sname
-------	-------

iii, R_1

instructor	cname
------------	-------

R_2

instructor	sname
------------	-------

But the third decomposition $D = \{R_1, R_2\}$ is in BCNF and also satisfies non-additive join property means if join is performed on R_1, R_2 based on instructor name, spurious tuples will not be generated.

First and second decompositions are also in BCNF, they will generate spurious tuples on join condition.

Algorithm: To decompose a non-BCNF relation into decomposition containing a set of relations where R_i in 'D' is in BCNF and always that decomposition satisfies non-additive join property.

input

output

$$D = \{R_1, R_2, R_n\}$$

Method.

$$D = \{R\}$$

repeat

choose a relation R

is not in BCNF

C \rightarrow B

Identity a non-trivial FD $X \rightarrow Y$ in
which X is not super key, then
add 'sub-relations' into D as $X \cup Y$ in
one subrelation of $R - (Y)$ into another subrelation
until each R_i in D is in BCNF.

fourth Normal form: It is based on multi-valued dependency (mvd)

MVD: X and Y are two sets subsets of attributes in R
mvd X (multidetermines) $\rightarrow\!\!\!\rightarrow Y$ holds a valid in a
relation 'state' r of R , if there exists two tuples t_1 &
 t_2 with $t_1[X] = t_2[X]$ then there exists another two
tuples t_3 and t_4 such that

$$i, t_3[X] = t_4[X] = t_2[X] = t_1[X]$$

$$ii, t_3[Y] = t_1[Y] \neq t_4[Y] = t_2[Y]$$

$$iii, t_3[Z] = t_2[Z] \neq t_4[Z] = t_1[Z]$$

where $Z = R - (X \cup Y)$

Trivial mvd if $X \cup Y = R$

Here the combination of all attributes is a super key.

Inference Rule:

IR1: Complimentation Rule for MVD's:

$$\text{If } x \rightarrow\!\!\! \rightarrow y \models x \rightarrow\!\!\! \rightarrow R - (x \cup y)$$

IR2: Augmentation Rule for MVD's:

If $x \rightarrow\!\!\! \rightarrow y$ and $w \geq z$ then $wx \rightarrow\!\!\! \rightarrow yz$

IR3: Transitive Rule for MVD's:

If $\{x \rightarrow\!\!\! \rightarrow y, y \rightarrow\!\!\! \rightarrow z\}$ then $\models x \rightarrow\!\!\! \rightarrow z - y$

IR4: Replication Rule for FD's & MVD's:

If $x \rightarrow y$ is a valid FD then $x \rightarrow y$ & vice versa.

IR5: Coalescence rule for FD's & MVD's:

If $x \rightarrow\!\!\! \rightarrow y$ and there exists

i. $w \cap y \neq \emptyset$

ii. $w \rightarrow z$

iii. $y \geq z$ then $x \rightarrow z$ is valid FD.

→ A relation R is said to be in fourth normal form based on a set of FD's and MVD's for every non-trivial MVD $x \rightarrow\!\!\! \rightarrow y$

Cmp:-

	ename	pname	dname (dependents)
	smith	product x	anna
	smith	product y	lena
off	smith	product x	lena
off	smith	product y	anna

$$F = \{ \text{ename} \rightarrow \text{pname}, \\ \text{ename} \rightarrow \text{dname} \}$$

Relation Emp is in BCNF not in 4NF because all attributes are key attributes in Relation R. The following procedure is adopted to convert non-4NF relation R into sub relations R_1, R_2, \dots to be included into decomposition D.

$$1. D = R$$

2. Choose any relation Q in D not in 4NF

3. While there is a relation Q not in 4NF,

{

choose MVD $\rightarrow Y$ in Q, violates 4NF rule

then

$$R_1 = XUY, R_2 = Q - Y$$

add R_1, R_2 of D

}

The above relation is decomposed in 4NF relations

$$R_1 = \boxed{\text{ename} \mid \text{pname}}$$

smith product x

smith product y

$$R_2 = \boxed{\text{ename} \mid \text{dname}}$$

smith anna

smith lena

$$F_1 = \{ \text{ename} \rightarrow \text{pname} \}$$

$$F_2 = \{ \text{ename} \rightarrow \text{dname} \}$$

so F_1 & F_2 are in 4NF

$(\text{ename} \cup \text{pname}) = R$ - so trivial no need to check for super key

Fifth Normal Form:

- Based on concept of Join Dependency (JD).
 - JD of relation state in of R has to satisfy non-additive join property
- for: JD (R_1, R_2, \dots, R_n) such that
- * $(\Pi_{R_1}(r), \Pi_{R_2}(r) \dots \Pi_{R_n}(r)) = r(R)$
- (natural join)
- JD (R_1, R_2, \dots, R_n) is trivial if one of R_i is R otherwise non-trivial.
 - A relation R is said to be in 5NF based on FD's, MVD's and JD's for each non-trivial JD of (R_1, R_2, \dots, R_n) each R_i must be super key.

Supply

sname	partname	projectname
smith	bolts	product x
smith	nuts	product y
scott	bolts	product x
scott	bolts	product y

The supply relation has no FD's and no MND's if it has all key attributes hence we say it is in BCNF not in 5NF. The above supply relation can be decomposed into subrelations:

$$D = \{ R_1, R_2, R_3 \}$$

with a 1 constraint each suppliers supplies parts, each supplier s will supply to project, each project will have atleast one part.

$$R_1 = \boxed{\text{sname} | \text{partname}}$$

smith bolts

smith nuts

scott bolts

$$R_2 = \boxed{\text{sname} | \text{projname}}$$

smith product x

smith product y

scott product x

scott product y

$$R_3 = \boxed{\text{partname} | \text{projectname}}$$

bolts product x

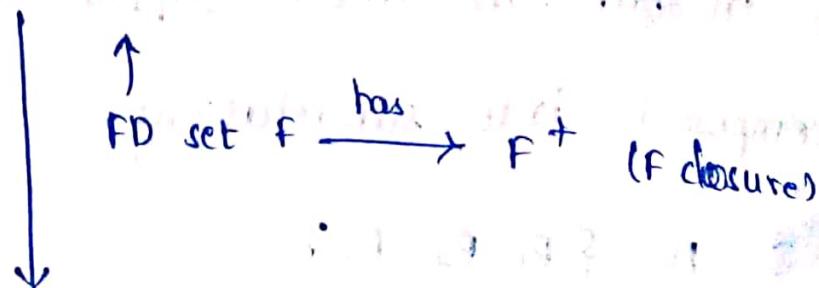
nuts product y

bolts product x

Each relation R_i in decomposition is in 5NF because if natural join is performed $\text{JD}(R_1, R_2, R_3) = r(R)$

Properties of Relational Decomposition

Universal Relation $R(A_1, A_2, \dots, A_n) \xleftarrow{\text{holds}} r(R)$



Decomposition

$D = \{R_1, R_2, R_3, \dots, R_n\}$ how to satisfy

- ① Attribute preservation property - $\bigcup_{i=1}^n R_i^a = R$
- ② FD preservation property

$$(\Pi_{R_1}(f), \Pi_{R_2}(f), \dots, \Pi_{R_n}(f))^+ = f^+$$

- ③ Non-additive join property or loss-less join property

$$*(\Pi_{R_1}(r), \Pi_{R_2}(r), \dots, \Pi_{R_n}(r)) = r(R)$$

Algorithm to test Non-additive Join Property:

Input: Universal Relation R, decomposition

$$D = \{R_1, R_2, \dots, R_n\}$$

Output: To test if D satisfies non-additive join property or not

Method: ① Create initial matrix S with a separate row 'i' for each R_i in D and a separate column 'j' for each attribute A_j in R.

② Initially set $S_{ij} = b_{ij}$,

③ for each row i representing R_i , for each attribute A_j in R { if R_i contains the attributes then set $S_{ij} = a_j$ for all the attributes in R_i in comparison with A_j of R }

④ Repeat the following until no changes are made with a matrix S.

for each FD $X \rightarrow Y$ in QF of R

{ if a row 'i' contains X as LHS attribute and Y with A_j in all the rows set with $S_{ij} = a_j$ for Y }

⑤ If any row is made with completely a_j 's symbols then D satisfies non-additive join property means if join is done on relations in D will not

generate spurious tuples.

⑥ If no row is completely made with symbols then D do not satisfies loss-less join pro

Example:

$$R = \{\underline{ssn}, \underline{ename}, \underline{pno}, \underline{pname}, \underline{ploc}, \underline{hours}\}$$

$$F = \{ \underline{ssn} \rightarrow \underline{ename},$$

$$\underline{pno} \rightarrow \underline{pname}, \underline{ploc}$$

$$(\underline{ssn}, \underline{pno}) \rightarrow \underline{hours} \}$$

$$D = \{ R_1, R_2, R_3 \}$$

$$R_1 = \{ \underline{ssn}, \underline{ename} \}$$

$$R_2 = \{ \underline{pno}, \underline{pname}, \underline{ploc} \}$$

$$R_3 = \{ \underline{ssn}, \underline{pno}, \underline{hours} \}$$

$$S = \begin{array}{c} \text{ssn} & \text{ename} & \text{pno} & \text{pname} & \text{ploc} & \text{hours} \\ \hline R_{1-1} & b_{11} & b_{12} & b_{13} & b_{14} & b_{15} & b_{16} \\ R_{2-2} & b_{21} & b_{22} & b_{23} & b_{24} & b_{25} & b_{26} \\ R_{3-3} & b_{31} & b_{32} & b_{33} & b_{34} & b_{35} & b_{36} \end{array}$$

$$S = \begin{bmatrix} a_1 & a_2 & b_{13} & b_{14} & b_{15} & b_{16} \\ b_{21} & b_{22} & a_3 & a_4 & a_5 & b_{26} \\ a_1 & b_{22} & a_3 & b_{34} & b_{35} & a_6 \\ a_2 & & a_4 & a_4 & a_5 & \end{bmatrix}$$

Row 3: In S is completely with symbols 'a'. So D satisfies loss-less join property.

→ If D contains only 2 relations $D = \{R_1, R_2\}$:

R has F and R_1 has f_1 & R_2 has f_2 if

i, $(R_1 \cap R_2) \rightarrow (R_1 - R_2)$ or

ii, $(R_1 \cap R_2) \rightarrow (R_2 - R_1)$ if, in F^+

then D satisfies NJP

→ If R is decomposed into $D = \{R_1, R_2, \dots, R_i, \dots, R_n\}$

and R_i is decomposed into $D_i = \{Q_1, Q_2, \dots, Q_k\}$

if D satisfies NJP and D_i satisfies LJP then

$D_2 = \{R_1, R_2, \dots, R_{i-1}, Q_1, Q_2, \dots, Q_k, R_{i+1}, \dots, R_n\}$ will

also satisfy LJP or NJP.

Relational DataBase Design Algorithms

1. functional Dependency Preserving 3NF decomposition

2. FDP & NJP 3NF decomposition

3. NJP BCNF Decomposition (near BCNF normal form)

① ILP - Relation R , FD F

Output - $D = \{R_1, R_2, \dots, R_n\}$ each R_i in 3NF and

D will satisfy FDP

method:

① Find minimal set G for F

② If $x \rightarrow A_1, x \rightarrow A_2, \dots, x \rightarrow A_n$ in g then create R_i in D with $\{x \cup A_1, A_2, \dots, A_n\} = g$

③ If any attribute A_j in R is not available in any R_i then include into relevant sub relation.

Ex: $R = \{ ssn, ename, salary, deptno, pno, pname, ploc \}$

$F = \{ ssn \rightarrow ename, salary, deptno,$
 $pno \rightarrow pname, ploc$

$\{ ssn, pno \rightarrow \{ salary, ename \} \}$

$G = \{ ssn \rightarrow ename,$

$ssn \rightarrow salary,$

$ssn \rightarrow deptno,$

$pno \rightarrow pname,$

$pno \rightarrow ploc,$

redundant $\left\{ \begin{array}{l} (ssn, pno) \rightarrow salary, \text{ if we remove } pno \\ (ssn, pno) \rightarrow ename, \text{ if we remove } pno \\ (ssn, pno) \rightarrow pname, \text{ if we remove } ssn \end{array} \right.$

Minimal set $G = \{ ssn \rightarrow ename,$
 $ssn \rightarrow salary,$
 $ssn \rightarrow deptno,$
 $pno \rightarrow pname,$
 $pno \rightarrow ploc \}$

So $R_1 = \{ ssn, ename, sal, deptno \}$

$R_2 \rightarrow \{ pno, pname, ploc \}$

$D = \{ R_1, R_2 \}$ each R_i is in 3NF and always

statisfy
N.F.

④ Here D satisfies FD preservation but not NNF so,

→ Algorithm: Step ① & ② are same.

③ if none of R_i in D contains key attributes then apply key identification algorithm identify set of key attributes and form a relation

④ If any attribute A_j in R is not available in any R_i then include in relevant sub relation

→ Null Values And Dangling Tuples; Alternative Relational DB design with NF's :-

Null Values:

Emp

ssn	ename	sal	sex	dno
12345	xx	10000	m	,
12346	yy	20000	f	2
12347	zz	15000	f	null

dept

dno	dname	dmgrssn
1	R	12345
2	A	12346

If we perform join on two relations emp, dept during natural join then the tuple will be missed.

During aggregate functions also we get inappropriate answers.

Dangling Tuples:

emp1

ssn	superssn	bdate	deptno	ssn	dno
12345	12345	01-dec-1980	1	12345	1
12346	12345	10-dec-1981	12346	2	2
12347	12346	10-dec-1982	12347	null	

Select * from emp, emp1 where emp.ssn = emp1.ssn / dno
 \rightarrow kle get 3 rows

emp1

ssn	dno
12345	1
12346	2
12347	null

if we do not want 12347 since no dept
then select * from emp, emp1 where emp.dno = emp1.dno
 \rightarrow kle, get 2 rows

So the tuple 12347 appears some times and
not appears some times. This is called as
Dangling Tuple:

Ch-2

Introduction to Query Processing & Optimization

Steps involved in processing high level language Query:

Query



scanning, parsing, validate

↓ Intermediate form of query

Query optimization

↓ executive plan

Query code generator

↓ code to execute query

Run time DB processor



Result

Scanning - Table names are available or not - If available
specified attributes are available or not
key words are scanned

Parsing - Given query is checked if semantically correct
or not (Syntax checker)

If the query is valid then it is sent to
Query optimizer module then exact tree/plan structure
is selected to execute in minimal time.

→ Translating SQL Queries into Relational Algebra Expressions

SQL → select ssn, ename
from employee
where sal > (select max(sal)
from employee e
where e.dno = 5);

RA → ① \exists $(\sigma_{dno=5}(\text{employee})) \rightarrow T$
intermediate {
query ② $\Pi_{ssn, name} (\sigma_{sal > T}(\text{employee}))$
form formed after
scanning, parsing &
validation

Algorithm for External Sorting:

External sorting is a sort-merge approach used to sort various DB files available in a large disk or secondary memory. It has 2 phases

1. Sorting phase:-

$$i \leftarrow 1$$

$$j \leftarrow b \quad (\text{size of DB file in blocks})$$

$$k \leftarrow n_B \quad (\text{size of buffer in main memory in blocks})$$

$$m \leftarrow m \leftarrow [j/k] \text{ cell}$$

while ($\ell \leq m$)

do {

 read k blocks of data into buffer or less
 than k

 sort records (internal sort)

$i \leftarrow \ell + 1$

}

2. Merging Phases

$i \leftarrow 1$

$p \leftarrow \log_{k-1} m$

$j \leftarrow m$

 while ($i \leq p$)

 do

 {

$n \leftarrow 1$

$a \leftarrow \lceil j/k \rceil$

 while ($n \leq a$)

 do

 {

 read $k-1$ subfiles, one block at a time

 merge and write as new file

$n \leftarrow n+1$

 }

$j \leftarrow a$

$i \leftarrow i+1$

}

Algorithm for SELECT Operation:

Q₁ - σ (emp)

scr=123456789

Q₂ - σ (emp)

dno=5 or dno=1

Q₃ - σ (emp)

sal<15000 and dno=5

Q₄ - σ (emp)

sal>30000

Various Algorithms for 'select' operation:

i. Based on simple condition

S₁ - linear search

S₂ - Binary search

S_{3a} - using primary index method.

S_{3b} - using hash key

S₄ - using primary index to retrieve multiple records

S₅ - using a clustered index to retrieve multiple records.

S₆ - using a secondary tree index (B+tree)
(for equality comparison of tuples)

ii) Complex selection:

S₇ :- Conjunctive selection

using an individual index

88: Conjugitive selection using composite index.

89: Conjugitive selection for intersection of record pointers

Algorithms to implement join operation:-

Relation R \rightarrow n number of tuples

Relation S \rightarrow m number of tuples

$R \bowtie S$ (RA notation)

$A=B$

→ natural or equi join

Select * from R, S where $R.A = S.B$ (SQL)

1. Nested loop join method:

2. Simple loop join method

3. Partition based hash join method → ① partition phase

4. Sort-merge join method. ② probing phase

sort the tuples in R & S on attributes A & B

set $i \leftarrow 1, j \leftarrow 1$

$T \leftarrow R \bowtie S$

while ($i \leq n$) and ($j \leq m$)

do {

 if $R[i].A > S[j].B$ then

 { output $(R[i], S[j])$ to T

$j \leftarrow j + 1$

}

 else if $(R[i].A < S[j].B)$ then

 { output $(R[i], S[j])$ to T

$i \leftarrow i + 1$

}

else

set $i \leftarrow j + 1$

S

($* R[i].A = S[j].B$)

output the concatenated tuple $\langle R[i], S[j] \rangle$

to T

set $j \leftarrow j + 1$

while $j \leq m$ and $R[i].A \neq S[j].B$

do {

Output the combined tuple

$\langle R[i], S[j] \rangle$ to T

$i \leftarrow i + 1$

y

$k \leftarrow i + 1$

while $k \leq n$ and $R[i].A = S[j].B$

do {

Output and $\langle R[i], S[j] \rangle$ to T

$k \leftarrow k + 1$

}

$i \leftarrow k, j \leftarrow l$

}

3

Algorithm for PROJECT Operation:

$T \leftarrow \Pi_{ssn}(\text{temp})$

$T_1 \leftarrow \Pi_{\text{salary}}(\text{temp})$

i.e $\boxed{T \leftarrow \Pi_A(R)}$

Implement Aggregate functions:

select max(sal)

from employee

$T \leftarrow \gamma_{\max, \text{sal}}(\text{temp})$

Algorithm for set operations:

$T \leftarrow R \cup S$

$T \leftarrow R \cap S$

$T \leftarrow R - S$

Using Heuristics in Query optimization:

Query tree or Query graph

select ssn, fname, dno, dname

from emp, dept

where emp.dno = dept.dno and dno=5;

