

Unit 1

DBMS provide ability to define, store and manipulate data.

1) Self-describing, 2) Independence of programs and data

3) Data Abstraction, 4) Support multiple views of data

5) Sharing of data and multi-user transaction processing.

There are 2 main categories of users of database

1) Actors on the scene. \Rightarrow Database Administrators, Database Designers, End users
System Analysts and Application Developers

2) Workers behind the scene \Rightarrow System Designers, Tool Developers, Operators and Maintenance Personnel.

Advantages - (i) Control Redundancy, (ii) Restrict unauthorized access
(iii) Persistent Storage, (iv) Efficient query processing
(v) Backup and Recovery Services, (vi) Multiple views to multiple users

Disadvantages - (i) Need fixed standards, (ii) Reduced Application development time
(iii) Need to maintain flexibility to change data structures.
(iv) Economies of scale (v) Current info must be available immediately

We shouldn't use database when

(i) High initial investment, (ii) DB, Application won't change, (iii) No multiple users

(iv) When real-time deadlines and DB overhead not affordable (v) Complex Data

Data Models = set of concepts to describe the structure of database, operations

for manipulating these structures and certain constraints

Conceptual = Way user perceives data

Physical = How data is stored in computer.

Implementation / Representational = In between both above. Mainly used

Three Schema Architecture

Internal = describe physical storage

Physical Data Model

\downarrow
Conceptual = structure and constraints for all users

Conceptual / Implementation

\downarrow
External = various user views.

Conceptual Data Model

Data Independence = when lower level schema changes only mappings to upper change

(i) **Logical DI** = when conceptual schema changes, external need not

(ii) **Physical DI** = when internal schema changes, conceptual need not.

DBMS Interfaces = Programming Language or User friendly.

Database System Utilities = Load data, back up data, monitor performance

sorting, user monitoring, data compression

Database System Environments

1) Centralized = Everything into 1. DBMS hardware, software, UI, Application Programs.

2) Client- Server = Client has proper UI, API to send to server via network
Server respond to client query.

Both must have proper client and server modules installed.

3) Three-tier = Extra layer called application layer added. as intermediate.

Stores web connectivity software and business logic

Client input → Partially Process → Server

Enhances security as client need this to access server and vice-versa

Entities - Relationship Model

Entity = Specific things or Objects in mini-World. Strong Weak

has attributes = Simple, Composite, Multi-Valued , Derived
 one of which is key

Entity set = collection of all entities

Value set = Domain of attributes

Relationship relates two or more distinct entities with specific meaning.

Relationship type = schema description (name, entities, type), 1:n, 1:1

Relationship set = Set of relationship instances

m:n, recursive

Cardinality ratio = ratio of maximum participation of participating entities

Alternative of this is (min,max) notation

Total participation = When ($\geq 1, \geq 1$) for the relation //

Partial participation = When ($0, \geq 0$) ✓

Weak entity = no key and identification depends exactly on another.

, Partial key

Unit 2

Relational Model :- This is what RDBMS follow to store data.

Relation = Table Rows = Tuples Column = Attributes (can be equal or unequal)

Key = Attribute that uniquely identifies each row in table = Primary Key

Schema = $\{R_1, R_2, \dots, R_n\}$ \Rightarrow Set of relations

Relation = $\{A_1, A_2, \dots, A_m\}$ \Rightarrow Set of attributes

Domain = data type or format that attribute has to follow

Relation State = Subset of Cartesian Product of domains of attributes. It is the populated table.

Relational Model was needed to remove all negatives of flat file. Flat file had a lot of redundant data and modifying data was tough. We split the single table db into relational db using normalisation.

Constraints of Relational Model

(1) Inherent or Implicit :- Based on data model itself. Eg. No lists as values

(2) Schema Based or Explicit :- Expressed using schema, provided by model.

(i) Domain Constraint :- Attribute value must be in domain specified

(ii) Key Constraint :- Key value must be unique for every row

(iii) Entity Integrity :- Key value can never be NULL

(iv) Referential Integrity :- Foreign Key must point to valid primary key.

(3) Semantic :- Cannot be expressed by model.

Insert, Update, Delete Violations

Insert can violate all 4 Schema constraints.

Delete can violate only referential. \Rightarrow It must be done with a certain extra

input of whether to RESTRICT, CASCADE or SET NULL to ensure

how delete is handled if referential integrity is broken.

Update = Insert + Delete.

Mapping of ER diagram to Relational

1) Mapping entities and weak entities to relations. Weak entities will have

a foreign key in them pointing back to their strong entity primary key

2) 1:1 relations are found. Assuming S has total participation and T isn't.

Primary key of T puts into S as foreign. (Can also merge S, T or create new R.)

3) 1:n relations are found. $S=1$ $T=n$. Put primary key (S) into T.

4) m:n relations are found. Create new relation R having primary key at both

5) Multi-valued attributes will also get a new relation with key same as primary key

6) N-ary relationship we create a new relation containing all primary keys.

Relational Algebra

Basic set of operations for relational model. It is closed. Result of operation is always a new relation.

Unary Relational Operators

(1) Select $\sigma_{\text{conditions}}$ (R)

This returns R' that has all the rows that match boolean condition specified.

(2) Project $\pi_{\text{attribute list}} (R)$

This returns R' with all rows but only specified attributes.

(3) Rename $P_{\text{New-name}}(\text{new-attributes names}) (R)$

This returns R' itself with new name and attributes.

Set theory operators

(1) Union $R \cup T$

All tuples either in R or T . At least $N_r + N_s$ max

(2) Intersection $R \cap T$

All tuples in both R, T . At least $N_r \cap N_s$ max

(3) Set Difference $R - T$

All tuples in R but not in T . At least $N_r - N_s$ max

(4) Cartesian Product $R \times T$

All possible combinations of R and T combining.

Subtangents elements of R and T will result in $N_r * N_s$ elements

Pre-requisite for Union, Intersection, Difference is that both R and T have to be type compatible \Rightarrow same no. of attributes, corresponding pair have same domain.

Their names in result will be same as R .

Binary Relational Operations

(1) Join $\bowtie R \bowtie_{\text{join conditions}} T$

Will return tuples that match the join condition.

$R.A_i > T.B_j$ or $<$, or $=$. If $=$ it means equi-join.

(2) Natural Join $* R * T$

Will perform equi-join on the basis of matching column names.

It will remove repeated Superfluous attributes.

(3) Outer Join $\bowtie R \bowtie T$

This will keep all the tuples. Matched $+ R$ with S as null $+ T$ with R as null.

Left \bowtie will keep only R . Right \bowtie will keep only T .

(4) Outer Union $\bowtie R \bowtie T$

This for union when not type compatible. All additional kept.

(5) Division $R \div T$ $R(A, B, \dots) \div T(A)$

This will return a tuple from R if it has all occurrences of A in T as part

of R . This answer can be found with the help of De Morgan's law.

Breaking the partitioned partition in a series of generalized process.

Query Tree Notation

Internal data structure to represent query. Estimate work involved, Optimizing it
Nodes = Operations Leaf nodes = Relations.

Additional Relational Operators

(1) Aggregate Functions $F_{\text{aggregate}} \rightarrow \text{attr} \dots (R)$
Sum, Count, Maximum, Minimum, Average etc.

(2) Grouping with Aggregation grouping attr $F_{\text{aggregate}} \rightarrow \text{attr}$
This will group results based on an attribute given. We need not apply any aggregate function and just group alone as well.

Not defined by relational algebra

(1) Recursive closure.

Recursive relationships that require iteration to find a hierarchy or something similar

Unit 3

We have to learn basic SQL commands here.

- Create Table table_name [small int | int | float | decimal] [CONSTRAINT constraint_name]

(Vari 1 DT Constraints,

: [small int | int | float | decimal]

Primary Key (Var i) [unique | primary key]

Foreign Key (Var j) references TableName (Var k)

) [constraint_name]

DTs = Numeric, Char, Bit, Boolean, Date, Timestamp, Interval [Increment/Decrement time]

Binary LOBs or Character LOBs [large page size 2^16 to 2^32]

* Internal LOBs also provide transactional support.

Constraints = Primary Key, Foreign Key [Specified separately]

Default, NOT NULL, Check (cond) [Specified with attribute]

We can name constraints and specify at end while creating table

- Drop command : This is used to drop named schema elements. We have to provide options of either cascade or restrict.

Alter command : This can change named schema elements

Alter Table Table_Name ADD/Drop Column/Constraint

Insert command : This is used to add values into tables.

(i) Insert into Table_name values (- - -) [inserts consider si prompt]

(ii) Insert into Table_name (- - -) QUERY [copy & paste from other query]

(iii) Create Table Table_name LIKE Table_1 (QUERY) WITH DATA

- Delete command : This is used to delete certain tuples

(i) Delete From Table_name Where Condition

- Update command : This is used to modify attribute values of selected tuple

(i) Update Table_name Set Var = Value Where Condition

- Three-Valued Logic

The two values used here are True, False, Unknown/Null

AND	T	F	N	OR	T	F	N	NOT
T	T	F	N	F	T	T	F	F
F	F	F	F	T	F	N	T	T
N	N	F	N	F	N	N	N	Unknown

T = True, F = False, N = Unknown/Null

AND operation result is True if all are True else False

OR operation result is True if any one is True else False

NOT operation result is True if value is False else True

True and False values are same so AND operation result is same

so AND operation result is same

SQL Retrieval Queries

Select <attribute> * = all Distinct = Unique values for attribute Aggregate variables

From <tables> as Alias / Another Name

[Where <conditions>]

[Nested | Sub or Correlated Query]

[Group By <grouping attributes>] Partition result into subsets. must be in select NULL is also

[Having <group condition>] - condition on group by

[Order By <attribute lists> DESC / ASC]

(Ans thorugh all the rows of both tables) Cartesian Product = 20 * 20 = 400

(i) If Where is not specified then Cartesian Product

(ii) Comparison Operators = Mostly usable in only Where clause

a) Like = Pattern matching for strings "RegEx" present in string

b) Between _ AND _ for variables

c) In used to check if value in set of values

d) Any / Some when we want other conditions other than equal. >, >=, <, <=

e) All (Query) when we want to compare with all rather than 1

(iii) Nested Query => Complete Query within another Query.

This could be uncorrelated or correlated. If uncorrelated it runs only once throughout execution of outer query.

Correlated Nested Query => Whenever a condition in the where clause of a nested query, it references attribute of a relation declared in outer query.

The correlated query is executed once for every tuple in outer query.

(iv) Functions in Where clause

a) Exist (Query) True = if 1 record at least False = if no records

b) Unique (Query) True = if all records unique False = if not all unique

(v) JOIN conditions in From Clause

We can say

From Table1 AnyType Of Join Table2 ON Condition

* Natural join has no condition necessary.

(vi) Constructs

a) With name_for_temp_table AS (Query | Insert | Update | Delete)

These create an auxiliary table for querying in larger queries.

They form Common Table Expressions.

The name_for_temp_table can be same as existing table and will be given higher precedence.

b) Case

When Condition Then Action

Else Action;

- **Assertions** = Used to specify additional types of constraints which is outside scope of built-in relational model constraints. This specifies a query that selects any tuples that violate desired condition. If used is to for Create Assertion Name Check (Condition using Query)

Condition will always be enforced and any violation raises an error

- **Triggers** = Specify automatic actions that will be performed when certain events or conditions occur.

Create Trigger T-name [constraint] [Information about trigger] [when] [event]

Before | After | Insert | Update | Delete | On (Attributes On Table)

For Each Row

When Condition must stated

Action

Event → Condition → Action

This gives use an object called **New** which can be used to check changes.

- **Views** = Single Table derived from other tables? It does not physically exist

Create View View-Name (Attributes) As Query

We can implement either query modification [Change query on view to tables]

which is very inefficient or view materialization where we create a temporary

view table when it is first queried. We would need a good strategy to

update this using immediate, lazy or periodic updates.

These can be used as an authorization mechanism.

- **Functions**

Create Function fn-name (param-list)

returns return-type

Must have return type

language plpgsql

Can't perform transactions.

as

Don't have transactional support.

\$\$ declare --variables \$\$

\$\$ begin --logic; end; \$\$

• Stored Procedures

(Create Procedure procedure_name([param_list]) [AS] language plpgsql) No compulsion to return type

Can make use of transactional support AS

\$\$ declare -- variable declaration \$\$

obliges \$\$ begin -- logic & end; \$\$

Cursors = Allows us to encapsulate a query and process each individual row at a time. To divide large result set to parts.

Declare → Open → Fetch row by row → Process → Close, next cursor

These are mainly used in functions and stored procedures.

Declare cursor_name [re]cursor; OR cursor_name cursor[attributes] for query;

Open cursor_name for query; OR Open cursor_name [(name := value, ...)]

Using Fetch [direction {from | in | ?}] cursor-name into target_variable
next, last, prior, first, Absolute_index, Relative_count.

Move [direction {from | in | ?}] cursor-name

Modifying Update table_name [set column = value]; Delete from table_name

DB Set column = value Where current of cursor_name;
Where current of cursor_name;

Closing Close cursor_name;

Online Analytical Processing

Used for processing historical data to get current results

processes to show the latest information available to facilitate user decisions

at fastest time-to-time basis all previous and all future data is required

Process of getting data to push information back into analysis

method of classification of data based on certain conditions

data is grouped into categories based on certain conditions

process of extracting useful information from data

mining data to find patterns, trends, anomalies, etc.

subset of data types etc.

Popularity, popularity

discovery of hidden patterns that can be used to predict future trends

process of extracting useful information from data

AB algorithm = mining etc

data mining = analysis of data

Unit 4

Design Guidelines for RDs

1) Each tuple should represent one entity. Attributes of different entities should not be mixed together in same relation.

2) Design a schema with no insert, delete, update anomalies.

3) Update anomaly = Data inconsistency from data redundancy and partial update

(a) Insert anomaly = Inability to insert due to absence of data.

(c) Delete anomaly = Unintended loss of data during delete.

3) Tuples should have minimal null values.

4) Relations should satisfy lossless join condition.

to avoid it no spurious tuples should be generated by natural join of all.

Functional dependencies = specify formal measures of goodness of relational designs. These are constraints that are derived from the meaning and inter-relationships of data attributes. $X \rightarrow Y$.

We may conclude that a FD may exist but not confirm. However, we can effectively say that certain FDs don't exist at all.

For $X \rightarrow Y$ if $t_1[X] = t_2[X]$ then $t_1[Y] = t_2[Y]$ is a must.

Given a set of FDs F , we can infer additional FDs that hold.

Armstrong Rules: 1) Reflexive = If Y is subset of X , $X \rightarrow Y$ always holds

Sound and Complete

2) Augmentation = If $X \rightarrow Y$, then $XZ \rightarrow YZ$

Can derive all others

3) Transitive = If $X \rightarrow Y$, $Y \rightarrow Z$, $X \rightarrow Z$ holds.

Decomposition $X \rightarrow YZ = X \rightarrow Y, X \rightarrow Z$. Union $X \rightarrow Y, X \rightarrow Z = X \rightarrow YZ$

Pseudotransitivity $X \rightarrow Y, Y \rightarrow Z$ then $X \rightarrow Z$

for $X \rightarrow Y = XU \rightarrow YV$ then $XU \rightarrow YU \cup YU \rightarrow Z = XN \rightarrow Z$

Closure of a set F of FD's is set F^+ of all FDs that can be inferred by F

or a set X of attr on FD's, F is X^+ of all attributes derived from X .

Equivalence of set of FDs = If every FD in F can be inferred from G and every FD in G can be inferred from F , then $F^+ = G^+$. F covers G , G covers F

Minimal covers of FDs = If every FD, RHS = 1 attribute [canonical form]

Can't replace $X \rightarrow A$ with $Y \rightarrow A$ where Y is subset of X

Can't remove any FD and have equivalent to F .

Algorithm: 1st) Get it to canonical form

2nd) For each $X \rightarrow A$, B in X , if $F - \{X \rightarrow A\} \cup (X - B) \rightarrow A = F$ remove B

3rd) For each $X \rightarrow A$, if $F - \{X \rightarrow A\} = F$ remove $X \rightarrow A$.

4th Normal Form = Multivalued Dependency

A multi-valued dependency $X \rightarrow\! Y$ exists if there is a set $Z = R - (X \cup Y)$

$$t_1[X] = t_2[X] = t_3[X] = t_4[X] \text{ and } Z = R - (X \cup Y)$$

then $t_1[Y] = t_3[Y]$, $t_2[Y] = t_4[Y]$ and Trivial if Z is subset of X or $X \cup Y = R$

$$t_1[Z] = t_4[Z], t_2[Z] = t_3[Z] \times \text{possibly sets}$$

For all non-trivial multivalued dependency in 4NF $\{X \rightarrow\! Y\}$

$X \rightarrow\! Y$, X is superkey then in 4NF $\{X \rightarrow\! Y\} = 0$

5th Normal Form = Project-Join (Normal) Form (NFIJ)

$$+ \forall \alpha ((\alpha - \beta) \leftarrow (\alpha \cap \beta))$$

and satisfies null pointer (9)

thus no valid and $\alpha \cap \beta = \emptyset$ no common part

and valid and $\alpha \cap \beta = \alpha$ $\alpha \cap \beta = \beta$

and valid and $\alpha \cap \beta = \alpha \cup \beta$ $\alpha \cap \beta = \emptyset$

Properties of Relational Decomposition

We perform this in an attempt to reduce redundancy of data.

a) Attribute Preservation Condition = Every attribute must appear once

b) Dependency-Preservation Property = $X \rightarrow\! Y$ holds $\Rightarrow X \rightarrow\! Y$ holds

A decomposition $D = \{R_1, \dots\}$ is dependency preserving if (9)

union of projections of F onto $R_i = F$ 4NF + union of sup. and sub.

$$[\pi_{R_1}(F) \cup \pi_{R_2}(F) \dots]^+ = F^+$$

c) Lossless Join property = no spurious tuples due to join \Rightarrow MUST PRESERVE

Algorithm to test for non-additive join property

1st) Initial matrix for each relation R_i having one row and one column j

for each attribute

2nd) Set all x to 1 first and in row 0 and 0 in col 0

3rd) For each row i in R_i and j in R_j if R_i has A_j

For each column j in R_j and i in R_i if R_j has A_i

if R_i has A_j then set all to 1

and if R_j has A_i make it a 1 if it is 0 = start again

4th) For each FD $X \rightarrow Y$ in F \rightarrow Repeat till whole loop had no change

All rows that have every part of X as 1

For all attributes of Y

if any row matched = 1
set all to 1

contd 4th) This means that every row which has X as l.a is checked to

see if any of the attributes of Y = l.a (and change fit in)

(XUY all rows if any such attributes exists $\Rightarrow [S] \sqsubseteq$)

5th) If any row has all l.a's. Property ✓ It is fit
else Property X $[S] \sqsubseteq [S] \sqsubseteq [S] \sqsubseteq [S]$

d) Testing Binary Decompositions

$D = \{R_1, R_2\}$ has lossless join if and only if $X_1 \nsubseteq X_2 \cup X_3$

The F.D $((R_1 \cap R_2) \rightarrow (R_1 - R_2))$ in F^+ or $((R_1 \cap R_2) \rightarrow (R_2 - R_1))$ in F^+

e) Successive Non-additive Join

If a decomposition $D = \{R_1, \dots, R_m\}$ has lossless join. and

$R_i = \{Q_1, Q_2, \dots, Q_n\}$ has loss less join

$\therefore D = \{R_1, \dots, Q_1, \dots, Q_n, \dots, R_m\}$ has lossless join

Synthesis in 3NF

Step 1) Find minimal cover G

Step 2) For each L.H.S. X of FD in G, find min. r.h.s. covering all

create relation schema D $\{ \{X \rightarrow A_1, X \rightarrow A_2, \dots\} \}$

where $X \rightarrow A_1, X \rightarrow A_2, \dots$ are only in G with X as LHS

Step 3) If none in D has key of R, add one more with key of R

We can get alternate 3NF by starting with same FDs. Due to minimal cover

Synthesis in BCNF

Given R , $F = FD$ and $D = R$ where $G = \text{minimal cover}$ of F

While there is relation Q in D not in BCNF do

do步 for each r.h.s. of Q if it is not a superkey of Q

Choose a Q in D not in BCNF, find FD $X \rightarrow Y$ violating BCNF

Replace Q in D by $Q - Y$ and $X \cup Y$ if $X \cup Y$ not in D

Problems with bad relations

1) Null Values

2) Dangling Tuple = Tuple that does not participate in natural join.

Symbol of last good tuple this happened

? on $X \rightarrow Y$ if X not r.h.s.

Example - After X joins Y this good tuple has H.A = ?

Then $X \rightarrow Y$ is violated by $X \rightarrow Z$

N = bad tuple for you Z

N is not good