

Introduction

- Data today is the most valuable resources.
- Database is a collection of interrelated data.
- Database Management System is a software which helps us manage the data.
- Database System = Database + DBMS
- During 1960's first DBMS came into picture which used N/w model.
- IBM then started using Hierarchical model.
- Later E. Codd came up with "Relational Model".
- In Relational Model data is stored in tables, so it is a very natural & intuitive model.
- Although many other models have emerged & are used since, Relational Model is still the dominant model after almost 50 years.
- We store our data (small) using File System, but to store large amount of data file system is inefficient.
 - One of the major problems with File Systems is Redundancy.
 - Another is inconsistency. File Systems aren't built for concurrent access, and this can lead to a lot of problems. FS doesn't provide a lot of features relating to restricting data access.
 - FS can't force constraints on the data stored.
- DBMS can take care of all this and more.

* Database Design Process

When a customer comes with a requirement for a DB, then the steps for its design would be :

(1) Requirement Analysis Process : Gather info.

(2) Convert the req. to some user friendly model (such as ER model) : Conceptual Design.

(3) Logical Design : Convert conceptual design to logical design using some logical model (such as Relational model) (implementation friendly)

(4) Physical Design : Implementation of logical design.

* Just like there are multiple models to represent Regular Languages, each with its own convention, notations, etc.

Similarly, data can be represented using various models, such as Relational, Object-Oriented, Graph, Document, ER, EER, SQL, etc.

DBMS can be based on many models. We study the most famous : Relational DBMS.

classmate

Date _____

Page _____

Relational Model

& Normalization

student for all rows

~~Relation \equiv Table.~~

- Relational Model: Table / Relation based model to represent the database.

E.F. Codd in his research paper of "Relational Model", gave:

- ① How to store the data.
- ② Terminology / definition (Keys, 2NF, 3NF, BCNF, constraints, etc.)
- ③ Query lang. to retrieve data: Relational Algebra &

Relational Calculus

- Later, based on this a practical query language was created by IBM, Structured Query Language (SQL).

Ex. Student

| S-no | Name | DoB | Phno |
|------|------|--------|------|
| 1 | Rax | 31/13 | 0123 |
| 2 | Daz | 01/212 | 4567 |

Student is the Name of the relation.

(Sno, Name, DoB, Phno) are the Columns or Attributes or Properties.

Each row is called Data or Tuple.

- Every attribute in a relation has a domain, which specifies the set of values that any tuple can have for that attribute.

* Domain of Attribute

For any attribute its domain is the set of "atomic" (indevisible) values, that the tuples can have for that attribute.

Ex. Say Name attribute is having Domain : String.

| | | | | |
|--------------|------------|----------|--------------------|--------------------|
| | ... | Name | ... | student of class 6 |
| Student Name | Amit Kumar | roll no. | student of class 6 | |

Can you retrieve just the surname? No. Only the whole name can be retrieved.

Ex. Say PhNo. is Domain : Integer

| | | | |
|--|----------|------|----------|
| | Ph No | ... | |
| | 257, 312 | | fructose |
| | anil | 3-ol | aniline |
| | | | an-E |
| d? No. It's not atomic for the domain of Ph N ex. | 1324 | 5115 | 503 |

So domain is the set of any & all ~~values~~ atomic values that the attribute can take.

* Schema vs Instance

Schema is the structure / outline of the table. In this we

Instance is the snapshot of the table.

Ex. For Student Table:

Student (

S.No: Integer,

Name: String,

DoB: Data,

Ph.No: Integer

) Schema

| Sno. | Name | DOB | PhNo. |
|------|---------|-------|-------|
| 01 | ABC XYZ | 01/12 | 123 |
| 02 | DEF WXY | 31/15 | 456 |
| 03 | GHI JKL | 6/7/8 | 788 |

↑
Instance

so, Schema consists of Table name, column names & column domains.
Instance consists of the tuples of the table.

- Degree (or arity) of a relation is the no. of attributes of its relational schema.

- Cardinality: Number of rows in any given instance of a relation.

* Why call a Table as Relation?

- In mathematics a rel. from set A to B, $R: A \rightarrow B$, is a subset of $A \times B$, i.e., $R \subseteq A \times B$. This is a binary rel'.

- Similarly a ternary rel' can be def'n as: $R \subseteq A \times B \times C$.

- say, A is set of all integers, B is set of all countries, C is set of all players of cricket.

$$A \times B \times C = \{(a, b, c) | a \in A, b \in B, c \in C\}$$

↳ 3-tuple.

In this all possibilities would be there, say, (01, USA, VK).

But we only want specific data which correctly represents

reality.

Say reality is: (a, b, c) means ' a ' is the rank of highest run scorer in ODI, ' b ' & their country ' c '.

So we can def. $R: (a, b, c) \in R$ iff

- We can extend this to n -ary relation: $R \subseteq A_1 \times A_2 \times \dots \times A_n$

$$A_1 \times A_2 \times \dots \times A_n = \{(a_1, a_2, \dots, a_n) \mid a_i \in A_i\}$$

n -tuple

- So, In a table, say,

| Domain: D_1 | D_2 | D_3 | Attribute / Column Header |
|---------------|---------|--------|---------------------------|
| Sr. No. | Name | Ph No. | |
| 01 | ABC-EFG | 0123 | Student Information |

every tuple $r_i \in D_1 \times D_2 \times D_3$

- Hence we call it Relation since it is actually a rel.

① In above example, $\text{Student} \subseteq D_1 \times D_2 \times D_3$.

② Every record $r_i \in D_1 \times D_2 \times D_3$.

③ Every instance I (set of records), $I \subseteq D_1 \times D_2 \times D_3$

- So, table is just a layman term we informally use, it actually is a relation.

Ex. Rel. $R(A_1 : D_1, A_2 : D_2, A_3 : D_3)$

$$|D_1| = 2$$

$$|D_2| = 3$$

$$|D_3| = 3$$

Max cardinality possible for an instance?

$$2 \times 3 \times 3 = 18. \text{ And the instn is } D_1 \times D_2 \times D_3.$$

No. of possible instances?

$$= \text{No. of subsets of } D_1 \times D_2 \times D_3$$

$$\Rightarrow 2^8$$

* Null

Ex. Suppose a student relation P_S to be designed with the

schema:

$\text{Student}(\text{s-name: string, PassportNo: string, Gender: string})$

Now, it is very much possible that a student doesn't have Passport no. Then what do we put here?

We put Null.

So NULL can mean, "not available". (It may/may not be available in future, but that doesn't matter).

Ex. Say we have Employee schema:

$\text{Employee}(\text{Id: string, Name: string, PassportNo: string, SpouseName: string})$

It may happen that an employee doesn't have a Spouse.

So NULL can also mean, "Not Applicable".

- So NULL is used to represent:
 - ① Not Known
 - ② Not applicable
 - ③ Not recorded, no information utilization
- etc.

Note: Every attribute's domain contains NULL. It can/may not be allowed always.

- * As is the property of relations, order doesn't matter & no duplicates allowed.

- So in any DB Relation, order of tuples doesn't matter, and no two rows can be exactly same in Relational Model.

- Order of attributes matters in Relational Model.

- Although NULL is part of domain of every attribute, we can restrict its use for attributes.

* Keys

- The purpose of storing data is to retrieve them later. But how do we optimize this. We want something to be an attribute (or set of attributes) which is unique for every tuple.

So we can query based on this attribute (set of attributes).

Ex. Identity nos, such as Aadhar, PAN, Social Security No, uniquely identify citizen of a country.

- So, a Key uniquely identifies each tuple in a relation.

- SuperKey:** A set of one or more attributes which uniquely identify each record in a relation.

Say S is superkey, and R is rel.

$$\forall s_i, s_j \in S, \exists a_i, a_j \in R : a_i(s_i) \neq a_j(s_j)$$

Ex. Student (s-name: String, State: Str, Dist: Str, PIN: Int, HouseNo: Int)

{sname}, {State}, {Dist}, {PIN}, {HouseNo} can't uniquely identify

{Say HouseNo. is unique in a PIN code}.

So, {sname, State} is not permitted as primary key.

{s={PIN, HouseNo}} can uniquely identify.

Then all supersets of S can also uniquely id. rows.

That's how student table forms a superkey.

- Superset of a superkey is also a Superkey.

- Can we guarantee that every relation has at least one superkey?

Yes. Since every two record must differ in at least one position. So in worst case, set of all attributes will act as Superkey.

~~(candidate for the) attributes still no hard group and not P~~
Ex. In our previous example suppose ~~the~~ the
superkey : {PIN, House No, S-name, State}

If we remove S-name, it is still Superkey.

If we remove state, it is still Superkey.

If we now try to remove either PIN or HNo, it won't be Superkey anymore.

So we can say {PIN, HNo} is a minimal superkey.

~~Minimal superkey~~ Candidate Key: A minimal superkey is called ~~Candidate Key~~ ~~CK~~, then ~~CK~~ candidate key.

Here C-key is not superkey.

Primary Key: From all possible Candidate Key, one of the CK is chosen as Primary Key. (It's a design choice, depends on the data stored in the relation).

Candidate Key [In case of Null Values]: Should identify all non-NUL records uniquely, and must be minimal.

- A candidate key with NULL value can't be chosen as Primary Key.

If PK comprises of multiple attribute, then no attribute of PK can have NULL values.

Any relation in Relational model, must have exactly one PK at any point of time. It can change over time, but must be exactly one of the possible CKs.

Simple CK: Single attribute CK.

Composite CK: Two or more attribute CK.

- Alternate Keys: Out of all CKs one is chosen at a time to be PK. All other remaining CKs are called Alternate Keys.

| Superkey | Candidate Key | Primary Key |
|---------------------------------|----------------------------|---|
| Should uniquely identify tuples | Should uniquely id tuples. | Should uniquely id tuples. |
| Should be minimal | Should be minimal | No attr. in it should have Null value, for any tuple. |

Note: Generally the term Key means Candidate Key.

Ex. Given some rel. $R(a_1, b_1, c_1, d_1, e_1)$; what is the max. no. of superkeys possible?

Q) What is the max. no. of superkeys possible?

$$\sum_{i=1}^n (-1)^{i+1} \binom{n}{i} \times 2^{n-i} \Rightarrow \sum_{i=1}^n (-1)^{i+1} n! \times 2^{n-i}$$

$\Rightarrow 2^n - 1$ (when every attr. is a CK) (when every attr. is a CK)
 Every non-empty selection of attr. is a SK) (including all attrs. other than one)

② What is the max. no. of CK possible?

$$\Rightarrow \binom{n}{m}$$

Ex. For $R(a_1, b_1, c_1, d_1, e_1)$ if $CKs = \{a_1, bc_1, bde_1\}$

No. of bks?

$$\text{For } a = 2^4$$

$$\text{For } bc = 2^3$$

$$\text{For } bde = 2^2$$

$$\text{For } abc = 2^2$$

$$\text{For } abde = 2^3$$

$$\text{For } bcde = 2^2$$

$$\text{For } abcde = 2^0 = 1$$

bks = $(2^4 + 2^3 + 2^2 - 2^2 - 2 - 2 + 1)$

$$\#SKs = 2^4 + 2^3 + 2^2 - 2^2 - 2 - 2 + 1$$

$$\Rightarrow 21$$

Ex. $R(a_1, a_2, \dots, a_n)$. Assume we know $a_1a_2\dots a_n$ is a CK.
 Other CKs are unknown.

(a) Max. no. of SK possible?

(b) Max. no. of CK possible?

proper

① Since $a_1a_2a_3$ is CK, then none of its \uparrow subset can be key. So no. of possible max. SKs = $2^3 - 1 = 7$

$$2^8 - (2^3 - 1) - 1 = 248$$

$$\Rightarrow 256 - 8 = 248$$

② W.R.t. max. no. of CKs is given by ${}^n C_{r_1, r_2, \dots}$ but any $L^{r_1, r_2, \dots}$ len combination ~~extreme case~~ which contains CK $a_1a_2a_3$ can't be CK itself.

So, max possible CKs:

$$= \text{total no. of } L^{r_1, r_2, \dots} - \left({}^5 C_1 \right) - 1$$

↑ { $a_1a_2a_3$ is CK already }

↑ all 4 len combo with a_1, a_2, a_3 already in it.

$$\Rightarrow 8C_4 - {}^5 C_1 + 1$$

* Database Instance: Collection of all relation's instance in database.

* SQL: A language for querying on RDBMS software.

Different vendors have their own variants.

So we follow the standard SQL-92.

SQL has two parts:

Deals w/ Schema \leftarrow ① Data Defⁿ Language (DDL): Used to define/ create/deletion of relation.

Deals w/ Data. \leftarrow ② Data Manipulation Language (DML): Used for various purpose, but we only study Query Language which is used for retrieval.

- Although SQL & Relational Model appear to be very similar, Relational Model is a theoretical model, i.e., it is defined on paper and it is mathematical ~~model~~ model (i.e. it has inflexible set of rules). SQL is loosely based on RM, and is a practical language and can vary a lot from Relational model.

| SQL | RM |
|---|--|
| <ul style="list-style-type: none"> Treats table as multiset of row records, i.e. duplicate rows allowed. | <ul style="list-style-type: none"> Treats relation as set of rows / records, i.e. duplicates not allowed. |
| <ul style="list-style-type: none"> Primary key not mandatory. | <ul style="list-style-type: none"> PK mandatory. |

* Integrity Constraints over Relations.

- Domain Integrity Constraint: For a attribute in the relation what ~~is~~ the set of all possible values that the attribute can take?

- Key Integrity Constraint: For any two rows in the relation the attribute a_i must be different. (a_i would be valid CK)

- Entity Integrity Constraint: For any row in the relation the attr. a_i must be unique & not null. (This would also make a_i a valid Primary Key)

- Ex. Suppose the following rel.

| SID | SName | District | State |
|-----|-------|-----------|-----------|
| 101 | ABC | Kanpur | UP |
| 102 | DEF | Bangalore | Karnataka |
| 103 | GHI | Hyderabad | Telangana |
| 104 | JKL | Kanpur | Bihar |

It is only logical to have the unique State for a district.
Say for Kanpur, we should only have UP, and not Bihar.

How can we implement such constraint.

None of above constraints talk about multi attr. constraints.

We can observe here: District determines state or

District implies State attr or

State is dependent on district.

Although State doesn't determine District.

Ex. Patna being district means state is definitely Bihar.

But Bihar being state doesn't necessarily mean Patna district.

- This type of Integrity Constraint is called Functional Dependency constraint. This is between attributes (within a table or b/w tables).

| Ex. Student-Course | | Course | | |
|--------------------|-------|--------|--------|--------|
| S-id | C-id | C-id | C-name | Inst. |
| 101 | CS120 | CS120 | OS | Vishal |
| 102 | CS122 | CS121 | DBMST | Deepak |
| 103 | CS125 | CS122 | Algo | Sachin |
| 104 | Null | CS123 | DM | Deepak |

Here for S-id 103, CS125 is there, but that course isn't available. For S-id 104, C-id is Null, which is fine, they might not have enrolled in any course.

So we can say, all the non-null values of Student-Course relation's C-id must come from Course's C-id.

This type of IC is called Referential IC. (Foreign Key constraint).

It can be within a table or within two tables.

* Domain Integrity Constraint

The relational schema specifies the domain of each attribute in the relation. Each instance of the relation for the attributes must have values from the specified domains.

Sometimes additional constraints are required on the domain.

Ex. Student (student_id, name, age)

s-id int,

Age int

But this doesn't prevent neg. value, and ensure non age value.

SQL-92 does allow this using CHECKS.

(In Relational Model, this is implicit.)

* Key Integrity Constraint (Uniqueness Constraint)

If α is CK then two tuples can't have same non-null value of α .

(In SQL this is achieved using "Unique" Keyword)

* Entity Integrity Constraint (Primary Key Constraint)

In Relational Model every table must have PK.

If α is PK then no tuple can have null or duplicate value for α .

With this constraint the aim is to identify every entity uniquely.

* Referential Integrity Constraint (Foreign Key Constraint)

An attribute x referring to some other attribute y must have values which are currently available there in y .

Ex. Student (s-id, Name, Cid)
Course (Cid, Cname, Instructor)

Student.Cid is referencing Course.Cid.

So, it should only take values which Course.Cid currently has.

Student here is called Child table & Course is Parent table based on ~~s~~ Student.Cid attribute.

Student.Cid is foreign key in Student table, because it is pointing to a key in a foreign table.

Ex. student (cid, name, cid) Course (id, name, monitor)

This is also possible.

(student) \rightarrow (course) student.monitor \rightarrow course.id

Ex. Employee (id, name, manager)

- We can't allow a non-key attribute to be referenced by some foreign key.
- So, foreign key always references CK.

- Foreign key may/may not be CK in its table but it always references a CK.

- It may/may not be null. Depends on design, and requirement.
- FK can consist of more than one attr. as well.

Ex. $R(A \times B \times C)$ (A AND B is CK) $T(x \times y \times z)$

So, $y \times z$ is the FK referring to CK AB.
(It's not same as y referring to A & z referring to B)

- FK can also point to CK within same table.

Ex. Suppose the following DB instance:

| A: | a | b | c | B: | x | y | z |
|-------|---|---|---|-------|---|---|---|
| r_0 | 2 | 1 | 3 | r_0 | 2 | 1 | 3 |
| r_1 | 1 | 2 | 3 | r_1 | 1 | 2 | 1 |
| r_2 | 9 | 2 | 2 | r_2 | 3 | 2 | 1 |

c references $x/1$ is direct

- If $A.r_0$ is deleted, is there problem? No.
- If $B.r_2$ is deleted, is there problem? Yes. $A.r_0$ & $A.r_1$ would be referencing to invalid value 3. So violates Ref. IC.
- If a row $(5, 1, 9)$ is added to A? Yes. 9 is not in B.x.
- If a row $(5, 1, 9)$ is added to B? No problem.

Violates Ref. I.C.

- In SQL-92, FK can be specified while creating table.

Ex. Enrolled (sid CHAR(20), cid CHAR(20),

grade CHAR(10),

PRIMARY KEY (sid, cid),

FOREIGN KEY (sid) REFERENCES Student)

{ suppose Student (sid: CHAR(20), sname)

Since we haven't specified which CK from Student table is to be referenced, by default PK would be used, i.e. Student.sid.

Note: In SQL-92, all such things as, PK, FK, domain constraint, key constraint can be specified using CONSTRAINT keyword.

Ex.

CREATE TABLE User (id: INT,

name: VARCHAR(30),

email: VARCHAR(50),

group: INT,

CONSTRAINT user-pk PRIMARY KEY(id),

CONSTRAINT group-fk FOREIGN KEY(group) REFERENCES Group,

CONSTRAINT email-key UNIQUE(email),

CONSTRAINT id-domain CHECK (id > 0),

Keyword name definition

If constraint is violated, error would show with constraint name.

- Violation of Referential IC happens when it is violated.
Suppose R.a references S.y

- ① Insertion to R.
- ② Deletion from S
- ③ Updation to R
- ④ Updation to S

- If violation happens when deleting/updating in referenced table then the possible actions are:

- (1) Cascade: Delete all such records in referencing table which referred to the particular record.
- (2) Set Null: Set the referencing field to null, when referenced value is deleted.
- (3) Set Default: Set some default value.
- (4) Restrict: Prohibit the violating operation.
(No Action in SQL-92)

In SQL-92 we can specify ON DELETE ACTION, for each reference.

Ex. FOREIGN KEY(sid) references student ON DELETE CASCADE

Note: In Schema representation of a relation attribute underlined means PK.

Ex. Enrollment (sid, cid, Grade)

So, sid cid is PK.

* Violation of IC by operations.

| Constraint | Operation | Violated | Remark |
|-------------|-----------|----------|-------------------------------|
| Domain | Read | No | out of constraint |
| | Add | Yes | out of domain value added |
| | Delete | No | out of constraint |
| | Update | Yes | out of constraint |
| Key | Read | No | out of constraint |
| | Add | Yes | Duplicate value |
| | Delete | No | old value will not be deleted |
| | Update | Yes | — |
| Entity | Read | No | out of constraint |
| | Add | Yes | Null or Duple. Value |
| | Delete | No | — |
| | Update | Yes | — |
| Referential | Read | No | out of constraint |
| | Add | Yes | In Referencing Relation |
| | Delete | Yes | In Referenced Relation |
| | Update | Yes | In both relation. |

* Functional Dependencies

Ex. We already saw an example earlier where:

District determined state OR

State was dependent on District.

i.e.

$\text{District} \rightarrow \text{State}$.

(similar to implication in logic.)

- The functional dependency $\alpha \rightarrow \beta$ on a relation R means,

$$\forall T_1, T_2 \in R \quad T_1. \alpha = T_2. \alpha \rightarrow T_1. \beta = T_2. \beta$$

i.e. for all rows T_1, T_2 in R ($T_1 \neq T_2$) if value of α is same for them then β value must also be same.

where, $\alpha, \beta \subset R$. Attributes.

- Similar to other constraints the DB owner would specify the FDs in the DB.

For example the owner may want District to determine state, or may want otherwise.

Ex. $\{\text{Aadhar No}\} \rightarrow \{\text{PAN No}\}$

also $\{\text{PAN} \rightarrow \text{Aadhar}\}$

Ex. $\{\text{Grade}\} \rightarrow \{\text{IsPass}\}$

but $\{\text{IsPass}\} \not\rightarrow \{\text{Grade}\}$

(Note: Think of $\alpha \rightarrow \beta$, as, if you know α , you should uniquely know β automatically)

Ex. $\{ \text{Class}, \text{Roll} \} \rightarrow \{\text{StudentName}\}$

(or $\text{CR} \rightarrow \text{N}$)

Here $\text{CR} \rightarrow \text{Name}$, but

Class by itself doesn't imply Name.

& Roll by itself also

(by design & also logically)

Ex. $R: A \quad B \quad C \quad D$

| | | | |
|---|---|---|---|
| 2 | 2 | 5 | 1 |
| 2 | 2 | 5 | 1 |
| 2 | 3 | 4 | 2 |
| 2 | 3 | 4 | 1 |

(A) $\text{AB} \rightarrow \text{C}$

(C) $\text{A} \rightarrow \text{C}$

(B) $\text{AB} \rightarrow \text{D}$

(D) $\text{B} \rightarrow \text{C}$

In relation R:

(A) $\text{AB} \rightarrow \text{C}$

(C) $\text{A} \rightarrow \text{C}$

(B) $\text{AB} \rightarrow \text{D}$

(D) $\text{B} \rightarrow \text{C}$

B, C, D are false.

A may/may not be true. For an ~~partial~~ FD to hold,

it must hold for all instances of relation.

For it to violate, it only needs one tuple.

- $\forall \alpha \subseteq R_{\text{Attrs}}, \alpha \neq \emptyset \quad \leftarrow \alpha \rightarrow \alpha \text{ (Trivial).}$

- FD is defined on relation schema R , (not on instance), so for all instances $r(R)$ of R , it must hold.

For an FD to be violated a single counterexample row is enough.

(So from a given instance $r(R)$ of relation R , we can comment on the FDs that don't hold for R , but not the ones that hold for R , except trivial ones.)

False antecedent, no FD

True antecedent

Ex. Consider $R(A, B, C, D, E)$ with FDs:

$A, B \rightarrow C$

$C, D \rightarrow E$

Suppose there are atmost 3 diff. values for each $A, B \& D$. What is max. possible no. of diff. values of E ?

Suppose for 3×3 comb. of AB we have all unique C .

$a_1 b_1 c_{11}$

$a_2 b_2 c_{22}$

$a_3 b_3 c_{33}$

$a_1 b_2 c_{12}$

$a_2 b_1 c_{21}$

$a_3 b_2 c_{32}$

$a_1 b_3 c_{13}$

$a_2 b_1 c_{23}$

$a_3 b_1 c_{31}$

$a_1 b_2 c_{11}$

$a_2 b_2 c_{21}$

$a_3 b_2 c_{31}$

$a_1 b_3 c_{12}$

$a_2 b_3 c_{22}$

$a_3 b_3 c_{32}$

So, 9×3 comb. of CD , & hence we'll have 27 unique E s.

So, 27 is max possible no. of diff. values of E .

Note: The concept of FDs were originally defined with no Null values.

• $\forall x, y, z \in R, A \in \mathcal{R} \quad (x \rightarrow y \wedge y \rightarrow z) \rightarrow (x \rightarrow z)$

$x, y, z \neq \emptyset$

[Transitive Rule]

• $\forall x, y, z \in R, A \in \mathcal{R}$ say $x = \{x_1, x_2, \dots, x_n\}, y = \{y_1, y_2, \dots, y_m\}$

① If $x \rightarrow y$ then $x \rightarrow y_1 \wedge x \rightarrow y_2 \wedge \dots \wedge x \rightarrow y_m$
 then, $x \rightarrow y_1 \wedge \dots \wedge x \rightarrow y_m$

$x \rightarrow y_2 \wedge$

$x \rightarrow y_3 \wedge$

\vdots
 $x \rightarrow y_n \wedge$

[Splitting / Combination Rule]

Basically on RHS of FD we can split (when more than one attr. on RHS) and combine (if LHS same)

We can't split on LHS.

Counterexample:

| A | B | C | |
|----------------|----------------|----------------|--|
| a ₁ | b ₁ | c ₁ | |
| a ₁ | b ₁ | c ₁ | |
| a ₁ | b ₂ | c ₂ | |
| a ₂ | b ₁ | c ₂ | |

Here, $AB \rightarrow C$ but neither $A \rightarrow C$, nor $B \rightarrow C$.

We can combine on LHS.

e.g., $(A \rightarrow C) \wedge (B \rightarrow C) \rightarrow (AB \rightarrow C)$

Proof: Say $A \rightarrow C \wedge B \rightarrow C$, but $\nmid A B \rightarrow C$ (a contradiction)

$A \quad B \quad (A \rightarrow C) \wedge (B \rightarrow C)$

$a \quad b \quad c_1 \quad c_2$ (a contradiction for the result)

$a \quad b \quad c_2$ (a contradiction)

but then $\nmid A \rightarrow C \wedge B \rightarrow C$ here \checkmark

\therefore Contradiction.

• Trivial FDs: FDs that always hold irrespective of the Relation.

For $\alpha, \beta \subseteq R.\text{Atts}$, $\alpha, \beta \neq \emptyset$

always holds if $\alpha \subseteq \beta$

$\forall x, y, z \in R.\text{Atts} \quad (x \rightarrow y) \rightarrow (xz \rightarrow yz)$

$x, y, z \neq \emptyset$

Proof: Say $x \rightarrow y$ holds & $\nmid xz \rightarrow yz$ doesn't hold.

$x \rightarrow y \quad ?z$ true by assumption

$x_1, y_1, z_1 \quad \text{①}$

$x_2, y_2, z_2 \quad \text{②}$

but then $x \rightarrow y$ doesn't hold.

\therefore Contradiction.

If $x \rightarrow y$; $wy \rightarrow z$ then $wx \rightarrow z$ { Proof by }

If $x \rightarrow y$; $w \rightarrow z$ then $xw \rightarrow yz$ { contradiction }

(Trick to remember, think of $x \rightarrow y$ as x beats y or X is more powerful than Y)

Closure of Set of Attributes: (Attribute Closure)

For a set of attributes X , its closure is given as:

$$X^+ = \{w \mid w \in R.\text{Atts}, X \rightarrow w\}$$

Ex. Suppose $R(A, B, C, D)$, with

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

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

$$(AD)^+ = (A, D)^+ = \{A, D, B, C\}$$

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

$$(AB)^+ = \{A, B, C\} \quad A \rightarrow B \rightarrow C$$

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

$$(AC)^+ = \{A, B, C\}$$

$$D^+ = \{D, B, C\}$$

$$(BC)^+ = \{B, C\} \quad \text{closed & no FD}$$

$$(BD)^+ = \{B, D, C\}$$

$$(CD)^+ = \{C, D, B\}$$

⋮

$$(XY \subseteq Z) \leftarrow (Y \rightarrow X)$$

Ex. If for a RBS if $X^+ = R.\text{Atts}$, then

(or R)

what can be said about X?

① Superkey

② Candidate Key

③ Primary Key

④ Foreign Key.

$$X^+ = R \text{ means } X \rightarrow R, \text{ p.e}$$

If I know X, we can uniquely identify R.
So, X must be SK.

It might/might not be CK & PK.

∴ if $X \subseteq R$ & $X^+ = R$, then X is SK of R.

X and X' stand X is X → X'. for having uniqueness of R.F.T.
(X must be primary key)

(1) For a relation R, if $X \subseteq R$ & $X^+ = R$, then X is SK of R.

(2) For a relation R, if $X \subseteq R$ & $X^+ = R$ & $\forall y \in X \setminus X^+ \neq R$ then

X is CK of R. (proper subset)

Ex. R(A, B, C, D) $FD = \{A \rightarrow B, B \rightarrow C\}$ Find all CKs of R.

$$A^+ = ABC$$

$$B^+ = BC$$

:

We can iterate, or we can observe.

Q1: A & D never appear on RHS of any FD, so for them to be determined they must be on LHS (trivial). And hence must also be in every CK.

$$(AD)^+ = ADBC$$

$$\therefore (AD)^+ = R$$

$\therefore AD$ is CK.

Now all superset of AD would be SK.

\therefore Only AD is CK. (AD is the only CK for original)

Ex. $\{AB \rightarrow C, AE \rightarrow D, D \rightarrow B\}$ R(A, B, C, D, E)

A, E must be in every key since they aren't on any RHS.

$$(AE)^+ = AE, D, B, C$$

$\therefore AE$ is the only CK.

\therefore All supersets are SK.

• Inferring FDs

$A \rightarrow B, B \rightarrow C \Rightarrow A \rightarrow C$

$A \rightarrow B, B \rightarrow C, C \rightarrow D \Rightarrow A \rightarrow D$

Ex. $R(A, B, C)$ $FD = \{A \rightarrow B, B \rightarrow C\}$

which other FDs can be inferred?

$A \rightarrow A, B \rightarrow B, C \rightarrow C$

$A \rightarrow B, B \rightarrow C \Rightarrow A \rightarrow C$

$A \rightarrow C, B \rightarrow C \Rightarrow AB \rightarrow C$

$AC \rightarrow C, AC \rightarrow A, \dots$ etc.

Ex. $R(A, B, C, D, E, F)$

$FD = \{AB \rightarrow C, BC \rightarrow AD, D \rightarrow E, CF \rightarrow B\}$

which FD can be inferred from above set

(A) $A \rightarrow A$ ✓ (Trivial)

(B) $A \rightarrow BC$ ✗

Check $A^+ - A^+ = A$. $\therefore A \nrightarrow BC$.

(C) $AB \rightarrow BDE$ ✓

$(AB)^+ = ABCDE \therefore AB \rightarrow BDE$

• Closure of FD set: For a FD set F , its closure

F^+ is the set of all FDs that can be inferred

from it.

Ex. $F = \{A \rightarrow B, B \rightarrow C\}$ $R(A, B, C)$

$A^+ = ABC$

So, $A \rightarrow A, A \rightarrow B, A \rightarrow C, A \rightarrow AB, A \rightarrow AC, A \rightarrow BC,$

$A \rightarrow ABC$

$$\therefore 2^3 - 1 = 7$$

$$B^+ = BC$$

$$\# \text{fds for } B = 2^2 - 1 = 3 \quad \# \text{fds} = 2^1 - 1 = 1$$

$$(AB)^+ = ABC$$

$$\# \text{fds} = 2^3 - 1 = 7$$

$$(AC)^+ = ABC$$

$$2^3 - 1 = 7$$

$$(BC)^+ = ABC$$

$$2^2 - 1 = 3$$

~~$$(ABC)^+ = ABC$$~~

$$\# \text{fds} = 2^3 - 1 = 7$$

$$\therefore \# \text{fds in } F^+ = 7 \times 4 + 3 \times 2 + 1 = 35.$$

FD Cover: A FD set F is said to cover another FD set G iff

$$G^+ \subseteq F^+$$

FD set Equivalence: Two FD set F & G are said to be equivalent iff

$$F \text{ Covers } G, \text{ i.e. } G^+ \subseteq F^+$$

$$\text{AND } G \text{ Covers } F, \text{ i.e. } F^+ \subseteq G^+$$

Ex. R(A, B, C, D) $F = \{A \rightarrow B, B \rightarrow C\}$ $G = \{A \rightarrow B, A \rightarrow C\}$

Which set covers which set

Checking for $G^+ \subseteq F^+$ (if F covers G):

$$A^+ = ABC \quad \{ \text{Found } A^+ \text{ using } F, \text{ & all FDs of } G \text{ are covered} \}$$

$$\therefore G^+ \subseteq F^+$$

Checking for $F^+ \subseteq G^+$

$$A^+ = ABC$$

$$B^+ = B$$

$\therefore B \rightarrow C$ not found So, $F^+ \not\subseteq G^+$.

$\therefore F$ covers G .

Ex. $F = \{A \rightarrow B, B \rightarrow ACD, CD \rightarrow B\} \quad R(A, B, C, D)$

$G = \{A \rightarrow CD, B \rightarrow CD, CD \rightarrow AB\} = F \text{ ref. shft}$

Checking $G^+ \subseteq F^+$

$$A^+ = ABCD$$

$$B^+ = BACD$$

$$(CD)^+ = CDBA$$

$\therefore F \text{ covers } G.$

$$F = I^{ABCD} \quad I^{ABCD} = \text{shft}$$

$$I^{ABCD} = \text{shft}$$

Checking $F^+ \subseteq G^+$

$$A^+ = ACD B$$

$$B^+ = BCDA$$

$$(CD)^+ = CDAB$$

$\therefore G \text{ covers } F.$

$F = G$. Hence F is irreducible.

For $A \rightarrow B$, $B \rightarrow C$, $C \rightarrow D$, $D \rightarrow A$ we have (A, B, C, D) .

From $A \rightarrow B$, $B \rightarrow C$, $C \rightarrow D$, $D \rightarrow A$ we have (A, B, C, D) .

From $A \rightarrow B$, $B \rightarrow C$, $C \rightarrow D$, $D \rightarrow A$ we have (A, B, C, D) .

From $A \rightarrow B$, $B \rightarrow C$, $C \rightarrow D$, $D \rightarrow A$ we have (A, B, C, D) .

From $A \rightarrow B$, $B \rightarrow C$, $C \rightarrow D$, $D \rightarrow A$ we have (A, B, C, D) .

From $A \rightarrow B$, $B \rightarrow C$, $C \rightarrow D$, $D \rightarrow A$ we have (A, B, C, D) .

From $A \rightarrow B$, $B \rightarrow C$, $C \rightarrow D$, $D \rightarrow A$ we have (A, B, C, D) .

From $A \rightarrow B$, $B \rightarrow C$, $C \rightarrow D$, $D \rightarrow A$ we have (A, B, C, D) .

From $A \rightarrow B$, $B \rightarrow C$, $C \rightarrow D$, $D \rightarrow A$ we have (A, B, C, D) .

From $A \rightarrow B$, $B \rightarrow C$, $C \rightarrow D$, $D \rightarrow A$ we have (A, B, C, D) .

From $A \rightarrow B$, $B \rightarrow C$, $C \rightarrow D$, $D \rightarrow A$ we have (A, B, C, D) .

From $A \rightarrow B$, $B \rightarrow C$, $C \rightarrow D$, $D \rightarrow A$ we have (A, B, C, D) .

From $A \rightarrow B$, $B \rightarrow C$, $C \rightarrow D$, $D \rightarrow A$ we have (A, B, C, D) .

From $A \rightarrow B$, $B \rightarrow C$, $C \rightarrow D$, $D \rightarrow A$ we have (A, B, C, D) .

From $A \rightarrow B$, $B \rightarrow C$, $C \rightarrow D$, $D \rightarrow A$ we have (A, B, C, D) .

From $A \rightarrow B$, $B \rightarrow C$, $C \rightarrow D$, $D \rightarrow A$ we have (A, B, C, D) .

From $A \rightarrow B$, $B \rightarrow C$, $C \rightarrow D$, $D \rightarrow A$ we have (A, B, C, D) .

- * **Prime Attribute:** In a relation R, an attribute is called Prime if it is member of some Candidate Key.
- * **Non-Prime Attribute:** In a relation R, an attr. which isn't part of any CK is called Non-Prime.

Ex. $R(A, B, C)$ $FD = \{A \rightarrow B, B \rightarrow C\}$

A isn't part of any RHS, so only A can determine A .

So, A is present in every key.

$$A^+ = ABC$$

So A is SK & since its minimal

so A is CK. $\therefore A$ is PK.

\therefore Only CK $\Rightarrow A$.

\therefore Prime Attr: A

Non-Prime Attr: B, C

- * **Trivial:** Something which always holds or never holds.

Trivial FD: In a relation R, an FD $\alpha \rightarrow \beta$, $\alpha, \beta \subseteq R$ is called trivial iff $B \subseteq \alpha$, i.e., B is subset of α . (since it always hold)

Non-Trivial FD: An FD is Non-Trivial, iff its non-trivial.

Ex. $AB \rightarrow BCD$

split RHS

Trivial

Non-Trivial

$\therefore AB \rightarrow BCD \equiv AB \rightarrow CD$.

Note: In further conversation about Normal forms, we will only consider Non-trivial FDs split on RHS with each FD having single attr. on RHS.

* Fully-Dependent FD: A FD $\alpha \rightarrow \beta$ is fully-FD dependent, iff $\forall \text{ non } \alpha - \{\alpha\} \nrightarrow \beta$.

β is fully dependent on α iff any proper subset α' ($\alpha' \subset \alpha$) of α can't determine β .

* Partial FD: A FD $\alpha \rightarrow \beta$ is partial FD iff there is some proper subset of α can determine β , i.e.

$$\exists \gamma \subset \alpha \quad \gamma \rightarrow \beta \quad \text{and } \gamma \neq \alpha$$

Ex. R(A, B, C, D, E) $\text{FD} = \{CD \rightarrow E; D \rightarrow A; A \rightarrow E\}$

① Is $CD \rightarrow E$ partial?

$$C^+ = C \quad D^+ = DA$$

so D can determine E

$\therefore CD \rightarrow E$ is partial because we don't need entire

CD to determine E.

② Is $D \rightarrow E$ partial?

No. we can't remove anything from LHS

③ $A \cdot BD \rightarrow C$. Partial or full? $\Delta \neq \Delta F$ different? $\Delta \neq \Delta$ modular

$BD \rightarrow C$ is not a valid FD on R, since $(BD)^+ = \{B, D, A, E\} \not\models C$

$$C \not\in (BD)^+ \wedge (A \rightarrow Y) \wedge (Y \rightarrow X) \not\Rightarrow C \rightarrow E$$

$\left\{ \begin{array}{l} A \rightarrow Y \\ Y \rightarrow X \end{array} \right\}$

④ Find CKs.

BCD must be part of every CK, as they don't appear on any RHS

$$(BCD)^+ = \{B, C, D, A, E\} \Leftarrow X \leftarrow A \leftarrow Y$$

\therefore Only ~~one~~ $\{B, C, D\}$ is CK. True

Ex. $R(A, B, C, D, E, F)$

$$FDs = \{ABC \rightarrow DEF, DE \rightarrow ABC, ABC \rightarrow F, DE \rightarrow F, AB \rightarrow D, E \rightarrow C\}$$

① $DE \rightarrow C$. Full / Partial

$$D^+ = D \quad E^+ = E \cup C$$

\therefore Partial.

② $ABC \rightarrow D$

$$A^+ = A \quad B^+ = B \quad C^+ = C \quad (AB)^+ = ABD$$

\therefore Partial.

• Diagram of FD



* Transitive FD : A FD $X \rightarrow A$ $\{x \in R, A \in A\}$ in relation R, iff
 $\{x \rightarrow y \rightarrow A\} =^+(A)$ $A \notin X$ $\forall x \in A$ below relation $x \rightarrow A$

& $\exists y \in R (x \rightarrow y) \wedge (y \rightarrow A) \wedge (y \not\rightarrow x) \wedge (A \notin y)$

$\{ \begin{matrix} x \\ \swarrow \uparrow \\ x \rightarrow y \rightarrow A \end{matrix} \}$

Some observations :

① $A \rightarrow X$. \Rightarrow prove $A \rightarrow X$ not transitive.

Proof : Assume $A \rightarrow X$.

$$Y \rightarrow A \rightarrow X \Rightarrow \{Y \rightarrow X\} =^+(A)$$

But $Y \not\rightarrow X$ $\{Y \rightarrow X\} \neq \emptyset$

\therefore Contradiction.

② Every Partial FD is also transitive FDs.

Ex. $R(A, B, C)$ \rightarrow $FD = \{A \rightarrow C\}$

① Is $AB \rightarrow C$ a TD?

Yes. $AB \rightarrow A$ & $A \rightarrow C$ $\therefore A \rightarrow AB \rightarrow C \notin \{A\} \wedge C \notin \{A\} + \{C \in A\}$
 (Also - Partial FD)

Ex. $R(A, B, C, D)$ $FD = \{A \rightarrow B, B \rightarrow C, C \rightarrow D\}$

① Is $A \rightarrow D$ partial?

No.

② Is $A \rightarrow D$ transitive?

$$A \rightarrow C \wedge C \rightarrow D \wedge C \rightarrow A \wedge D \notin \{A\} \wedge D \notin \{C\}$$

OR

$$A \rightarrow C \rightarrow D \wedge D \notin \{A\} \wedge D \notin \{C\}$$

\therefore Yes. $A \rightarrow D$ is transitive.

* Normal Forms:

A relation (relational schema) R , can be in the following

Normal forms: depending upon fulfillment of rules of each

- 1NF } Defⁿ wrt. Domain
- 2NF } Defⁿ wrt Non-Trivial FDs
- 3NF } Functional Dependencies.
- BCNF } w/ RHS containing single attr. (A,B,C,D) \rightarrow A

* First Normal Form: A relation R is in 1NF iff the domains of all attributes of R are atomic.

In Relational Model every relation is in 1NF, since we consider domain of every attr. as atomic.

(So trivial for 1NF. Defined wrt Domain). A \leftarrow B

* Second Normal Form: A relation R is in 2NF iff every non-prime attr. is fully dependent on every Candidate Key.

OR R is in 2NF iff no non-prime attr. is partially dependent on some Candidate Key (or just Key).

So. FD of the form: Some CK \rightarrow Non-prime attr., must not be partial.

Ex. $R(A, B, C, D, E)$. $FD_R = \{A \rightarrow B, BC \rightarrow E, ED \rightarrow A\}$

CD not on any RHS.

$$(CD)^+ = \{C, D\}$$

$$(ACD)^+ = \{A, C, D, B, E\} \quad \therefore CK$$

$$(BCD)^+ = \{B, C, D, E, A\} \quad \therefore CK$$

$$(ECD)^+ = \{E, C, D, A, B\} \quad \therefore CK$$

$$\therefore CK_R = \{ACD, BCD, ECD\}$$

\therefore Prime attr = {A, B, C, D, E}

Nonprime attr = {}.

\therefore Nonprime attr are not there

type of ~~partial~~ partial dependency doesn't exist

$\therefore R$ is in 2NF.

Ex. $R(A, B, C, D)$ FD = { $A \rightarrow C, B \rightarrow A, A \rightarrow D, AD \rightarrow C$ }

B not on RHS.

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

$\therefore CK_B = \{B\}$. contradiction. A cannot depend on B.

Prime Attrs = {B}

Non Prime Attrs = {A, C, D}

All FDs of form $CK \rightarrow NPA$:

$B \rightarrow A, B \rightarrow C, B \rightarrow D$. contradiction. A, C, D are prime attrs.

None of them are partial FDs.

$\therefore R$ is in 2NF.

L

Ex. $R(A, B, C, D)$ FD = { $C \rightarrow D, CD \rightarrow A, AB \rightarrow C, BD \rightarrow A$ }

B not on RHS.

$$B^+ = \{B\}$$

$$(BA)^+ = \{B, A, C, D\}$$

$$\therefore CK_B = \{AB, BC, BD\}.$$

\therefore No non-prime attr.

Hence no $CK \rightarrow NPA$ partial FD.

$\therefore R$ is in 2NF.

$$A^+ = \{A, B, C, D\} = ^+(ABCD)$$

$$B^+ = \{B, C, D\} = ^+(BCD)$$

$$C^+ = \{C, D\} = ^+(CD)$$

Ex. $R(A, B, C, D)$ $FD = \{AB \rightarrow C, AC \rightarrow B, BD \rightarrow A\}$

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

$$(AC)^+ = \{A, C, B, D\}$$

$$(AD)^+ = \{A, D\}$$

$$(BC)^+ = \{B, C\}$$

$$(BD)^+ = \{B, D, A, C\}$$

$$(CD)^+ = \{C, D\}$$

For any other comb. there is a proper subset which is Key.

$$\therefore CKs = \{AB, AC, BD\}$$

\therefore No NPA.

i.e., In 2NF.

$$(AB \rightarrow C) \text{ is not a partial FD as } A \rightarrow C$$

Partial FD for AB.

Ex. $R(A, B, C, D)$

$$FD = \{AB \rightarrow CD, B \rightarrow C\}$$

AB not on RHS.

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

$$\therefore CKs = \{AB\}$$

$$\text{Prime attr} = \{A, B\}$$

$$NPA = \{C, D\}$$

All FDs of CK \rightarrow NPA form:

$$AB \rightarrow C, AB \rightarrow D$$

. $AB \rightarrow C$ is partial FD, $\therefore B \rightarrow C$.

$\therefore R$ is not in 2NF.

So, basically 2NF is violated when some proper subset of some CK can determine some non-prime attr.

Ex. $R(A, B, C, D, E)$

$$FD = \{AB \rightarrow D, C \rightarrow B, E \rightarrow A\}$$

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

$$\therefore CKs = \{CE\}$$

$$PA = \{C, E\}$$

$$NPA = \{A, B, D\}$$

FDS of form CK \rightarrow NPA.

$$CE \rightarrow A, CE \rightarrow B, CE \rightarrow D$$

$\therefore CE \rightarrow A$ is partial FD.

$$(E \rightarrow A)$$

$\therefore R$ isn't in 2NF.

Ex. $R(A, B, C, D, E)$ w.r.t FDs = $\{AC \rightarrow BDE, A \rightarrow B, D \rightarrow E\}$

AC not on RHS.

$$\begin{aligned} AC^+ &= \{A, B\} \\ (AC)^+ &= \{A, B\} \end{aligned}$$

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

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

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

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

$$PA_S = \{A, C\} \quad A \rightarrow B \rightarrow E$$

$$NPA_S = \{B, D, E\} \quad D \rightarrow E$$

$$\{A, C\} = ^+(AC)$$

All FDs of form CK \rightarrow NPA:

$$AC \rightarrow B, \quad AC \rightarrow D, \quad AC \rightarrow E,$$

$\therefore AC \rightarrow B$ is Partial FD ($\because A \rightarrow B$)

$\therefore R$ not in 2NF.

Ex. Can 2NF have partial dependency?

Yes, it is possible.

Ex. $\{AB \rightarrow CD, CD \rightarrow AB, B \rightarrow D\}$

This is in 2NF.

$AB \rightarrow D$ is Partial FD & allowed.

Only CK \rightarrow NPA type of Partial FD disallowed.

Ex. R is in 2NF iff some NPA is fully dependent on all CKs. T/F.

False. All \oplus NPA, not some.

Ex. R is in 2NF iff all NPA are fully dependent on all SuperKeys. T/F?

False. Only on all CKs.

Ex. Can a non-prime attr. determine some prime attr.?

No. Proof: Assume $A \rightarrow B$, A is non-prime & B is prime.

Since B is prime, it is part of some CK, X_B . $X \subseteq R$.

Since $A \rightarrow B$,

$A \rightarrow X \rightarrow B$

So, A X is CK. A is prime. (From above) Hence X

\therefore Contradiction.

Ex. Can a non-prime attr. determine another non-prime attr.?

Yes. Ex. $AB \rightarrow CD$

$C \rightarrow D$

Ex. Can a prime attr. determine another prime attr.?

Yes

$A \rightarrow AB$ (prime) $\rightarrow AB \rightarrow CD$ (prime)

PAIRS of cardinality

Third Normal Form: A relation R is in 3NF iff every non-prime attr. is non-transitively dependent on every Candidate Key.

So 3NF is violated by R, iff some non-prime attr. is transitively dependent on some CK.

3NF hence is violated when: (All ~~non-prime~~ ^{obvious} attrs. are ~~non-prime~~)

$(CK \rightarrow \text{Non superkey} \rightarrow \text{Non Prime Attr}) \wedge (NPA \notin CK)$

$\wedge (NPA \in \text{Non SK})$

So, Such FDS in R cause violation:

Non-Superkey \rightarrow Non-prime attr

$\Rightarrow \text{Ex} \quad (\text{Not a SK } \rightarrow \text{NPA}) \text{ after applying some rule } \rightarrow$
 (Some CK $\xrightarrow{\text{Transitively}} \text{Some NPA}$)
 $\Rightarrow \text{Ex} \quad \text{for 1NF it is needed to have some}$

Proof:

Suppose α is Non SK $\Rightarrow \alpha \text{ is NPA.}$

α can't determine any CK, since it's Non SK.

If $\alpha \rightarrow A$, then

all CK $\rightarrow B$. $B \rightarrow \alpha$.

$\Rightarrow B \rightarrow \alpha \rightarrow A$ after applying some rule
 \Leftarrow
 (since α is Non SK)

$\therefore B \xrightarrow{\text{Transitively}} A.$

Ex. Is $\text{Some NPA} \rightarrow \text{Some other NPA}$ a problem to 3NF?

a problem to 3NF?

Yes. If $A \rightarrow B$, and A is non SK, then it's a problem.

That is, Non SK \rightarrow NPA type of FD.

Also for some CK,

$\text{CK} \rightarrow A \rightarrow B \quad \& \quad B \neq A \wedge B \neq \text{CK}$

Ex. Is NonPrime (A) \rightarrow Prime (B) a problem to 3NF?

a problem to 3NF?

(NonPrime A)

Never possible.

Not a SK



Ex. Is $\text{Non-Prime} + \text{Prime} \xrightarrow{\alpha} \text{Non-Prime}$ a problem
for 3NF?

Yes. $\text{Non-SK} \rightarrow \text{Non-Prime}$

Also,

$$(CK \rightarrow \alpha \rightarrow A) \wedge (A \notin CK, A, A \notin \alpha)$$

So, 3NF violations include FDs of type:

① $NPA \rightarrow NPA$.② Partial dependency of form $CK \rightarrow NPA$, or. } violation of
(Proper subset of CK) $\rightarrow NPA$. } 2NF③ Prime + Non-Prime

Non-SK

Ex. $R(A, B, C, D)$ with FDs = $\{A \rightarrow C, B \rightarrow A, A \rightarrow D, AD \rightarrow C\}$.

Check if in 3NF.

B not on RHS

$$B^+ = B, A, C, D.$$

$$\therefore CK_B = \{B\}$$

$$\therefore PA = \{B\} \text{ and } NPA = \{A, C, D\}$$

A \rightarrow C violates 3NF. ($\text{Non-Prime} \rightarrow \text{Non-Prime}$)

∴ R not in 3NF.

Ex. $R(A, B, C, D)$ with FDs = $\{C \rightarrow D, CD \rightarrow A, AB \rightarrow C, BD \rightarrow A\}$

B not on RHS.

$$B^+ = B,$$

$$(AB)^+ = A, B, C, D$$

$$(CB)^+ = C, B, D, A$$

$$(CD)^+ = D, B, A, C.$$

$$\therefore CK_B = \{AB, CB, DB\}$$

$$\text{Prime} = \{A, B, C, D\}$$

 $\therefore \text{No NPA}$ $\therefore \text{In 3NF.}$

Ex. $R(A, B, C, D)$ $\text{FD} = \{ BD \rightarrow C, AB \rightarrow D, AC \rightarrow B, BD \rightarrow A \}$

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

$$(AC)^+ = \{A, C, B, D\}$$

$$(AD)^+ = \{A, D\}$$

$$(BC)^+ = \{B, C\}$$

$$(BD)^+ = \{B, D, C, A\}$$

$$(CD)^+ = \{C, D\}$$

$C_K = \{AB, AC, BD\}$ (prime attr. by min. dependency)

thus Prime = $\{A, B, C, D\}$ (by modus operandi)

\therefore No NPA

\therefore 3NF.

Ex. In 3NF, do we have transitive dependency?

Yes. Prime Attr. \rightarrow Prime Attr. (by modus operandi)

Ex. $R(A, B, C, D, E)$

$$AB \leftrightarrow CD$$

$$A \rightarrow C$$

Here, $AB \rightarrow C$ is transitive dependency.

Ex. In 3NF, a NPA can be transitively determined by OK?

Yes. There exists a $C_K \rightarrow PA$ FD for all prime attrs.

Ex. In 3NF, a NPA can be ~~transitively~~ determined by another NPA?

$$No. (A, B, C, D) = NPA$$

$$A, B, A \rightarrow^+ (CA)$$

$$A, B, D \rightarrow^+ (AD)$$

$$A, A, D \rightarrow^+ (AD)$$

- A partial dependency implies transitive dependency,

$$\alpha \xrightarrow{\beta} \beta \rightarrow \alpha$$

where $\beta \subset \alpha$.

& β is Non SK (\because Proper subset)

So, any PD of form $\text{CK} \xrightarrow{\text{Par.}} \text{NPA}$ would also be a Transitive Dependency.

Hence $\neg R$ is not in 2NF $\rightarrow R$ is not in 3NF.

$\Rightarrow R$ is in 3NF $\rightarrow R$ is in 2NF

\therefore If a relation R is in 3NF, it is also in 2NF.

Ex. $R(A, B, C, D)$ $\text{FD} = \{B \rightarrow C, AC \rightarrow D, ABD \rightarrow C, BCD \rightarrow A\}$

Is R in 3NF? 2NF?

B not on RHS.

$$B^+ = B, C,$$

$$AB^+ = A, B, C, D \quad \text{for } 3NF \quad \text{for } CK_S = \{AB, DB\}$$

$$CB^+ = C, B \quad \text{for } 2NF \quad \text{for } PAH_S = \{A, B, D\}$$

$$DB^+ = B, D, C, A,$$

$$NPAH_S = \{C\}$$

Not in 3NF.

$$AB \rightarrow B \rightarrow C \quad \text{&} \quad C \notin AB \quad \text{&} \quad C \notin B.$$

\therefore Trans. dependency $AB \rightarrow C$.

Not in 2NF.

$AB \rightarrow C$, but $B \rightarrow C$. \therefore Partial dependency.

\therefore 1NF

Ex. $R(A, B, C, D)$

$$FD = \{AB \rightarrow C, ABD \rightarrow C, ABC \rightarrow D, AC \rightarrow D\}$$

AB not on RHS.

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

$$\therefore CK_S = \{AB\}$$

$$PAH_S = \{A, B\}$$

NOT 3NF. $AB \rightarrow AC \rightarrow D$

$$NPAH_S = \{C, D\}$$

In 2NF.

Ex. R(A, B, C, D) FDs = {ABD ⊆ C, A → B, AB → C, B → A}

D not on RHS, and $A \leftarrow B \leftarrow C \leftarrow D$

$$D^+ = D, AD, BD$$

$$AD^+ = ADB, \therefore CK's = \{AD, BD\}$$

$$BD^+ = BDA, \text{ PATH} = \{A, B, D\}$$

$$CD^+ = CD, \text{ NPATH} = \{C\}$$

Not in 3NF. $\leftarrow AD \rightarrow AB \rightarrow C$ is 3NF

\leftarrow \exists B of \exists A \rightarrow \exists C of \exists B

Not in 2NF. $BD \rightarrow C$, but $B \rightarrow AC$

\exists A of \exists B \rightarrow \exists C of \exists B, which is 2NF

∴ In 1NF

Ex. R(A, B, C, D) FDs = {ACD ⊆ B, AC ⊆ D, D ⊆ C, AC ⊆ B}

A not on RHS

$$A^+ = A, AA, ABA, ABD$$

$$AB^+ = AB, ABB, AABA, AABD$$

$$AC^+ = ACDB, ACD, ACDB$$

$$AD^+ = ADCB, ADC, ADCB$$

$$\text{NPATH} = \{B\}, A, B, C, D = \{BD\}$$

All possible nonSK = {C, D, A, CD}

None determines any NPATH.

∴ no NonSK → NPATH dependency.

∴ In 3NF.

NonSK is determined by the CKs = {BD}

BD is determined by the CKs = {ACD}

ACD is determined by the CKs = {AC}

AC is determined by the CKs = {A}

A is determined by the CKs = {A}

C → D → B → A → C → A → D

relation R, iff

- Boyce-Codd Normal Form: In every non-trivial FD $\alpha \rightarrow \beta$ ($\alpha, \beta \subseteq R$), α ~~is not~~ $\rightarrow \beta$ SuperKey then R is in BCNF.

3NF is saying that NonSK \rightarrow NonPrime shouldn't happen. i.e either LHS is SK or RHS is prime.

BCNF is saying that LHS should be SK.

| | 3NF - A | BCNF - B |
|------------------------------|---------|----------|
| $SK \rightarrow NPA_{Hr}$ | ✓ | ✓ |
| $SK \rightarrow P_{Attr}$ | ✓ | ✓ |
| $NonSK \rightarrow P_{Attr}$ | ✓ | ✗ |
| $NonSK \rightarrow NPA_{Hr}$ | ✗ | ✗ |

\therefore BCNF is inherently 3NF.

So, if R is in BCNF then R is in 3NF as well.

Ex.: When will a relation be in 3NF but not BCNF?

If the relation has no FD of form $NonSK \rightarrow NPA_{Hr}$, but there is some FD of the form $NonSK \rightarrow P_{Attr}$ then relation is in 3NF, not BCNF.

Ex. R(A, B, C, D, E, F) $FD = \{A \rightarrow BC, C \rightarrow F, D \rightarrow E\}$

AP not on RHS

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

$$\therefore CK = \{AD\}, PA_{Hr} = \{A, D\}, NPA_{Hr} = \{B, C, E, F\}$$

Not in BCNF. $\{A \rightarrow BC, A \text{ isn't SK}\}$

Not in 3NF. $AD \rightarrow A \rightarrow B$ $B \not\in AD, A \not\in B$.

Not in 2NF. $AD \rightarrow B$, but $A \rightarrow B$.

$\therefore R$ is in 1NF.

Ex. $R(A, B, C, D, E)$ $FDs = \{A \rightarrow C, B \rightarrow D, AB \rightarrow E\}$

R is in which NF?

A, B not on RHS.

$$AB^+ = A, B, C, D, E$$

$$\therefore CKs = \{AB\} \quad PAttrs = \{A, B\} \quad NPAttrs = \{C, D, E\}$$

Not in BCNF. $\{A \rightarrow C, A$ not SK}

Not in 3NF. $\{AB \rightarrow A \rightarrow C\}$ \leftarrow $C \not\in AB \& C \not\in A$

Not in 2NF. $\{AB \rightarrow C, \text{but } A \not\rightarrow C\}$

$\therefore R$ is in 1NF.

Ex. Which of the following is false?

(A) A relation w/ 2 attrs is in BCNF.

(B) A rel. w/ every key having single attr is in 2NF.

(C) A prime attr. can be transitively dependent on a key in 3NF.

(D) A prime attr. can be transitively dependent on a key in BCNF.

A. BCNF says only $SK \rightarrow X$ type of dependencies allowed.

For rel. w/ 2 attrs (say A, B) possible non-trivial FDs are:

$\{A \rightarrow B, B \rightarrow A\}$. For any subset of this:

(1) $\{A \rightarrow B\}$, $CK = \{A\}$, $PA = \{A\}$, $NPA = \{B\}$ $\therefore A \rightarrow B$ is SK.

(2) $\{B \rightarrow A\}$, B is SK.

(3) $\{A \rightarrow B, B \rightarrow A\}$ $\therefore A, B$ are SK.

\therefore Always in BCNF.

B. If every key is w/ single attr, $CR \rightarrow NPAttrs$
~~mean~~ can never be partial. \therefore Always in 2NF.

C. 3NF disallows FD of the form $\text{NonSK} \rightarrow \text{NPA}$. So, $\text{NonSK} \rightarrow \text{PA}$ & $\text{SK} \rightarrow \text{PA}$ are allowed. So PAttr. transitivity dependent on some CK is allowed.

D. Suppose it is allowed. $\text{CK} \xrightarrow{\text{Trans.}} \text{PAttr.}$ So there must be some α such that,

$$\text{CK} \xrightarrow{\leftarrow} \alpha \xrightarrow{\leftarrow} A \quad \text{A} \not\subseteq \text{CK} \quad \text{A} \not\subseteq \alpha$$

But this is BCNF. So, α must be SKey for $\alpha \rightarrow A$ to be valid FD, & if α is SK, then $\alpha \rightarrow \text{CK}$.
 \therefore contradiction.

$\therefore D$ is false.

* Natural Join

Ex. Emp

| SSN | Name | DoB | Pay-Grade | | Paygrade | Salary |
|-----|------|-----|-----------|---|----------|--------|
| | | | 1 | 2 | | |
| 101 | ABC | 213 | 2 | 1 | 2 | 70 |
| 253 | DEF | 312 | 3 | | 3 | 60 |
| 192 | GHI | 916 | 2 | 1 | 4 | 80 |
| 211 | JKL | 619 | 4 | | | |

So say we want to find DoB of (SSN=192). This only requires one table, Emp. DoB is 916.

What is the salary of SSN=253. This requires both tables.

First we check the Paygrade in Emp for SSN=253. Then from Paygrade table, we find salary for the found Paygrade. This is the natural way of joining two tables.

Natural join of two tables: join two tables with same attributes.

(1) Find the common attributes (attr. with same name).

(a) Equate the values. ~~Match~~ Match same values for

(3) In combined relation write common attributes

once ~~ABA & A~~ \rightarrow \leftarrow \rightarrow

Ex. Previous example, joined.

| SSN | Name | DoB | Payscale | Salary |
|-----|------|-----|----------|--------|
| 101 | ABC | 213 | 2 | 70 |
| 253 | DEF | 312 | 3 | 80 |
| 192 | GHI | 916 | 2 | 70 |
| 211 | JKL | 619 | 4 | 80 |

Symbol for Natural Join: \bowtie

Only the tuples for which common attributes match are allowed in the joined table.

| A | B | C | | B | C | D | | A | B | C | D |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 3 | 3 | 2 | X | 5 | 1 | 6 | ⇒ | 6 | 4 | 3 | 9 |
| 6 | 4 | 3 | | 1 | 5 | 8 | | 3 | 5 | 1 | 6 |

- If b/w two tables if there are no common attributes, then combine every row of first table with every row of second.

* Redundancy & Anomalies

- We can notice, that imposing functional dependency constraints results in data redundancy.

Ex. * District \rightarrow State

| | SSN | Distr. | State |
|---|-----|--------|-------|
| 1 | | Jaipur | Raj |
| 2 | | Jaipur | Raj |
| 3 | | Kota | Raj |

Redundancy in state attr.
no redundancy, since district is different.

- Although not explicitly stated yet but FD constraints aren't limited to a single relation. Initially there is no tables at requirement stage. There are attrs. that should be there in the DB & the FDs that ~~are~~ they should follow. Based on it & other requirements relations are designed.
So FDs ~~are~~ constraints must hold in the DB not within a single table.

- Update Anomaly: If there is redundancy, say due to $X \rightarrow Y$, and you want to update Y corresponding to $X = x_1$, then you'll have to update at multiple places (x_1, y_1) to (x_1, y_2) . This could lead to consistency, if the data isn't updated in all corresponding places.

- **Insertion Anomaly:** suppose the following

Schema:

$S(Eid, Name, Salary, Dept)$

Now you want to add a new department. This can't be done without adding an employee.

- **Deletion Anomaly:** say the last employee of Dept = Music resigns you delete the data, then with the employee, music department is also gone.

- These anomalies are happening because we are trying to put all attrs. in a single table.

The anomalies can be solved by the process of decomposition.

- Some qualities of good database design =

Making sure that semantics of the attributes is clear in schema.

- Reducing the redundant info. in table.
- Reducing the null values in tuples.
- Disallowing the possibility of generating spurious tuples.

* Decomposition

- For a relational schema R , R_1, R_2, \dots, R_n is called as decomposition of R , iff

$$\bigcup_{i=1}^n R_i = R$$

i.e. attributes must be preserved, no less no more.

Ex. Which of the following is/are decomposition of $R(A, B, C, D)$:

(a) $\{AB, CD\}$

(b) $\{AB, BCD\}$

(c) $\{A, B, C, CD\}$

(d) $\{A, BC\}$

Binary Decomposition: Decomposition into exactly two parts.

Lossless Decomposition: A decomposition which will never give wrong information.

Lossy Decomposition: A decomposition which may/may not give wrong information. (i.e. for some DB instance it will give wrong info).

| Ex. R : | a | b | c | R_1 : | a | b | R_2 : | b | c |
|-----------|---|---|---|---------|---|---|---------|---|---|
| | 1 | 2 | 3 | → | 1 | 2 | | 2 | 3 |
| | 2 | 2 | 3 | | 2 | 2 | | 1 | 2 |
| | 3 | 2 | 2 | | 3 | 2 | | 0 | 2 |

Suppose you want a for c=2.

Original: 3

Decomp: 1, 2, 3

Wrong Info.

Original table has three tuples.
If we natural join $R_1 \bowtie R_2$, then we'll get $3 \times 2 = 6$ tuples.

The extra tuples are called Spurious tuples.

So, basically in Lossy Decomposition we don't lose tuples, rather we get spurious tuples.

Hence this is loss of information, not tuples.

Ex. Same as before.

| $R_1:$ | a | b | $R_2:$ | a | c |
|--------|---|---|--------|---|---|
| 1 | 2 | | 1 | 3 | A |
| 2 | 2 | | 2 | 3 | B |
| 3 | 2 | | 3 | 2 | C |

So if we do $R_1 \bowtie R_2$, we get exactly R refinement form.

∴ For a lossy decomposition R_1, R_2, \dots, R_n of relation R
is such that:

$$R \subseteq \bigcup_{i=1}^n R_i$$

For lossless decomposition R_1, R_2, \dots, R_n of relation R

$$\bigcup_{i=1}^n R_i = R$$

In a binary decomposition $\{R_1, R_2\}$ we can check
Lossy vs Lossless :

A ~~lossy~~ binary decomposition is lossless if the common attr. $R_1 \cap R_2$ is superkey of at least one of the schema, p.e.

$$[(R_1 \cap R_2) \rightarrow R_1] \vee [(R_1 \cap R_2) \rightarrow R_2]$$

→ lossless binary decomp.

else the binary decomposition is ~~lossless~~ lossy.

Lossless decomposition is a.k.a Lossless Join Decomposition.

a.k.a Non-Additive Join Decomposition.

Ex. $R(A, B, C, D, E)$ $FDs = \{A \rightarrow BC, CD \rightarrow E, B \rightarrow D, E \rightarrow A\}$

1. $R_1 = \{A, B, C\}$ $R_2 = \{A, D, E\}$ lossless?

$$R_1 \cap R_2 = A$$

Checking $A^+ = ABCDE$

A can determine B, C (through A)

A is SK in R_1 (also in R_2)

This is enough to say R_1, R_2 is lossless decomp.

2. $R_1 = \{A, B, C\}$ $R_2 = \{C, D, E\}$

$$R_1 \cap R_2 = C$$

$$C^+ = C$$

Not SK in either.

∴ Lossy.

Ex. $R(A, B, C, D, E) FDs = A \rightarrow BC, CD \rightarrow E, B \rightarrow D, E \rightarrow A$.

$R_1(A, B, C)$

$R_2(A, D, E)$

$$R_1 \cap R_2 = D$$

$$D^+ = D$$

∴ ~~SK~~ in neither R_1 nor R_2 .

So Lossy.

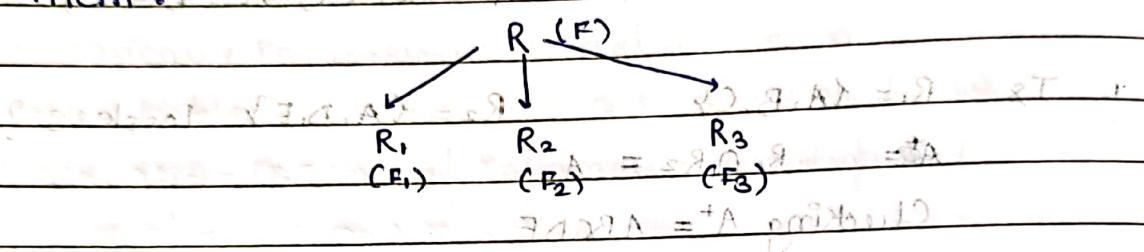
$$(A \rightarrow B, B \rightarrow A) \Rightarrow A$$

$$(A \rightarrow D, D \rightarrow A) \Rightarrow A$$

Lossless decomposition can't be done.

Dependency Preserving Decomposition:

Initially the database has some attributes & FDs on them. Say set of FD on database is F . Then the database/universal relation is broken down into multiple relations, each of which have set of FDs which only hold on them.



A decomposition of relation R with FD set F into multiple relations, $R_1(F_1); R_2(F_2); \dots; R_n(F_n)$ is called Dependency Preserving iff,

$$(F_1 \cup F_2 \cup \dots \cup F_n) \equiv F$$

i.e.

$$\bigcup_{i=1}^n F_i \text{ covers } F \text{ & } F \text{ covers } \bigcup_{i=1}^n F_i$$

Ex. $R(A, B, C, D, E)$ $FD_R = \{A \rightarrow BC, CD \rightarrow E, A \rightarrow D, E \rightarrow A\}$

$$R_1(A, B, C) \quad R_2(A, D, E)$$

Residual Velocity on Bus

$$A^+ = ABC + DE$$

$$B^+ = BD$$

$$C^+ = C$$

$$D^+ = D$$

$$E^+ = EABCDP$$

$$\text{In } R_1 : A^+ = ABC / B^+ = B / C^+ = C / D^+ = D / E^+ = EAD$$

$$F_1 = \{A \rightarrow BC, BC \rightarrow A\}$$

$$F_2 = \{A \rightarrow DE, E \rightarrow AD\}$$

(we are only writing Non-Trivial FDs)

Checking Preservence:

$(F_1 \cup F_2)^+$ doesn't cover F

$$(CD)^+ = CD$$

Hence $CD \rightarrow E$ isn't covered

$$\{A-B, A-C\}^+ = \{A\} \quad \{E, A, C\}^+ = \{E\}$$

∴ Not Dependency preserving.

$$\{A, B, C\}^+ = \{A\} \quad \{E, A, C\}^+ = \{E\}$$

Even if a decomposition is lossless but not dependency preserving, it doesn't really loose dependency.

Say R is decomposed into R_1 & R_2 & it is not dependency preserving.

Still when $R_1 \bowtie R_2$, all the dependencies are restored since the decom. is lossless.

It's called non-dependency preserving since without joining the dependency can't be enforced.

That means if we do a join only the valid rows (which follow all FDs) would be there but without the join the individual tables ~~can't~~ can have data that doesn't follow all the dependencies.

Ex. $R(\text{Sid}, \text{State}, \text{District})$

$$FD_R = \{ \text{Dist} \rightarrow \text{State}, \text{Sid} \rightarrow \text{Dist}, \text{Sid} \rightarrow \text{State} \}$$

$R_1(\text{Sid}, \text{State})$

$$F_1 = \{\text{Sid} \rightarrow \text{State}\}$$

• This is lossless but not dependency preserving.

In R the tuple $\langle 213, \text{Haryana}, \text{Jaipur} \rangle$ would not be allowed.

In $R_1 \langle 213, \text{Haryana} \rangle$ & in $R_2 \langle 213, \text{Jaipur} \rangle$ is allowed.

In $R_1 \bowtie R_2$, the tuple won't appear. So although an lossless join dependencies aren't lost, individually the lost FDs can't be enforced on R_1 & R_2 .

- In any decomposition lossless property is necessary, dependency preserving is desirable.

Ex. $R(A, B, C, D, E)$ $FD_3 = \{A \rightarrow BC, CD \rightarrow E, B \rightarrow D,$

$E \rightarrow A\}$

① $R_1(A, B, E) \quad R_2(B, C, E) \quad R_3(C, D, E)$

$R_1 \cap R_2 = \{BE\}$

$BE^+ = B, E, D, A, C$

$\therefore R_1, R_2$ is lossless decompos.

$F_1 = \{A \rightarrow BE, E \rightarrow A\} \quad F_2 = \{E \rightarrow BC, BC \rightarrow E\} \quad F_3 = \{B \rightarrow D, CD \rightarrow B\}$

$(F_1 \cup F_2 \cup F_3)$ covers $F_{(3)}$ ad from given condition

F obviously covers $F_1 \cup F_2 \cup F_3$ ab. no. 3 remove last

last \therefore This is dependency preserving

Ex. $R(A, B, C, D)$ $FD_3 = \{A \rightarrow B, B \rightarrow C, C \rightarrow D, D \rightarrow A\}$

$R_1(A, B) \quad R_2(B, C) \quad R_3(C, D)$

$F_1 = \{A \rightarrow B, B \rightarrow C\} \quad F_2 = \{B \rightarrow C, C \rightarrow D\} \quad F_3 = \{C \rightarrow D, D \rightarrow A\}$

$(F_1 \cup F_2 \cup F_3)^*$ covers F

\therefore Dependency preserving.

- How much should the database be decomposed depends on the Normal form. As we got from $1NF \rightarrow 2NF \rightarrow 3NF \rightarrow BCNF$, redundancy due to FDs keep decreasing with 0% in BCNF (0% due to FDs, there can be other type of redundancy).
 $R_1 \times R_2 = \{ \text{mild form} \}$ {such as MVD}

~~Normalization~~: Decomposing a relation R such that each decomposed relation is in a higher normal form than R , and the dependency must be lossless.

~~Non-Binary Decomposition~~: NB Decomposition of a relational schema R into R_1, R_2, \dots, R_m , $m \geq 3$, such that

$$R_1 \cup R_2 \cup \dots \cup R_m = R \text{ with no anomalies}$$

e.g. $R = R_{\text{dept}} \cup R_{\text{emp}}$

We can check if a Non-Binary decomp. is Dependency Preserving or not by using the def'n of Dependency Preserving decomp.

To check if a binary decomp of $R: R_1, R_2$, is lossless we check

$$R_1 \cap R_2 \rightarrow R_1 \text{ is } (A, A) \text{ and } (B, A)$$

$$\text{OR } (A, B) \text{ and } (B, A)$$

$$R_1 \cap R_2 \rightarrow R_2$$

This doesn't work for NonBinary decomp.

We use "Chase Algorithm" to check lossless join property of NBdecomp.

Chase Algorithm: Used to check if a decomposition (binary or non-binary) is lossless or not.

Inputs: Relational Schema $R(A_1, A_2, \dots, A_n)$, a set of FDs F & a decomposition $\rho = (R_1, R_2, \dots, R_k)$

- Construct a matrix of size $k \times n$ where each row represents one subschema & each column represents one attribute.
- Mark $M[i][j]$ if A_j is in subschema R_i with a_{ij} , else with b_{ij} .
- For each FD $X \rightarrow Y$ in FD set $(X, Y \subseteq R)$ if there are two or more rows such that for all attrs in X , the rows have same value, then in those rows we equate the values for all attrs. In Y , if one of the equal val. for X rows, if $A_j (A_j \in Y)$ has a_{ij} entry then make all b_{ij} in the other rows entry as a_{ij} . Else if one of them is b_{ij} make all of them b_{ij} .
- Repeat on the FD set until no changes to M .
- If now any row looks like a_1, a_2, \dots, a_n then decomp is lossless, else lossy.

Ex. $R(A, B, C, D)$ FDs = $\{A \rightarrow B, B \rightarrow C, C \rightarrow D\}$.

$R_1(A, B, C)$ $R_2(C, D)$.

| | A | B | C | D |
|----------------|-----------------|-----------------|----------------|---------------------------------|
| R ₁ | a ₁ | a ₂ | a ₃ | b ₁ , a ₄ |
| R ₂ | b ₂₁ | b ₂₂ | a ₃ | a ₄ |

This is final table.

We can see R_1 is a₁, a₂, a₃, a₄.

∴ Lossless decomp.

Ex. $R(A, B, C, D)$ $FD_S = \{A \rightarrow B, B \rightarrow C, C \rightarrow D, D \rightarrow A\}$

$R_1(A, B)$

$R_2(B, C)$

$R_3(C, D)$

| | A | B | C | D |
|-------|-------------------------------|-------------------------------|-------------------------------|-------------------------------|
| R_1 | a ₁ | a ₂ | b ₁ a ₃ | b ₁ a ₄ |
| R_2 | b ₂ a ₁ | a ₂ | a ₃ | b ₂ a ₄ |
| R_3 | b ₃ a ₁ | b ₃ a ₂ | a ₃ | a ₄ |

\therefore Lossless ($\{a_1, a_2, a_3, a_4\} = \{a_1, a_2, a_3, a_4\}$).

Ex. $R(A, B, C, D, E)$! $FD_S = \{A \rightarrow C, B \rightarrow C, C \rightarrow D, DE \rightarrow C, CE \rightarrow A\}$

$R_1(A, D)$ $R_2(A, B)$ $R_3(B, E)$ $R_4(C, D, E)$ $R_5(A, E)$

| | A | B | C | D | E |
|-------|----------------|----------------|-------------------------------|----------------|----------------|
| R_1 | a ₁ | | b ₁ a ₃ | a ₄ | |
| R_2 | a ₁ | a ₂ | b ₁ a ₃ | a ₄ | |
| R_3 | a ₁ | a ₂ | b ₁ a ₃ | a ₄ | a ₅ |
| R_4 | a ₁ | | a ₃ | a ₄ | a ₅ |
| R_5 | a ₁ | | b ₁ a ₃ | a ₄ | a ₅ |

(assume all \Rightarrow initially empty cells have b_{ij})

We can stop at this point, since R_3 has a₁a₂a₃a₄a₅ (not because no further changes).

\therefore Lossless decomposition.

Ex. $R(A, B, C, D)$ $FD_S = \{A \rightarrow B, B \rightarrow C, C \rightarrow D, D \rightarrow B\}$ $R_1(A, B)$ $R_2(B, C)$ $R_3(C, D)$

| | A | B | C | D |
|-------|----------------|----------------|----------------|----------------|
| R_1 | a ₁ | a ₂ | a ₃ | a ₄ |
| R_2 | | a ₂ | a ₃ | a ₄ |
| R_3 | | a ₂ | a ₃ | a ₄ |

\therefore Lossless

Ex. $R(A, B, C, D, E) \rightarrow R_1(A, B, C) \rightarrow R_2(B, C, D) \rightarrow R_3(A, D, E)$
 $R_1(A, B, C) \rightarrow R_2(B, C, D) \rightarrow R_3(A, D, E)$

Check if the above decomp. is lossless & dependency preserving with following FD set:

$$(A) B \rightarrow E \& CE \rightarrow A$$

Checking DP:

$$F_1 = \{BC \rightarrow A\} \quad F_2 = \{\} \quad F_3 = \{CE \rightarrow A\}$$

∴ So, not preserved.

Checking lossless:

| | A | B | C | D | E |
|-------|----------------|----------------|----------------|----------------|---|
| R_1 | a ₁ | a ₂ | a ₃ | a ₄ | b ₁ , b ₂ , b ₃ , b ₄ |
| R_2 | a ₁ | a ₂ | a ₃ | a ₄ | b ₅ , b ₆ |
| R_3 | a ₁ | | a ₃ | | a ₅ |

So, ~~lossy~~ lossy.

$$(B) AC \rightarrow E \& BC \rightarrow D$$

$$F_1 = \{\} \quad F_2 = \{BC \rightarrow D\} \quad F_3 = \{AC \rightarrow E\}$$

∴ Dependency preserved.

| | A | B | C | D | E |
|-------|----------------|----------------|----------------|----------------|---------------------|
| R_1 | a ₁ | a ₂ | a ₃ | a ₄ | a ₅ |
| R_2 | a ₁ | a ₂ | a ₃ | a ₄ | full to date and SW |
| R_3 | a ₁ | | a ₃ | | a ₅ |

∴ Lossless

$$(C) A \rightarrow D, D \rightarrow E, B \rightarrow D$$

$$F_1 = \{\} \quad F_2 = \{B \rightarrow D\} \quad F_3 = \{A \rightarrow E\}$$

Not DP.

| | A | B | C | D | E |
|----------------|----------------|----------------|----------------|-----------------|----------------|
| R ₁ | a ₁ | a ₂ | a ₃ | b ₁₄ | a ₅ |
| R ₂ | | a ₂ | a ₃ | a ₄ | a ₅ |
| R ₃ | a ₁ | | a ₃ | b ₁₄ | a ₅ |

∴ Lossless.

(D) $A \rightarrow D$, $CD \rightarrow E$ & $E \rightarrow D$.

$$F_1 = \{ A \} \quad F_2 = \{ CD \} \quad F_3 = \{ AC \rightarrow E \}$$

Not DP.

| | A | B | C | D | E |
|----------------|----------------|----------------|----------------|-----------------|----------------|
| R ₁ | a ₁ | a ₂ | a ₃ | b ₁₄ | a ₅ |
| R ₂ | | a ₂ | a ₃ | a ₄ | a ₅ |
| R ₃ | a ₁ | | a ₃ | b ₁₄ | a ₅ |

∴ Lossy.

Ex. R(L,M,N,O,P) FD_R = {M → O, NO → P, P → L, L → MN}

$$R_1(L, M, N, P) \quad R_2(M, O)$$

(1) Is the above decomp. lossless?

(2) Is the above decomp. dependency preserving?

(3) What is the highest normal form satisfied by above decomp?

$$(1) R_1 \cap R_2 = M \quad M^+ = MO \quad M \rightarrow R_2$$

∴ lossless.

$$(2) F_1 = \{ L \rightarrow MN \text{ OR } P \rightarrow L \} \quad F_2 = \{ M \rightarrow O \}$$

NO → P isn't covered by $F_1 \cup F_2$.

∴ Not dependency preserving.

(3) For R₁:

$$CK_R = \{ L, P \}$$

∴ R₁ is BCNF

(since all LHS are key) ←

For R₂:

$$CK_R = \{ M \}$$

∴ R₂ is BCNF

∴ Highest NF

satisfied by decomp
is BCNF.

∴ R₂ has 2 attrs.

* Minimal Cover of FDs

(a.k.a Canonical Cover, aka irreducible FD set)

most appropriate name

~~Let's say we define stronger FD (non-standard term) as:~~

Among 2 FDs, F_1 & F_2 (F_1 is stronger than F_2) iff

$$F_1 \rightarrow F_2 \text{ (or } F_1 \sqsupseteq F_2\text{)}$$

$$\text{Ex. } A \rightarrow C \\ (f_1)$$

$$A \rightarrow BC \\ (f_2)$$

$$\text{Here, } f_2 \rightarrow f_1$$

$\therefore f_2$ is stronger

$$\text{Ex. } AB \rightarrow C \\ (f_1)$$

$$A^+ = ABC \\ (\text{using } f_2)$$

$$\therefore f_2 \rightarrow f_1$$

So, f_2 is stronger

$$\text{Ex. } A \rightarrow B \\ (f_1)$$

$$A \rightarrow C \\ (f_2)$$

$$A^+ = AC \\ (\text{using } f_2)$$

$$A^+ = AB \\ (\text{using } f_1)$$

so neither implies the other

Can't compare

For an FD set F , its closure F^+ is defined as

$$F^+ = \{g \mid g \text{ is a FD} \wedge F \sqsupseteq g\}$$

closure

So, F^+ includes all the FDs & g that can be inferred from F.

- We can go the other way around for a FD set F, and get a set of only essential FDs, i.e., the FD set which is minimal / irreducible / canonical & covers F.

Extraneous Attributes on LHS of a FD (Redundant)

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

here, we don't have need for B, A can imply C.

So, ~~$\{A \rightarrow C\}$~~ $\rightarrow F$

A FD $\alpha A \rightarrow \beta B$, ($\alpha \subseteq C_R, \beta \in R$) has A as extraneous (redundant) variable iff $\alpha \rightarrow \beta$, or $B \subseteq \alpha^+$.

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

In $AB \rightarrow C$ ~~$\{AB\}$~~

• Removing B, $A^+ = AB$

• Removing A, $B^+ = B$ ~~to remove redundant column~~

So, both A, B are essential.

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

For ~~$AB \rightarrow C$~~ $A^+ = ABC$ $B^+ = B$ ~~so, A determines C~~

i.e., A can determine C. ~~A solution to part~~

So, B is redundant.

$\therefore \{A \rightarrow C, A \rightarrow BC\}$

Ex. $F = \{ AC \rightarrow G, ACD \rightarrow B, BC \rightarrow D, CG \rightarrow B, CG \rightarrow D, CE \rightarrow A, CE \rightarrow G, D \rightarrow E, D \rightarrow G \}$

For $AC \rightarrow G$: $A^+ = A$, $C^+ = C$. \therefore No extra attr. in LHS. \therefore No extraneous.

For $ACD \rightarrow B$: $A^+ = ACGB$ (transitive closure) \therefore D is extraneous.

$$A^+ = A \quad C^+ = C$$

\therefore Both not extraneous. \therefore $AC \rightarrow B$.

For $BC \rightarrow D$: $B^+ = B$, $C^+ = C$.

\therefore No extra.

For $CG \rightarrow B$: $C^+ = C$, $G^+ = G$. \therefore No extra.

For $CG \rightarrow D$: $C^+ = C$, $G^+ = G$. \therefore No extra.

For $CE \rightarrow A$: $C^+ = C$, $E^+ = E$.

For $CE \rightarrow G$: $C^+ = C$, $E^+ = E$. \therefore No extra.

For $CG \rightarrow D$: $C^+ = C$, $G^+ = G$. \therefore No extra.

$\therefore FD = \{ AC \rightarrow G, AC \rightarrow B, BC \rightarrow D, CG \rightarrow B, CG \rightarrow D, CE \rightarrow A,$

$CE \rightarrow G, D \rightarrow E, D \rightarrow G \}$

Removing Redundant FD

First make sure that on RHS we only have single attributes.

In a FD set F , a FA $X \rightarrow A$ ($X \subseteq R$, $A \in R$) is redundant iff $F - \{ X \rightarrow A \}$ still infers $X \rightarrow A$ (or still infers $A \in X^+$).

Basically, if using all FDs except $X \rightarrow A$ if we can show that X^+ contains A , then $X \rightarrow A$ is redundant.

Transitive closure of A is A^+ .

$F \rightarrow A \iff A \subseteq X^+$

Ex. $F = \{AC \rightarrow G, CD \rightarrow B, BC \rightarrow D, CG \rightarrow B, CG \rightarrow D, CE \rightarrow A, CE \rightarrow G, D \rightarrow E, D \rightarrow G\}$, all attributes are in LHS.

For $AC \rightarrow G$: $AC^+ = AC$

For $CD \rightarrow B$: $CD^+ = CDEG \setminus A$ $\therefore CD \rightarrow B$ is redundant.

~~FOR $BC \rightarrow D$. PROBLEMS~~ towards redundant attribute removal

Now FD is $F - \{CD \rightarrow B\}$. Further checking will be done with it.

For $BC \rightarrow D$: $BC^+ = BC$.

For $CG \rightarrow B$: $CG^+ = CG \setminus DEA$

For $CG \rightarrow D$: $CG^+ = CG \setminus BDEA \therefore CG \rightarrow D$ is redundant.

Now FD is $F - \{CD \rightarrow B, CG \rightarrow D\}$

For $CE \rightarrow A$: $CE^+ = CEGBD$

For $CE \rightarrow G$: $CE^+ = CEAGBD \therefore CE \rightarrow G$ is redundant.

FD is $F - \{CD \rightarrow B, CG \rightarrow D, CE \rightarrow G\}$

For $D \rightarrow E$: $D^+ = DG$

For $D \rightarrow G$: $D^+ = DE$

$\therefore FD = \{AC \rightarrow G, BC \rightarrow D, CG \rightarrow B, CE \rightarrow A, D \rightarrow E, D \rightarrow G\}$

Note: Minimal Covers for a FD is not unique. It depends on order of removing extraneous vars. on LHS & order of removing redundant FDs.

Steps to find minimal cover of FD:

- Break FDs if needed, such that every FD on RHS has only one attr.
- For every FD, check if one or more variable is redundant & remove the redundant variables on LHS.
- For all remaining FDs, check if they are redundant, and remove them.

- Minimal Cover doesn't mean minimum no. of FDs, rather it means that no FD can further be removed from the cover.

- Order of removing redundant RHS vars & removing redundant FDs doesn't matter.

Ex. R(A, B, C, D, E, F, G)

$$FD_1 = \{ AC \rightarrow G, D \rightarrow E, D \rightarrow G, BC \rightarrow D, CG \rightarrow B, CG \rightarrow D, A \leftarrow D, A \leftarrow G, ACD \rightarrow B, CE \rightarrow A, CE \rightarrow G \}$$

$$ACD \rightarrow B, CE \rightarrow A, CE \rightarrow G$$

S1: Removing extraneous on LHS

$$AC \rightarrow G : A^+ = A \cup D, C^+ = C \cup D \Rightarrow A \cap C = \emptyset \Rightarrow A \leftarrow D$$

$$D \rightarrow E$$

$$D \rightarrow G : D^+ = D \cup A \cup C \Rightarrow D \cap A = D \cap C = \emptyset \Rightarrow D \leftarrow A, D \leftarrow C$$

$$BC \rightarrow D : B^+ = B \cup C \Rightarrow B \cap C = \emptyset$$

$$CG \rightarrow B : C^+ = C \cup G, G^+ = G \cup D \Rightarrow C \cap G = \emptyset \Rightarrow C \leftarrow D$$

$$CG \rightarrow D : C^+ = C \cup G, G^+ = G \cup D \Rightarrow C \cap G = \emptyset \Rightarrow C \leftarrow D$$

$$ACD \rightarrow B : AC^+ = A \cup C \cup D \cup B \cup F \Rightarrow D \text{ is redundant} \Rightarrow A \leftarrow D$$

$$CE \rightarrow A : C^+ = C \cup E \Rightarrow C \cap E = \emptyset \Rightarrow C \text{ is primary}$$

$$CE \rightarrow G : C^+ = C \cup E, G^+ = G \cup D \Rightarrow C \cap G = \emptyset \Rightarrow C \leftarrow D$$

$$\therefore F = \{ AC \rightarrow G, D \rightarrow E, D \rightarrow G, BC \rightarrow D, CG \rightarrow B, CG \rightarrow D, AC \rightarrow B, CE \rightarrow A, CE \rightarrow G \}$$

S2: Removing red. FDs.

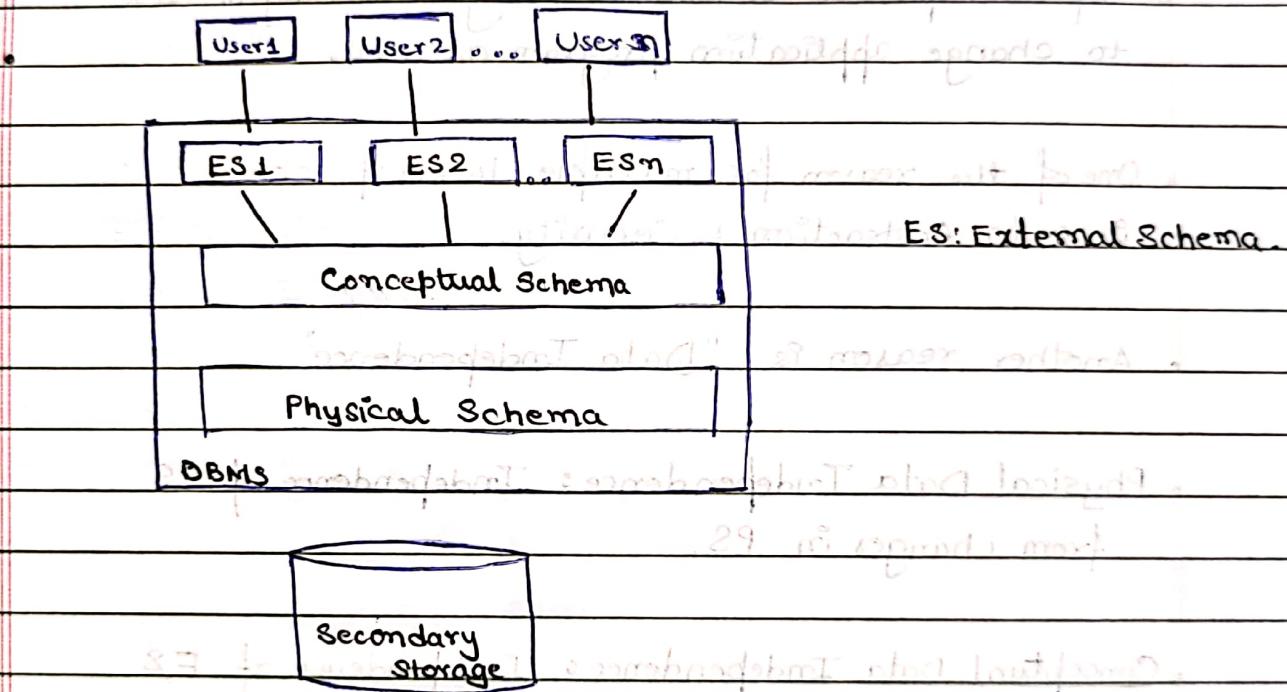
$$\text{Red. } FD_2 = \{ CG \rightarrow B, CE \rightarrow G \}$$

$$\therefore \text{Minimal Cover: } \{ AC \rightarrow G, D \rightarrow E, D \rightarrow G, BC \rightarrow D, CG \rightarrow B, AC \rightarrow B, CE \rightarrow A \}$$

* Three levels of Schema.

- In a database, not all the data be visible to every users. Different users require different level of access & visibility to data.

Ex. In a bank's database, a user must not be able to access other customer's transactions.



Physical Schema: Tells how the data is actually stored on the disk, sorted or unsorted, indices, datastructure, etc. (Internal schema)

Conceptual Schema: Tells what ~~data~~ is the meaning (structure, interpretation, constraints, etc) of the data. (Logical schema)

External Schema: Tells about the use of data by various DB users. (~~different users~~)

Data is actually stored on secondary storage.

- External schema isn't stored on disk. It is derived as per requirement from conceptual schema.

Ex. Application program can require data consisting of attr. from multiple relations, or based on some filters, or some specific ordering.

- If conceptual schema is changed, we might have to change application program as well.

- One of the reason for multiple levels of schema is Data Abstraction & Security.

- Another reason is "Data Independence".

- Physical Data Independence: Independence of CS from changes in PS.

- Conceptual Data Independence: Independence of ES from changes in CS.

Although possible, it is rare that change to PS results in change at CS level. This is ensured by DBMS.

So DBMS provides us the facility to change physical changes without requiring changes at CS level & Application level.

- DBMS also takes up the responsibility to ensure that changes made at CS level should rarely result in changes at ES level. This is more difficult to achieve.
[Conceptual / Logical Data Indep.]

Normalization

(3NF & BCNF Decomposition)

* Redundancies

- Trivial FDs never create redundancies.
- Non-Trivial FDs may/maynot create redundancies.

Ex. Student (sid, name, dist, state) $\text{dist} \rightarrow \text{state}$.

Say there are 500 students from Jaipur, then all of them will have state Rajasthan. So the repetition of state value is redundancy.

Suppose we decompose the relation as:

S (sid, name, dist) R (dist, state).

Now $\langle \text{dist}, \text{state} \rangle$ tuple will be stored once & state won't be redundantly specified.

So, decomposition is a soln to redundancy.

Why does $(\text{dist} \rightarrow \text{state})$ causes redundancy?

① Repetition of dist causes repetition of state.

② Since dist is not SK in rel. so dist can repeat.

Ex. R(A,B,C) FDs = { $A B \rightarrow C$, $C \rightarrow B$ }

Which FD may create redundancy?

CKs = { $A B$, $A C$ }

So, AB is SK, it won't repeat.

C is not SK, it can repeat, & repetition of C causes repetition of B.

$\therefore C \rightarrow B$ may create redundancy

So, non-trivial FD with LHS as superkey can't cause redundancy.

- Non-trivial FD $\alpha \rightarrow \beta$ can cause redundancy iff α is not SK.

- In BCNF every non-trivial FDs have SK on LHS. So, BCNF is free from FD caused redundancy.
- 3NF may have FD caused redundancies.

Ex. Courses (sid, cid) Movies (csid, Movie)

One student can take many courses & one course can be taken by many students. Many-to-many relationship. There is no redundancy in this table.

Similarly, no redundancy in Movies table. It is also many-to-many relationship.

If we merge them in one student table.

Student (sid, cid, Movie)

C: sid cid and cid is not Movie; not sid from Movie

s₁ -> c₁ s₁ -> c₂ s₁ -> c₃ A

s₂ -> c₁ s₂ -> c₂ s₂ -> c₃ B

s₃ -> c₁ s₃ -> c₂ s₃ -> c₃ C

s₄ -> c₁ s₄ -> c₂ s₄ -> c₃

When merged:

S: sid cid Movie

s₁ c₁ A

s₁ c₁ B

s₁ c₁ C

s₁ c₂ A

s₁ c₂ B

s₁ c₂ C

So, we can clearly see there is redundancy. Same info

repeats multiple times. But there's no non-trivial FDs in Student relationship. So, where is this redundancy coming from? This happened because more than one many-to-many relationship, in a single table.

(Actually this creates a new type of dependency, called Multivalued Dependency (MVD))

This dependency says, every movie causes repetition of courses, & vice-versa.

Redundancy

~~FD caused redundant part of R~~

~~(removed by BCNF)~~

~~MVD caused~~

~~(removed by 4NF)~~

* BCNF Decomposition

- A relation R is in BCNF iff for every non-trivial FD of R $\alpha \rightarrow \beta$, α is PK of R or β is trivial.
- A relation R with two attributes must be in BCNF.
- A set of relations $S = \{R_1, R_2, R_3, \dots, R_n\}$ is in BCNF iff $\forall i R_i$ is in BCNF.
- If a relation R is not in BCNF, we have to decompose it into subrelations such that each subrelation is in BCNF.

Ex. $R(A, B, C, D)$ $R_1(A, B) \sqcup R_2(C, D)$.

This is BCNF decomposition. But it is lossy. $= R_1 \cap R_2 = \emptyset$.

- To check if a decomposed set of relation is in BCNF, we have to check if all relations are BCNF.

(So, if asked to decompose to BCNF, we must give a lossless decomposition. If decomposition is already given to check, then it may be BCNF but not lossless.)

Note: Question asks to find BCNF decomposition:

(1) Must be lossless.

(2) Try for dependency preservation, if possible.

Question asks if given decomposition is BCNF, then it doesn't matter if lossless or lossy.

(The same is true for 2NF & 3NF decomposition also)

BCNF Decomposition Algorithm:

Guarantees lossless BCNF Decomposition

For a relation R & FD set F, if R is in BCNF, stop.

If R is not in BCNF, then for the violating FD: $X \rightarrow \alpha$ ($X \subset R$, $\alpha \in R$ & the FD is non-trivial), split R into

$$R_1 = X^+ \quad R_2 = (X - X^+) \cup \{X\}$$

~~Ex. R(A, B, C) & FD = {A → B}~~

This is guaranteed lossless, since $R_1 \cap R_2 = X$ & $X \rightarrow R_1$.

Ex. R(A, B, C) & FD = {A → B}

$$CK_R = \{AC\}$$

So, $A \rightarrow B$ violates BCNF.

$$\therefore R_1(A, B)$$

$$R_1 = A^+ = \{A, B\}$$

$$R_2(A, C)$$

$$R_2 = (A, B, C) - \{AB\} \cup \{A\} = \{A, C\}$$

Ex. $R(A, B, C, D)$ $F = \{A \rightarrow B; A \rightarrow C\}$

$$CK_3 = \{AD\}$$

so, $A \rightarrow B$ is violation: $A^+ = \{A, B, C\} \setminus \{AD\} = \{B, C\}$

$$R_1 = \{A, B, C\} = A^+$$

$$R_2 = (R - A^+) \cup \{A\} = \{A, D\}$$

$$\therefore R_1(A, B, C)$$

$$F_1 = \{A \rightarrow B, A \rightarrow C\}$$

Both are in BCNF.

$$R_2(A, D)$$

$$F_2 = \{\}$$

Ex. $R(A, B, C, D)$

$$CK_3 = \{AD\}$$

so, $B \rightarrow C$ is violating.

$$R_1 = B^+ = \{B, C\}$$

$$R_2 = (R - B^+) \cup \{B\} = \{A, B, D\}$$

$$R_1(B, C)$$

$$F_1 = \{B \rightarrow C\}$$

This is in BCNF.

$$R_2(A, B, D)$$

$$F_2 = \{A \rightarrow B\}$$

$$CK_3 = \{AD\}$$

Not in BCNF. ($A \rightarrow B$)

$$\Delta R_{21} = A^+ = \{AB\}$$

$$R_{22} = \{AD\}$$

$$\therefore R_1(B, C)$$

$$R_{21}(A, B) = \{A\}$$

$$R_{22}(A, D) = \{D\}$$

are all in BCNF.

Note: You may get different BCNF decomposition based on the order in which we resolve violation.

Ex. $R(A, B, C, D, E)$

$$F = \{AB \rightarrow C, C \rightarrow D, D \rightarrow E\}$$

$$CK_1 = \{AB, C\}$$

So, $C \rightarrow D$ violates BCNF.

$$R_1 = C^+ = CDE$$

$$R_2 = ABC$$

$$R_1(C, D, E)$$

$$F_1 = \{C \rightarrow D, D \rightarrow E\}$$

~~R₂~~ is in BCNF. R₁ is not.

$$R_2(A, B, C) = ?$$

$$F_2 = \{AB \rightarrow C\}$$

$$R_{11} = D^+ = DE$$

$$R_{12} = CD$$

$$R_{11}(C, D)$$

$$R_{12}(D, E) = (R_2(A, B, C))^+ = ?$$

All are in BCNF.

$$(A, B, C) < ? \text{ and } (A, B) < ?$$

$$(B \leftarrow A) = ?$$

$$(A \leftarrow B) = ?$$

Ex. $R(A, B, C, D, E)$

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

$$CK_1 = \{A, E, BC, CD\}$$

R Violation: $AB \rightarrow D$

$$R_1 = B^+ = BD \quad F_1 = \{B \rightarrow D\}$$

$$R_2 = ABCE \quad F_2 = \{A \rightarrow BCE, E \rightarrow A, BC \rightarrow E\}$$

$$CK_2 = \{A, E, BC\}$$

So, both are in BCNF.

Ex. $R(A, B, C)$, $F = \{AB \rightarrow C, C \rightarrow B\}$

$$CKs = \{AB, AC\}$$

R violation: $C \rightarrow B$ is not valid w.r.t. AB .

$$R_1 = C^+ = \{C, BC\} \text{ Both in BCNF.}$$

$$R_2 = \{A, C\}$$

$$F_1 = \{C \rightarrow B\}, F_2 = \{\} \text{ since } C \rightarrow B \text{ is not valid.}$$

This is lossless, but not dependency preserving.

BCNF

A relation R with FD set F can always be decomposed into some lossless set of tables $\{R_1, R_2, \dots, R_n\}$, but it is not guaranteed that some dependency preserving[↑] decomposition exists for R .

So, for all relations[↑] lossless BCNF decomposition always exist.

For some relations, no lossless dep. preserving BCNF decomposition doesn't exist.

* 3NF Decomposition

For 3NF Decomposition, we'll use "Synthesis Algorithm". It guarantees lossless & dependency preserving decomposition.

Although the BCNF decompo. algo. also is 3NF decompo., but it doesn't guarantee

Synthesis Algo : Input: Relation R , FD set F

- Find minimal cover of F , say F_c .

- In F_c , if for two or more FDs have LHS same, merge them.

$$\alpha \rightarrow A, \alpha \rightarrow B \Rightarrow \alpha \rightarrow AB.$$

So, you'll have FDs in F_c of form $\alpha \rightarrow B$. ($\alpha, B \subseteq R$).

- For every FD f_i in F : $\alpha \rightarrow f_i$ create a new relation $R_i(\alpha, \ldots, f_i)$
- If none of the relations R_i contains a CK, then take any one CK & make a new relation for it.
- Out of all new relations, delete any such R_i , for which there is some $R_j = R_i \cap R_j$, until no other relation can be deleted.

Ex. Give 3NF decomposition for relation A

$R(A, B, C, D, E)$ + LHS: $(A, B) \rightarrow C, (B, C) \rightarrow D, A \rightarrow D$ rightmost attr. same fact

$F = \{AB \rightarrow C, C \rightarrow B, A \rightarrow D\}$ remove redundant attr.

Checking 3NF:

$CKs = \{AFC, AEB\}$ N.P. Atts = {D} partial dependency

Not even in 2NF. $AEB \rightarrow D$ is partial dependency

Applying Synthesis algo:

S1: Finding Minimal Cover:

$$\begin{array}{lll} AB \rightarrow C & AB \rightarrow C & AB \rightarrow C \\ C \rightarrow B & \xrightarrow[LHS]{\text{minim.}} C \rightarrow B & \xrightarrow[\text{Useless}]{\text{FD Removal}} C \rightarrow B \end{array}$$

$A \rightarrow D$ contract $A \rightarrow D$ info $A \rightarrow D$

S2: $R_1(A, B, C) \quad R_2(C, B) \quad R_3(A, D) \quad R_4(A, E)$

S4: $R_1(A, B, C) \quad R_2(B, C) \quad R_3(A, D) \quad R_4(A, C, E)$ CK.

S5: Remove $R_2(R_2 \subset R_1)$

$R_1(A, B, C) \quad R_2(A, D) \quad R_3(A, C, E)$

Note that Synthesis Algo. doesn't care if the original relation R is in 3NF or not.

This algo is good, since algorithmically, checking if R is in 3NF or not is more expensive than just applying Synthesis Algo.

Hence it's not called 3NF Decompo. Algo, but 3NF Synthesis Algo.

So, for ALL relation R, it is always possible to get a:

- (1) Lossless BCNF decomposition.
- (2) Lossless Dependency Preserving 3NF decomposition.

Ex. $R(A, B, C, D, E, F)$, $F = \{A \rightarrow BCD, BC \rightarrow DE, B \rightarrow D, D \rightarrow A\}$

(We won't check for 3NF)

S1: Minimal Cover.

$$\begin{array}{l} A \rightarrow B \quad B \rightarrow C \rightarrow D \\ A \rightarrow C \quad BC \rightarrow E \\ A \rightarrow D \quad B \rightarrow D \quad A \rightarrow D \\ D \rightarrow A \end{array} \quad \begin{array}{l} A \rightarrow B \quad B \rightarrow D \\ A \rightarrow C \quad B \rightarrow E \\ A \rightarrow D \quad B \rightarrow E \\ D \rightarrow A \end{array} \quad \begin{array}{l} A \rightarrow B \quad B \rightarrow D \\ A \rightarrow C \quad B \rightarrow E \\ A \rightarrow D \quad B \rightarrow E \\ D \rightarrow A \end{array}$$

S2: $A \rightarrow BC, B \rightarrow DE, D \rightarrow A$ S3: $R_1(A, B, C, D), R_2(B, D, E), R_3(A, D)$

S4: $CKS = \{FA, FB, FD\}$

$\Rightarrow R_1(A, B, C, D) \cup R_2(B, D, E) \cup R_3(A, D) \cup R_4(A, F)$.

S5:

$R_1(A, B, C, D) \quad R_2(B, D, E) \quad R_3(A, D) \quad R_4(A, F)$.

Ex. $R(A, B, C, D)$, $F = \{AB \rightarrow C, BC \rightarrow D, CD \rightarrow A, AD \rightarrow B\}$

S1: $AB \rightarrow C$

$AB \rightarrow C \quad AB \rightarrow C$

$BC \rightarrow D$

$BC \rightarrow D \quad BC \rightarrow D$

$CD \rightarrow A$

$CD \rightarrow A \quad CD \rightarrow A$

$AD \rightarrow B$

$AD \rightarrow B \quad AD \rightarrow B$

$AD \rightarrow B \quad AD \rightarrow B$

$CKS = \{AB, BC, AD, CD\}$

$S_3: R_1(A, B, C) \quad R_2(B, C, D) \quad R_3(A, C, D) \quad R_4(A, B, D)$

$S_{4,5}: R_1(A, BC) \quad R_2(A, C, D)$

$R_3(A, B, D) \quad R_4(A, B, D)$

* 2NF Decomposition

- In a relation R , 2NF violation is caused by presence of FDs of form $C \rightarrow NPA$

Proper subset of CK $\rightarrow NPA$

as this type of FD makes some $C \rightarrow NPA$ a partial dependency.

- For a relation R with FD set F , if any $FD \alpha \rightarrow A \subset A$ in F^+ violates 2NF, then decompose: $\{ \alpha \in R, A \in R \} \rightarrow A$

$$R_1 = \alpha^+ \quad \text{and} \quad A \subset A \quad \alpha \rightarrow A \quad A \subset A$$

$$R_2 = (R - \alpha^+) \cup \alpha \quad A \subset A \quad A \neq A$$

DP.

This guarantees lossless 2NF decomposition.

Ex. $R(A, B, C, D, E), F = \{ AB \rightarrow C, C \rightarrow B, A \rightarrow D \}$

$$CK_S = \{ ABE, ACE \} \quad N.P.AHr = \{ D \}$$

Not in 2NF, due to $A \rightarrow D$.

2NF Violation: $A \rightarrow D$

$$R_1 = A^+ = AD \quad F_1 = \{ A \rightarrow D \} \quad (\text{In 2NF})$$

$$R_2 = ABCE \quad F_2 = \{ AB \rightarrow C, C \rightarrow B \}$$

$$CK_{S_2} = \{ ABE, ACE \} \therefore \text{No NPAHr}, \text{ so in 2NF}$$

$A \in T_1$

$A \in T_2$

$A \in T_3$

So, for every rel. R we can guarantee a lossless dependency preserving 2NF decomposition.

- * In a database, only insertion or modification can violate a FD.
- * A ~~depe~~ decomposition is Dependency Preserving iff all the FDs in the database can be enforced ~~any~~ without any join, i.e. with help of any individual table (not two or more).
- * In general we prefer 3NF, (not BCNF), since we can guarantee a lossless D.P. Decomposition.

1. What is the name of the person in the picture?

2. Who is the person in the picture?

3. What is the person's job?

4. What is the person's name?

5. What is the person's job?

6. What is the person's name?

7. What is the person's job?

8. What is the person's name?

9. What is the person's job?

10. What is the person's name?

11. What is the person's job?

12. What is the person's name?

13. What is the person's job?

14. What is the person's name?

15. What is the person's job?

16. What is the person's name?

17. What is the person's job?

18. What is the person's name?

19. What is the person's job?

20. What is the person's name?

21. What is the person's job?

22. What is the person's name?

23. What is the person's job?

24. What is the person's name?

25. What is the person's job?

26. What is the person's name?

27. What is the person's job?

28. What is the person's name?

Queries

- Schema tells us the name of relation, the name of the attributes and their types.

Instance on the other hand has the actual tuples of the relation.

Schema is also called Intension.

Instance is also called Extension.

- E.F. Codd in his paper gave "Relational Model", along with two Query Languages (both theoretical):

- (1) Relational Algebra
- (2) Tuple Relational Calculus.

- Query Languages: Specialised languages designed specifically to ask questions related to data in the database.

- The Query Languages of interest for Relational Model are:

- (1) Relational Algebra
 - (2) Tuple Relational Calculus
 - (3) Structured Query Language.
- Theoretical
- Practical.

- Procedural QL vs Declarative QL

A function tells us the mapping of input to output. An algorithm gives step-by-step instruction for getting the output for a input.

Ex. f: Sort Inp: 2, 1, 0, 3 → Out: 1, 2, 3

Algo: MergeSort : Do this, then that, then another.

Similarly, Procedural QL require the user to tell what they want & how to get it. Relational Algebra is Procedural QL.

Appendix

Declarative QL only requires the user to tell what they want. Relational Calculus is such language. SQL too.

So, for user, Declarative QL is better. For DBMS software, Procedural QL is better.

In SQL actually, the user query (which is declarative) is converted to equivalent relational algebra query (which is procedural) by Query Processor & further optimized by Query Optimizer.

* Relational Algebra

| Tables to Relational Algebra | Algebra and ADT |
|-----------------------------------|---|
| Schemas | Variables |
| Instances | Values |
| Operators: Basic ↓ Derived: | Operations: Basic: +, - ↓ Derived: ×, ÷, etc. |

$$\text{Ex. } R_1 \cup R_2$$

Expr. is eval. based on instances

$$10 \leftarrow$$

$$\text{Ex. } a(b+c)$$

Expr. will be eval. based on value of a, b.

For some query Q in RA, inputs are instances, and output would be some instance. Query Q itself is written in terms of schema & operators.

Sometime for complex queries we require intermediate results, that also is instances.

- In a relation, the position of attributes are fixed and they can be referred to by their positions.

Ex. $R(A, B, C)$

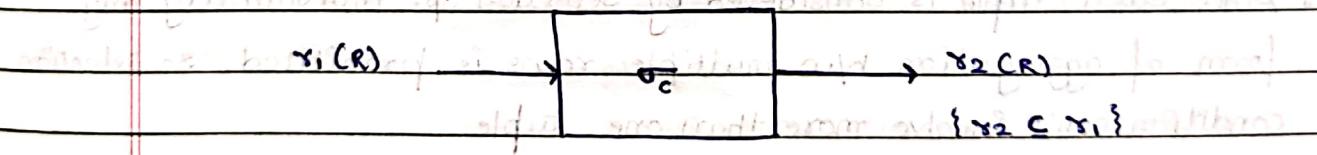
$$R.A \equiv R.1, R.B \equiv R.2, R.C \equiv R.3.$$

- Basic Operators in RA: Select (σ), Project (Π), Cartesian Product (\times), Set Difference ($-$), Union (\cup), Rename (ρ)

- Derived Operators in RA: Set Intersection (\cap), Natural Join (\bowtie), Conditional Join (\bowtie_c), Division (\div), etc.

Selection Operator (σ_c)

- Unary operator, which, based on, some condition c , selects tuples from input instance & discards the others.



σ_c here selects tuples from $r_1(R)$ which satisfy the condition c .

Ex. $R(\text{roll}, \text{name}, \text{gpa})$

- $\sigma_{\text{gpa} > 7.5}$: Select all tuples with $\text{gpa} > 7.5$.
But this isn't complete query.
- $\sigma_{\text{gpa} > 7.5}(R)$: Select all tuples from R with attr $\text{gpa} > 7.5$.
The output instance would have schema: $R(\text{roll}, \text{name}, \text{gpa})$.
- $S = \sigma_{\text{gpa} > 7.5}(R)$: Select all tuples from R (or instance of schema R) with $\text{gpa} > 7.5$. The outp. instance would have schema: $S(\text{roll}, \text{name}, \text{gpa})$.

So selection operator without assignment doesn't change schema. Also it doesn't change the degree, but cardinality of output may differ from input.

Selection operator applies on each tuple individually in the input instance.

The condition σ_c in σ_c is propositional logic formula, and can have connectives such as \wedge, \vee, \neg . It can have \neq conditions on attributes of a tuple, such as $>, <, \geq, \leq, \neq, \neg (x=y)$.

Ex. User (uid, name, age, pop).

Q: Select users with popularity at least 0.9, and age under 10 or above 12.

$$\sigma_{(pop \geq 0.9) \wedge (age < 10 \vee age \geq 12)} (\text{User})$$

Since each tuple is considered by selection op. individually any form of aggregation b/w multiple rows is prohibited. So selection condition can't involve more than one tuple.

$$\sigma_{pop > \max(pop)} (\text{User})$$

This is wrong. $\max(pop)$ requires more than one tuple at a time.

So, it does horizontal partitioning.

Projection Operator (Π_a)

- Unary operator which projects specified attributes from a given relation. If after projection, some tuples are duplicate it will only keep single copy of that tuple.

| Ex. R: | A | B | C | | A | C | |
|--------|----------|----|---|--|----------|---|--|
| | α | 10 | | $\rightarrow Q: \Pi_{(A,C)}(R)$ | α | 1 | |
| | α | 20 | | $\rightarrow \Pi_{(A,C)}(R) \rightarrow \beta$ | β | 1 | |
| | B | 30 | 1 | $\rightarrow S: \Pi_{(C,A)}(R)$ | β | 2 | |
| | B | 40 | 2 | | | | |

- So, it is doing vertical partitioning.
- Again, assignment such as, $S = \Pi_{(A,C)}(R)$, would change the name of output's schema, from R to S.
- Importantly, any operation in RA on schema & instances stored in DR, doesn't affect the original tables, rather all outputs are temporary copy.
- If not assigned to any var, the result of a query won't be further usable. (If assigned, the var name can be used in further queries.)
- The output instance would have attributes in the order as specified in projection operator's attribute list.

Ex. $R(A, B, C)$ $S = \Pi_{(C, A)}(R)$ $S(C, A)$

- So, the degree of output of $\Pi_{(a)}$ would be same as $|a|$ & may differ from imp's ~~can~~ degree. Cardinality of out. may also differ from imp's (in case of duplicate tuples).

Ex. To make sure the cardinality of projected instance is necessarily same as input instance, what can you say about the attribute list of proj. operators?

The attribute list must be some superkey of the input's schema.

Ex. $R(a, b, c)$. FDs = $\{a \rightarrow b, bc \rightarrow a\}$. Which of the following queries are guaranteed to return same no. of tuples for any input instance r of R ?

- (i) $\Pi_a(r)$
- (ii) $\Pi_{a,c}(r)$
- (iii) $\Pi_c(r)$
- (iv) $\Pi_b(r)$
- (v) $\Pi_{a,b}(r)$
- (vi) $\Pi_{a,b,c}(r)$

$CK_R = \{ac, bc\}$ so, (ii), (vi) are guaranteed to have same no. of ~~other~~ tuples as input.

Properties of Selection & Projection:

1. Selection is commutative: $\sigma_A(\sigma_B(R)) \equiv \sigma_B(\sigma_A(R))$

since ultimately we have tuples which pass both conditions.

2. Projection is not commutative: $\Pi_{C_1}(\Pi_{C_1, C_2}(R)) \not\equiv \Pi_{C_1, C_2}(\Pi_{C_1}(R))$

3. If $X, Y \subseteq R$; $X \subseteq Y$, then, $\Pi_X(\Pi_Y(R)) \equiv \Pi_X(R)$

Ex. $R(\text{sid}, \text{surname}, \text{rating}, \text{age})$

Select the name & rating of all students with rating > 8 .

$\Pi_{(\text{surname}, \text{rating})} (\sigma_{\text{rating} > 8}(R))$

Q. We can perform the two operations on relation R:

(1) Select tuples based on C

(2) Project attributes of list A

In either order iff C uses attr. only in the list A.

i.e. $\sigma_C(\Pi_A(R)) \equiv \Pi_A(\sigma_C(R))$

iff all attr. condition C uses are present in list A.

Ex. In general which transformation (substitution) is valid:

(A) $\Pi_A(\sigma_C(R)) \rightarrow \sigma_C(\Pi_A(R))$ (Subs. LHS with RHS)

(B) $\sigma_C(\Pi_A(R)) \rightarrow \Pi_A(\sigma_C(R))$

B is valid. If $\sigma_C(\Pi_A(R))$ is a valid query, it means C uses only attr. available in A; so, $\Pi_A(\sigma_C(R))$ would also be valid since

R is superset of A. (This is a substitution of which)

Union Operation (+), Set Difference (-), Intersection (n)

• Binary operations which just behave like set operations. Since instance of any relation simply is a set of tuples, these operations behave simply as set operations.

• For these operations b/w two relations, R & S, they must have compatible schema (a.k.a. Union Compatible Schema).

Relation R is Union Compatible with relation S, if

(1) Degree of both relations is same.

(2) The domain of ith attribute of both relation must be same.
(Name of attr. may differ)

| A | B | C | D |
|---|---|---|---|
| 1 | 2 | 3 | 4 |
| 1 | 2 | 3 | 4 |
| 2 | 1 | 4 | 3 |
| 3 | 4 | 1 | 2 |
| 4 | 3 | 2 | 1 |

| RUS: | A | B | R-S: A | B | S-R: B | D |
|----------|---|-----------------------------------|--------|-------------------|--------|---|
| α | 1 | A first off α with insight | 1 | B | 3 | |
| α | 2 | plus addition of B | 1 | final addition of | | |
| B | 1 | + 100 = 100 | (100) | 50 | | |
| B | 3 | then with addition after No. | 99 | 11 | | |

SNR: (8) D

$$\alpha_B \leftarrow ((\alpha_B)_{AT})^{\frac{1}{2}} \leftarrow ((\alpha_B)_{AT})^{\frac{1}{2}} \cdot \pi - (\alpha)$$

$$f((\alpha \beta)_{ATT}) \leftarrow ((\alpha)_{ATT}) \sqcup (\beta)$$

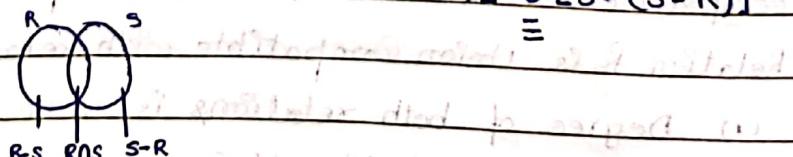
- The schema of output instance 8 will be same as the schema of first operand (A,B,C) so it will have 3 columns.
 - Union is commutative (we don't worry about output schema, but only output data).

• Set Difference is not Commutative : $R - S \neq S - R$
(since data may differ)

- Intersection is also commutative (although schema of o/p may differ, but data is same): $R \cap S = S \cap R$

Intersection is actually a derived operator:

$$RNS \equiv [(R+S) - (R-S)] - (S-R) \equiv [R - (R-S)] + [S - (S-R)]$$



- Cross Product (\times)

- Binary operation

Ex. Student.Course (sid, sname, cid, cname).

Get • Sid, Sname of all students who have registered in "DB" course but not in "Algo" course.

$$R = \Pi_{(sid, sname)} (\sigma_{cname = "DB"} (S \cdot C)) - (\Pi_{(sid, sname)} (\sigma_{cname = "Algo"} (S \cdot C)))$$

- Cross Product (\times)

- Binary operation .

- $R \times S$, pairs every tuple of R with every tuple of S .
- The schema of output instance for $R \times S$ would consist of all attrs. of R with all attrs. of S in order.
- suppose $R(A, B)$, $S(C, B, D)$, then there will be confusion in output of $R \times S$, say $X(A, B, C, B, D)$.

Convention 1: $X(R.A, R.B, S.C, S.B, S.D)$ or

$X(A, R.B, C, S.B, D)$

Convention 2: $X(A, 1, C, 4, D)$ or

$(A, 1, 2, 3, 4, 5)$.

(we generally follow convention 1)

| Ex. R: | | | | Rxs: | | | |
|--------|--|--|--|------|---|-----|-----|
| | | | | A | B | R.B | S.B |
| | | | | 1 | 1 | a | + |
| | | | | 1 | 1 | b | - |
| S: | | | | 1 | 1 | a | - |
| | | | | 1 | 1 | a | + |
| | | | | 1 | 1 | b | - |
| | | | | 1 | 1 | a | - |

- For $R \times S$,

$$|R \times S| = |R| \times |S|$$

$$\& \deg(R \times S) = \deg(R) + \deg(S)$$

- Cross product is commutative: $R \times S = S \times R$, since same data, even though diff. schema.

• Rename Operator (ρ_s) -

- Unary operator

- Renames the given input & after creating its copy.

$\rho_s(R)$: action of ρ to output grove string - Ex:

$\rho_s(R)$: $R: [A \times B] \rightarrow s: [A \rightarrow C]$

- Renaming can be done for schema, and optionally for attributes also.

Both schema renaming & attr. renaming are optional.

Ex. $R(A, B, C)$

Change Schema Name (copy in main m/m): $\rho_s(R)$

Change Schema & all attr. (\rightarrow): $\rho_{s \rightarrow} (R)$

Change Sch. name & some attr (\rightarrow): $\rho_{s \rightarrow 1, c \rightarrow 2} (R)$

Change only attr name (\rightarrow): $\rho_{(c \rightarrow 1, c \rightarrow 2)} (R)$

Since renaming just copies the rel. R to main m/m & doesn't affect on disk, we can skip renaming of ρ (this won't create two copies of ρ in main m/m).

Ex. $\text{Follows}(\text{uid}_1, \text{uid}_2)$. ($\langle x, y \rangle \in \text{Follows}$ if x follows y .)

(a) Find all such pairs $\langle x, y \rangle$ where x follows y & y follows x .

$$\Pi_{(F_1.\text{uid}_1, F_1.\text{uid}_2)} \left(\sigma_{(F_1.\text{uid}_1 = F_2.\text{uid}_2) \wedge (F_1.\text{uid}_2 = F_2.\text{uid}_1)} \left(P_{F_1}(\text{Follows}) \times P_{F_2}(\text{Follows}) \right) \right)$$

(b) Find uid of those who follow at least two people.

$$\Pi_{F_1.\text{uid}_1} \left(\sigma_{(F_1.\text{uid}_1 = F_2.\text{uid}_1) \wedge (F_1.\text{uid}_2 \neq F_2.\text{uid}_2)} \left(P_{F_1}(\text{Follows}) \times P_{F_2}(\text{Follows}) \right) \right)$$

(c) Find uid of those who follow exactly one person.

$$P_{(\text{uid})} \left(\overline{\Pi}(\text{Follows}) \right) - P_{(\text{uid})} \left(\begin{array}{l} \text{Follows} \\ \text{At least two} \\ \text{people} \end{array} \right)$$

Join Operations

- Usually performing cross product (\times) such as $A \times B$ results in a lot of useless tuples & rarely tells us the required information.
- So, cross product is generally followed by selection (σ) or projection (Π) operations.
- Since cross product followed by selection is so common, we've created new operator : $\text{Join}(\bowtie)$.
- Cross Prod followed by Selection = Conditional Join
 $\{ \sigma_C(R \times S) \} \quad \{ R \bowtie_C S \}$
- There are two types of Join Operators:
 - Inner Join (or Join) : Conditional Join / Theta Join (\bowtie_C), EquiJoin, Natural Join (\bowtie)

(2) Outer Join: Left OJ (\bowtie_L), Right OJ (\bowtie_R),

Full OJ (\bowtie_F)

- Conditional Join (\bowtie_C)
 - Pairs rows of R with rows of S which satisfies the condition(s): $R \bowtie_C S$
 - Practically the DBMS, doesn't convert \bowtie_C to $\sigma_C(x)$. There are join algorithms which directly perform the operation.

| Ex. R: | A | B | C | S: | D | E |
|--------|---|---|---|----|---|---|
| | 1 | 2 | 4 | 3 | 1 | 5 |
| | 2 | 5 | 5 | | 6 | 2 |
| | 3 | 8 | 6 | 9 | | |

| $R \bowtie_C S:$ ($R.B < S.D$) | A | B | C | D | E |
|-------------------------------------|---|---|---|---|---|
| | 1 | 2 | 4 | 3 | 1 |
| | | 1 | 2 | 4 | 6 |
| | 2 | 5 | 5 | 3 | 1 |

So, R has a "dangling row", since it doesn't pair with any row of S, (3, 8, 6).

Equi Join

Conditional Join with only equality conditions.

Ex. $R \bowtie_C (R.A = S.C) \wedge (R.B = S.E) S$

$R \bowtie_C (R.A = S.C) \wedge (R.B = S.E) S$

No other conditional operator is allowed.

Natural Join (\bowtie)

- Performs join based on only the name of attributes (not the domain).
- Say $R(A, B, C)$ & $S(B, C, D)$ & we write: $R \bowtie S$, then it means join R & S based on equality of all common attributes, here. (B, C) .

The schema of $R \bowtie S$ would be (A, B, C, D) , (i.e. no redundant attribute).

So, $R \bowtie S$ would behave as $R \bowtie S$
 $\quad\quad\quad (R.B = S.B) \wedge$
 $\quad\quad\quad (R.C = S.C)$

- If no common attr. b/w R & S by name: $R \bowtie S = R \times S$
- $(R \bowtie S) \bowtie T \equiv R \bowtie (S \bowtie T)$

Division (\div)

- In Number Theory $\frac{N}{D} = Q$, Q is the largest integer such that

- The Division operation in Relation Algebra behaves quite similarly, for relation R & S , $R/S = Q$, where Q is the largest relation such that $Q \times S \subseteq R$.

$$\text{Ex. } R(a,b) = \{(1,1), (1,4), (2,1), (2,2), (2,3), (2,4), (3,1), (3,3), (3,4)\}$$

$$S(b) = \{(1), (3)\}$$

$Q(R/S)$ has to satisfy: $Q \times S = R$. \therefore Schema for Q must be (a) .

$$R/S = \{(2,1), (3,1)\} = Q$$

Let's check. $Q(a) \times S(b) = \{(2,1), (2,3), (3,1), (3,3)\}$

So, this is subset of R .

Also, this Q is largest possible.

$$\text{Ex. } R(A, B, C, D, E) = \{(\alpha, a, \alpha, a, 1), (\alpha, a, \beta, a, 1), (\alpha, a, \beta, b, 1), \\ (\beta, a, \gamma, a, 1), (\gamma, a, \beta, a, 1), (\gamma, a, \beta, b, 1), (\gamma, a, \beta, b, 1)\}$$

$$S(D, E) = \{(a, 1), (b, 1)\}$$

$$Q = \{(\alpha, a, \beta), (\gamma, a, \beta)\}$$

- For R/S to be defined, the attribute set S must be subset R of attribute set S .

Ex. Division operation is valid for:

- (A) $R(A, B, C) \div S(B, C)$ and student
- (B) $R(A, B) \div S(B, C)$ and student

Ex. For any relation R, S what is $(R \times S) \div S$

Ex. If $Y \times S \subseteq R$, then $R/S = Y$?

No. Y should be the largest set, such that

$$Y \times S \subseteq R \quad (\text{Ans}) \Rightarrow Y = \text{TM}(E \times S)$$

$\therefore (Y \times S) \subseteq R$

- We can rewrite Division operator in terms of basic operators.

$$R/S = \prod_{\substack{(R-S) \\ \text{ways}}} (\text{R}) - \left\{ \prod_{\substack{(R-S) \\ \text{ways}}} \left[\left(\prod_{\substack{(R-S) \\ \text{ways}}} (R \times S) \right) - R \right] \right\}$$

attribute set

(Healy Implementation)

Let $S[i]$ represent i^{th} tuple in S .

$$R/S = \prod_{i=1}^{n(S)} \prod_{(R-S)} (\text{R}) \quad (n(S)) = n(S)$$

$$\text{OR } R - \{S[i]\} \quad (i, \text{const}) = 2^n$$

$$R/S = \bigcap_{\text{v.bies}} \prod_{(R-S)} \left(\overline{\prod_{(R-S)} (R \times S)} (R) \right)$$

(Maier's Implementation)

$$(R \times S) \cap (R \times T) = (R \times (S \cap T))$$

$$\{ (1, a), (2, b), (3, c) \} \cap \{ (1, a), (2, b), (4, d) \} = \{ (1, a), (2, b) \}$$

$$\{ (1, a), (2, b), (3, c) \} = \{ (1, a) \}$$

$$\{ (1, a), (2, b), (3, c), (4, d) \} = \emptyset$$

Ex. Given R, S , $|R|=m$, $|S|=n$, ($m, n > 0$), then what is the max. & min. cardinality of $R \div S$. using IT

Min: 0, when no $(R-S)$, such that all R tuples are satisfied.

Max: $\left\lfloor \frac{|R|}{|S|} \right\rfloor$. When all values in $(R-S)$ schema satisfy all R tuples.

$$\text{Ex. } R = \{(A, 2), (A, 1), (A, 3), (B, 1), (B, 2)\}$$

$$\begin{matrix} (A, B) \\ S(B) \end{matrix} = \{(1), (2)\}$$

$$R/S = \{(A), (B)\}$$

$$\therefore |R/S| = \left\lfloor \frac{5}{2} \right\rfloor = 2$$

$$|R|=m, (m>0)$$

Ex.. In R/S , if $|S|=0$, then $|R \div S|$?

Healy's
We can use Kratzer's defⁿ. $\prod_{(R-S)}(R) - \left[\prod_{(R-S)} \left(\prod_{R-S}^{} (R) \times S - R \right) \right]$

The cardinality $|\prod_{(R-S)}(R) \times S| = |\prod_{(R-S)}(R)| \times |S| = |\prod_{(R-S)}(R)| \times 0 = 0$.

The card. $|R-S| = 0$.

$$\therefore \text{Ultimately, } \prod_{(R-S)}(R) - \left[\dots \right] = \prod_{(R-S)}(R).$$

{0 card.}

So, at max, this can be $= |R| = m$.

at min, this can be $= |\prod_{(R-S)}(R)| = 1$

If $|R|=0$

Then $|R/S| = 0$.

- Outer Join

- It includes the dangling tuples with help of NULL.

- Left Outer Join (IX)

$R \text{ IX } S$ includes the dangling tuples of R .

- Right Outer Join (XI)

$R \text{ XI } S$ includes the dangling tuples of S .

- Full Outer Join (XJ)

$R \text{ XJ } S$ includes the dangling tuples of both R & S .

Ex. $R: A$

| A | B |
|---|---|
| 1 | 2 |
| 3 | 4 |

$S: B | C$

| B | C |
|---|---|
| 2 | 5 |
| 6 | 9 |

| A | B | C |
|---|---|---|
| 1 | 2 | 5 |

| A | B | C |
|---|---|------|
| 1 | 2 | 5 |
| 3 | 4 | NULL |

| A | B | C |
|------|---|---|
| 1 | 2 | 5 |
| NULL | 6 | 9 |

| A | B | C |
|---|---|------|
| 1 | 2 | 5 |
| 3 | 4 | NULL |

- By default Outer Join = Full Outer Join (IX , XI , or XJ)

- By default, we first do natural join then add dangling tuples for Outer Join.

- We can also have conditional outer join. Include all rows which satisfy condition, then include dangling tuples (based on left, right, full).

* Tuple Relational Calculus

In RA, we write queries which tell, what we want & the order in which operations must be performed to get the results. So, RA is a procedural language.

TRC queries just tell what we want, and not the order of the operations to be performed to get that. So TRC is Declarative Language.

In TRC, we represent tuples as variable, so every variable in TRC can be replaced with tuple.

To which relation the tuple belongs can be indicated as:

$t \in R$ or $R(t)$

The attribute of tuple can be represented as : $t.A$ or $t[A]$.

The queries in TRC are of the form: $\{t \mid P(t)\}$ where $P(t)$ is predicate on t .

There are two ways of writing TRC queries:

(1) We first come up with logic, then get results.

(2) We first picture the results, then get derive logic for it.

Ex. Write SELECT operation of RA (\rightarrow) in TRC.

say, $R(A, B)$ RA Query: $\sigma_{B=17}(R)$

TRC Query: $\{t \mid t \in R \wedge t[B] = 17\}$

Ex. Sailors (sid, sname, rating, age)

Boats (bid, bname, color)

Reserves (sid, bid, day)

(1) Find all Sailors w/ rating > 7.

RA: Rating > 7 (S)

TRC: {t | t.es AND t[rating] > 7}

Ex. Write Projection o/pⁿ int TRC.

say, R(A, B, C)

RA: TA(R)

RA: TTA, C(R)

TRC: {t.A | t ∈ R} → TRC: {t.A, t.C | t ∈ R}

This is way 1.

Now, suppose r(R) is an instance:

| | A | B | C |
|----------|----|---|------------------|
| α | 10 | 1 | TTA(B α) |
| α | 20 | 1 | |
| B | 30 | 1 | |

why (α) is there in r' ? Because there exists

a tuple in r such that $t.A = \alpha$.

Similarly for (B,) tuple in result r' .

So,

$$\{t | \exists_{x \in r} (x.A = t.A)\}$$

there is no info, where tuple t belongs to, but since with 't' we only use attr. A, we can say schema of t is (A,) only.

Suppose we wanted query TTA, C(r):

$$\{t | \exists_{x \in r} [(x.A = t.A) \wedge (x.C = t.C)]\}$$

∴ Schema of t is (A, C)

Ex. $R(A, B)$ $\sigma_{B=17}(R)$ Way 1: $\{t \mid t \in R \wedge t.B = 17\}$ Way 2: $\{t \mid \exists x \in R (t.A = x.A \wedge t.B = x.B) \wedge t.B = 17\}$ ↑ this t represents the tuple in result table.Schema of t : (A, B) Ex. $R(A, B) \cap (\sigma_{B=17}(R))$ Way 1: $\{t.A \mid t \in R \wedge t.B = 17\}$ Way 2: $\{t \mid \exists x \in R (t.A = x.A \wedge x.B = 17)\}$ Ex. $R(A, B), S(C, D)$ (R & S are union compatible)① $R \cup S$ W1: $\{t \mid t \in R \vee t \in S\}$ W2: $\{t \mid \exists x \in R (t.A = x.A \wedge t.B = x.B) \vee \exists x \in S (t.A = x.C \wedge t.B = x.D)\}$ ② $R \cap S$ W1: $\{t \mid t \in R \wedge t \in S\}$ W2: $\{t \mid \exists x \in R (t.A = x.A \wedge t.B = x.B) \wedge \exists x \in S (t.A = x.C \wedge t.B = x.D)\}$ ③ $R - S$ W1: $\{t \mid t \in R \wedge t \notin S\}$ W2: $\{t \mid \exists x \in R (t.A = x.A \wedge t.B = x.B) \wedge \forall x \in S (t.A \neq x.C \wedge t.B \neq x.D)\}$ Ex. $R(A, B) \quad S(C, D, E) \quad R \times S$ W1: $\{t_1, t_2 \mid t_1 \in R \wedge t_2 \in S\}$ W2: $\{t \mid \exists x \in R \exists y \in S (t.A = x.A \wedge t.B = x.B \wedge t.C = y.C \wedge t.D = y.D \wedge t.E = y.E)\}$

Note: In Way 2, the $P(t)$ predicate answers the question, "why the tuple t is in the resultant relation?"

Ex. $R(A, B) \quad S(B)$ $R \sqsubseteq S$

$W_1: \{ t_1 \forall x \in S \exists y \in R (x.B = y.B \wedge t_1.A = y.A) \}$

$W_2: \{ t_1.A \sqsubseteq t_2 \wedge \forall x \in S \exists y \in R (x.B = y.B \wedge t_1.A = y.A) \}$

Ex. $R(A, B) \quad S(A, D)$ $R \bowtie S$

$W_1: \{ t_1.A, t_1.B, t_2.D \mid t_1 \in R \wedge t_2 \in S \wedge t_1.A = t_2.A \}$

$W_2: \{ t_1 \mid \exists x \in R (t_1.A = x.A \wedge t_1.B = x.B) \wedge \exists x \in S (t_1.A = x.A \wedge t_1.D = x.D) \}$

- The benefit of way 2 of TRC queries is that there's always only one free variable.

Safe vs Unsafe Query

Tuples in the database instance are real data.

If a query gives some data which isn't present in the database (or can't be derived from it), then the query is unsafe.

Else the query is safe.

Also, if query gives infinite result, then its also unsafe (since database instance is finite).

Ex. $\{ t_1 \mid t_1 \notin R \}$

Query gives inf. result \rightarrow Query is unsafe

{ One way }

Expressive Power of TRC

We have shown that for every Rel. Algebra operator we have equivalent TRC Query (safe query), so,

$RA \equiv \text{Safe TRC}$

in terms of expressive power.

{ TRC (safe + unsafe) has higher expressive power than RA }

* Structured Query Language (SQL)

Declarative Query Language.

- It has multiple subsets based on query type:
 - Data Definition Language:** For creation, deletion & modification of defⁿ of tables & view. Constraints can also be included here.
 - Data Manipulation Language:** For insertion, deletion, modification of data. Also to query / retrieve data.

(Note: In GATE, retrieval of data is of concern.)

- Although motivated from theoretical query lang., such as RA & TRC, SQL has some major differences.
 - Two rows couldn't be entirely same.
 - There is exactly one PK (at a time).
 - There is ≥ 1 CK & ≥ 1 SK for relation.
 - Relational is a set of tuples.
 - Don't have NULL in table (except for outer join).

Set

Multiset

- Unordered.

- Duplicates not allowed

- Duplicates allowed

$$\{1, 2, 2\} = \{1, 2\}$$

$$\{1, 2, 2\} \neq \{1, 2\}$$

$$|\{1, 2, 2\}| = |\{1, 2\}| = 2$$

$$|\{1, 2, 2\}| = 3$$

In SQL, Table is equivalent to multiset.

In SQL,

- Two rows may have exactly same value (duplicates allowed).
- May/Maynot have PK, CK, SK.
- NUL values are allowed.

Note: Although SQL keywords are case-insensitive, we'll write them in BLOCK letters.

Basic SQL Query:

SELECT A₁, A₂, A₃... ← Attribute Selection

FROM R₁, R₂, R₃... ← Conceptually CrossProd of R₁ × R₂ × R₃...

WHERE -condition- ← [Optional] tuple selection condition

(SELECT, FROM are mandatory)

Ex. R(A,B), S(C,D)

① **SELECT A, D**

FROM R, S

WHERE (A>10) AND (C>5)

② **SELECT DISTINCT A, D**

FROM R, S

WHERE (A>10) AND (C>5)

① similar to : $\Pi_{A,D}(\overline{A>10} \wedge C>5)(R \times S)$, but ① is SQL query & will allow duplicates.

② is exactly same as: $\Pi_{A,D}(\overline{A>10} \wedge C>5)(R \times S)$, since it doesn't allow duplicates by use of "DISTINCT".

In SQL also, if you say (X,Y,...) is primary key, then none of them can have NULL value. (Entity IC)

And obviously, it can't have exactly same two tuples.

Ex. Suppose in above query, R has PK A & S has PK C,D, do we still need DISTINCT in query to avoid duplicate tuples?

Yes. Still dup. tuples can be there in result (without DISTINCT). Ex.

R: A B S: C D

110 1

21 1

120 2

22 1

Ex. Same Query, But "SELECT A,C,D..." Do we need DISTINCT keyword to avoid duplicate.

No. Since A is PK in R & C,D is PK in S, no duplicate tuple will result from this query.

Note: ~~SELECT A,A2,A3...~~ This can be thought of as:
 FROM R₁,R₂,... projecting A₁,A₂,A₃ from the
~~WHERE cond~~ crossprod. R₁ × R₂... based on
 condition 'cond'.

Putting * after SELECT, means project all attributes in R₁ × R₂ × ...

Ex. SELECT * FROM Student S;

would give all tuples, all attributes of Student table.

When we run SQL data retrieval query, original data doesn't change. It exists on disk. It is copied to main m/m & op's are performed on the copy.

Once we write "... FROM Student S...", i.e. alias Student to S, we can't use "Student" in the query, only the alias "S" can be used.

Ex. ~~"SELECT Student.Name FROM Student S"~~ is wrong.

The way we understand SQL query evaluation is not how it actually happens.

Ex. SELECT A FROM R,S WHERE C.

Our Understanding: ① X = R × S → ② Y = σ_C(X) → ③ Z = π_A(Y)

This isn't what SQL does internally, although this is how we understand it.

But the actual op " will produce the same result
(but it uses more optimized algorithms).

Ex. Sailors (sid, sname, rating, age)

Boats (bid, bname, color)

Reserves (sid, bid, day).

(1) Find name & age of all sailors.

`SELECT sname, age FROM Sailors;`

(2) Find all sailors above rating 7.

`SELECT * FROM Sailors WHERE rating > 7;`

DISTINCT : Whenever used after `SELECT` keyword, it works on the entire tuple for each tuple in result.

Ex. `SELECT DISTINCT A,C FROM ...`

would compare (A,C) for each tuple to be distinct, not A, not C separately.

i.e. two tuples are non-distinct t_1, t_2 , iff:

$$t_1[A] = t_2[A] \text{ AND } t_1[C] = t_2[C]$$

. AS : Optional keyword for aliasing.

"... FROM Sailors S" \equiv "... FROM Sailors AS S"

Helpful when same table used two or more times in single query.

Ex. contd.

(3) Find name of sailors who have reserved boat with id 103

`SELECT sname FROM Sailors S, Reserves R WHERE`

$$S.sid = R.sid \text{ AND } R.bid = 103;$$

(4) Find Sailor names who have reserved at least one boat.

SELECT sname FROM Sailors S, Reserves R WHERE S.sid = R.sid;

(5) SELECT S.sid FROM Sailors S, Reserves R WHERE S.sid = R.sid;

In this query, does adding DISTINCT make any difference?

Yes. Although S.sid is PK, but after cross-prod with R, we have multiple entries repetition of S.sid.

| Ex. S: | Sid | sname... | R: | sid | bid... |
|--------|-----|----------|----|-----|--------|
| | 212 | A | | 212 | 103 |
| | 231 | B | | 212 | 104 |
| | | | | 231 | 103 |

(6) Find sname of Sailors who've reserved red boat.

SELECT sname FROM Sailors S, Reserves R, Boats B WHERE

S.sid = R.sid AND R.bid = B.bid AND B.color = "red";

(7) Find color of boats reserved by sailor named "Lubber".

SELECT color FROM Sailors S, Reserves R, Boats B WHERE

S.sname = "Lubber" AND S.sid = R.sid AND R.bid = B.bid

In SQL Union is performed using Keyword UNION, difference

with EXCEPT (or MINUS in ORACLE), and intersection using

INTERSECT. These keywords also require union compatible

table. Although the input tables may have duplicates, the output will only have unique tuples.

Ex. R: A S: B

a₁

a₂

a₂

a₃

a₂

a₃

a₂

a₃

• $(\text{SELECT } * \text{ FROM } R) \text{ UNION } (\text{SELECT } * \text{ FROM } S);$

R-S: A

a₁

a₂

a₃

• $(\text{SELECT } * \text{ FROM } R) \text{ INTERSECT } (\text{SELECT } * \text{ FROM } S);$

R-S: A

a₁

a₂

• $(\text{SELECT } * \text{ FROM } R) \text{ EXCEPT } (\text{SELECT } * \text{ FROM } S);$

R-S: A

a₁

There are also versions of these operations which don't remove duplicates. UNION ALL, INTERSECT ALL, EXCEPT ALL. They perform frequency based operations.

• UnionAll \Rightarrow A with. IntersectAll \Rightarrow A

a₁, a₂, a₃, a₄, a₅, a₆, a₇, a₈

a₁, a₂, a₃, a₄, a₅, a₆, a₇, a₈

a₂ {at most two a₂'s are}

a₂ {common. S doesn't have a}

a₃ {third.}

• Except All \Rightarrow A

a₁, a₂, a₃, a₄, a₅, a₆, a₇, a₈

a₂ {Two a₂ get removed, one}

remaining.

• Can't write R UNION S directly. Invalid.

$(\text{SEL } * \text{ FROM } R) \text{ UNION } (\text{SEL } * \text{ FROM } S);$

Note: Take care of union, intersection & except when performing on Key vs non-Key attributes.

Ex. [Contd.]

(8) Find sid of sailors who've reserved a red or a green boat.

① $\text{SELECT sid FROM Reserves R, Boats B WHERE (R.bid = B.bid)}$
 $\text{AND (B.color = "red" OR B.color = "green");}$
 {This will give non-distinct}

② $(\text{SELECT sid FROM } \dots \text{ WHERE R.bid = B.bid AND B.color = "red"})$
 UNION

$(\text{SELECT sid FROM } \dots \text{ WHERE } \dots \text{ AND B.color = "green"})$
 {This will give distinct results}

① with DISTINCT keyword would be equivalent to ②.

(9) Find sid of sailors who've reserved both a red & a green boat.

$(\text{SELECT sid FROM Reserves R, Boats B WHERE (R.bid = B.bid)}$

$\text{AND B.color = "red")}$

INTERSECTS

$(\text{SELECT sid } \dots \text{ FROM } \dots \text{ WHERE } \dots \text{ AND B.color = "green"})$

(10) Find Sname of sailors who've reserved both a red & a green boat.

NOW the previous approach won't work, since Sname isn't key attr.

Suppose $\langle \text{sid:112, sname:A} \rangle$ has hired red boat & $\langle \text{sid:223, sname:A} \rangle$ has hired green. If we intersect the (Sname hired red) & (Sname hired green), we would get A in result, which is wrong.

So, one approach is using "nested query", which will come later.

Another is to intersect $\langle \text{sid, sname} \rangle$ pairs.

$(\text{SELECT sid, sname FROM Sailors S, Res R, Boats B WHERE S.sid = R.sid AND R.bid = B.bid AND B.color = "red")$

INTERSECT $(\dots \text{ B.color = "green"})$