

# Bases de Datos (BD) =

Purpose: (Problems to solve from filesystem)

- Data Redundancy and Inconsistency
- Difficulty of accessing data
- Data isolation (Multiple files and formats)
- Integrity problems (Hard to create, modify constraints)
- Atomicity of updates (Solved with transactions)
- Concurrent access by multiple users
- Security problems
- Low data independence

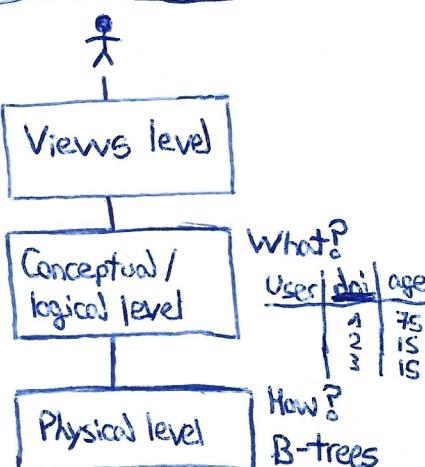
## Definition:

- Database : Set of interrelated data / data + metadata: Information about the data
- DataBase Management System: Is the software application that interacts with the user, applications and the database to capture and analyze data.
- It allows users to create, update, remove, read data in a database by an interface:
  - + Retrieve = Select
  - + Modify = - insert - delete - update

Data Models (Types): (They describe: Data, Relationships, Semantics, Constraints)

- \* Relational Model: Implementation → Defines how the data is stored (in tables (relations), with tuples (rows), attributes (columns), with keys, constraints, ...)
- \* Entity-Relationship data model: Design → Defines what the database should look like (Entities (real-world objects), attributes (properties) of the entities) and relations (connections between entities))
- \* Object-based data models
- \* Semi-structural data model (XML)
- \* Network model
- \* Hierarchical model

## Abstraction levels:



## SQL query language:

- Non procedural (Declarative)
- Takes several tables as input
- Returns always one as output
- Is NOT turing machine equivalent

## Transaction management:

- A transaction is a collection of operations that are performed as a block (cannot be divided). If during the execution a failure happens, all changes are undone (rollback) to maintain db consistency.
- Ensures that the db remains consistent despite system failures
- Properties ACID: Atomicity, Consistency, Isolation and Durability

## Schemas and Instances:

- Schemas are the structure = (Fixed)
  - Logical: What data is stored and how is related
  - Physical: How data is stored (files, tables, ...)
- Instances are the data at a given point (Changes)
- \* Independence of physical data: Means that we can change the physical schema (level) without affecting the logical one. And vice versa.

## Data Definition Language (DDL):

Defines the DB structure (schema) {CREATE, DROP, ALTER}

## Data Manipulation Language (DML):

Works with the data inside the tables (instances) {Insert, Select, Update, Delete}

- Types:
- Pure → For theoretical and research purposes (More procedural / imperative)
  - Commercial → Used in real-world db. systems (More declarative)

## DB Users:

- Native (Final Users): No knowledge about db.
- Application Programmers: Use python, C#, Java with db connectivity
- Sophisticated: Query the db directly without writing application programs
- Specialized: Work with non-traditional db like CAD.

## Data Base Administrator (DBA):

- Person who has control over the database.

Entities = (Instances) (A specific person, company, event)

↳ Have attributes: Value, domain

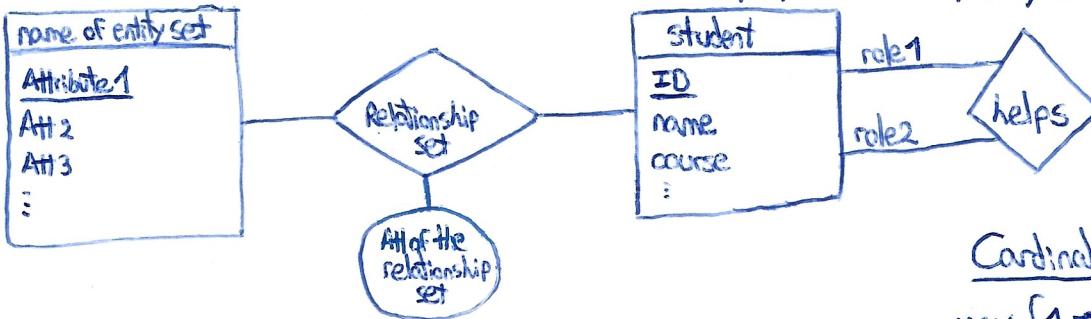
Entity Sets = (Classes) (The group of entities of the same type)

Relationship: Connection of 2 or more entities

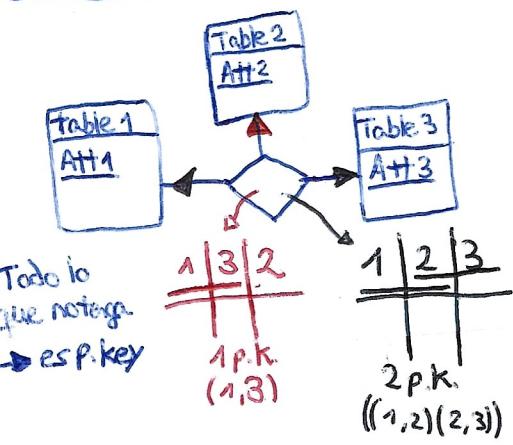
Sets: Collection of similar relations between entities.

## E-R Diagram:

If an attribute is underlined is a primary key: Identifies uniquely each member/entity of the set/Table.

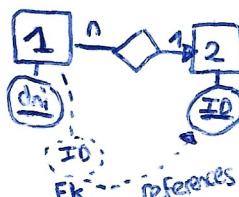


## Non-binary Rela. Sets:

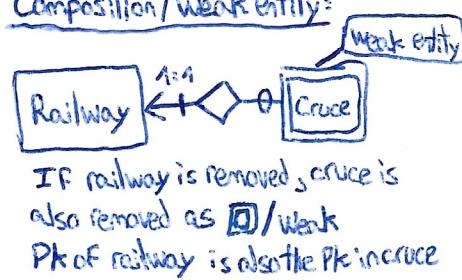


To do  
give notes.  
→ esp. p.key

## Foreign keys:



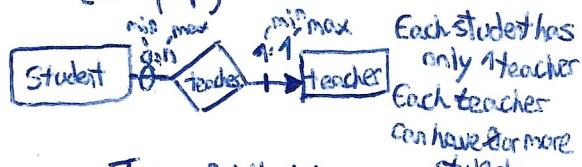
## Composition / Weak entity:



## Cardinality:

• MAX  $\{1 \rightarrow 1 : (\rightarrow) \text{ or } (1)\}$   
 $\{+1 \rightarrow n : (-) \text{ or } (n)\}$

• MIN  $\{\emptyset : (+)\}$   
 $\{1 : (+)\}$



## Types of Attributes:

xxxx Descriptive

xxxx Multivalue  
(is a table)

xxxx Identificative

## Generalization / Specialization:



The father's abstract

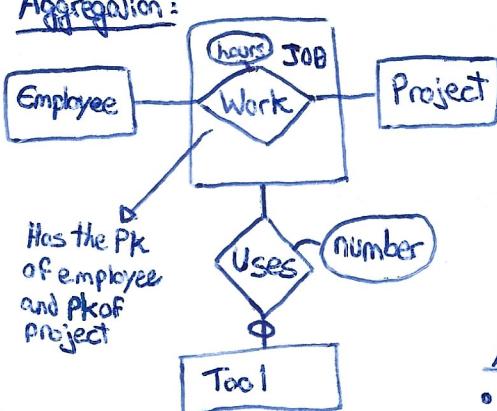


The father is not abstract

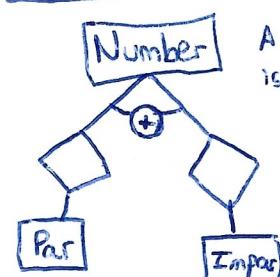
IS A → Cannot be both types

IS As → Can be both types of children

## Aggregation:



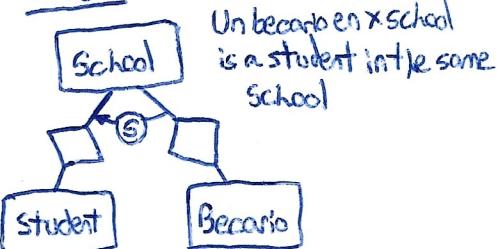
## Exclusion:



## Life cycle of a Data Base:

- Requirement Analysis (of the universe of discourse)
- Conceptual Design (E-R model / view integration)
- Logical Design (E-R → Relational model)
- Refinement Through use (optimization / denormalization)
- Physical Design (creation of indexes)
- Implementation (constraints, security and applications)
- Testing, monitoring, and maintenance

## Subset:



They are in account so removed



## Bases de datos (BD) =

### Relational Algebra:

Selection  $\sigma_p(r) \rightarrow$  where

Projection  $\pi_{a_1, a_2}(r) \rightarrow$  Select

Cross Join  $r \times s$

Natural Join  $r \bowtie s$  (Derived from cross Join)

Rename  $P_x(r)$

Union  $r \cup s$

Difference  $r - s$

Intersection  $r \cap s = r - (r - s)$

### ARMSTRONG RULES:

Reflexivity  $A \rightarrow A \quad B \subseteq A$

Augmentation  $A \rightarrow B \quad wA \rightarrow wB$

Transitivity  $A \rightarrow B, B \rightarrow C \quad A \rightarrow C$

Decomposition  $A \rightarrow BC \quad A \rightarrow B, A \rightarrow C$

Union  $A \rightarrow B, A \rightarrow C \quad A \rightarrow BC$

Pseudotransitivity  $A \rightarrow B, Bw \rightarrow C \quad Aw \rightarrow C$

### Definition of $V = W \quad r(R)$ :

$\forall t_1, t_2 \in r \quad (t_1[V] = t_2[V] \rightarrow t_1[W] = t_2[W])$

Definition Superkey: Candidate key + any other attribute (or not)

Any  $X$  where  $X^+ = R$  (all possible attributes) (The last superkey is the set of all attributes)  
Definition Candidate Key: candidate key

A minimal superkey where no attributes can be removed without losing  $X^+ = R$   
(there can be several candidate keys; one may be chosen as primary key)

### Functional Dependencies (constraints):

$$F = \left\{ \begin{array}{l} A \rightarrow ABC \\ \vdots \\ DE \rightarrow FG \end{array} \right. = \text{Original set of functional dependencies}$$

$F^+$  is an equivalent closure

$$F_c = \left\{ \begin{array}{l} PC \rightarrow \text{city} \\ \text{city} \rightarrow \text{province} \\ \text{DNI} \rightarrow \text{name}, PC \end{array} \right. = \text{Minimal cover closure (without redundancy)}$$

$X^+$  are the reachable attributes from  $X$ , and from the reachable of those, ...

### Closure $F^+$ of $F$ given $X$ :

$R = (A, B, C, \dots) \rightarrow$  All possible attributes

$X = (A, B) \rightarrow$  Set of  $R$  (attributes that I have)

$$F = \left\{ \begin{array}{l} AB \rightarrow C \\ C \rightarrow DE \end{array} \right. \quad (A, B)^+ = \{A, B, C, D, E, \dots\}$$

IF  $X^+ = R \rightarrow X \in F^+$   
And  $X$  is a superkey

### Steps for 3NF:

0.  $F$  should be  $F_c$

1. Each dependency  $X \rightarrow A$  is a table:  $R_i = (X, A)$  if  $X \rightarrow A, X \neq B$  then  $R_i = (X, A, B)$

2. Remove right attributes by transitivity  $A \rightarrow BC \rightarrow A \rightarrow B$   
 $B \rightarrow C \rightarrow B \rightarrow C$

3. Remove left attributes by pseudotransitivity  $A \rightarrow B \rightarrow A \rightarrow B$   
 $AB \rightarrow C \rightarrow A \rightarrow C$

### Boyce-Codd Normal Form:

For each  $R_i$  either:

- $X \rightarrow Y$  is trivial
- $X$  is a superkey

### 3NF:

For each  $R_i$  either:

- $R_i$  is in BCNF
- All  $Y$  is prime (part of a candidate key)

Is more important to Preserve Dependencies than Remove Redundancy

## Decomposition:

$$\Gamma(R) \rightarrow \Gamma_1(R_1), \Gamma_2(R_2), \dots, \Gamma_n(R_n)$$

Types:

• **Lossy** (very bad): Phantom rows appear  $r < \Gamma_1 \bowtie \Gamma_2 \bowtie \Gamma_3 \bowtie \dots \bowtie \Gamma_n$

• **Lossless** (don't lose info):  $\Gamma = \Gamma_1 \bowtie \Gamma_2 \bowtie \dots \bowtie \Gamma_n$   $R = R_1 \cup R_2 \cup R_3 \dots \cup R_n$

## Dependency Preservation: (DP)

$$F \leftrightarrow F_1 \cup F_2 \cup F_3 \cup \dots \cup F_n$$

Example of relational algebra:

Select DNI

$$\begin{array}{l} \text{From Person} \\ \text{Where DNI} > 10 \\ \text{OR name} \neq \text{'Paco'} \end{array} \rightarrow \prod_{(DNI)} (\delta_{(DNI > 10)} (\text{Person}) \cup \delta_{(\text{name} \neq \text{'Paco'})} (\text{Person})) = \prod_{(DNI)} (\delta_{(DNI > 10 \text{ } || \text{ } \text{name} \neq \text{'Paco'})} (\text{Person}))$$

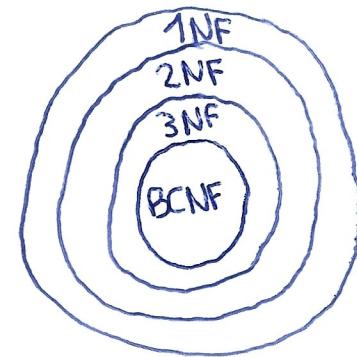
Select e.Name AS EmpName, d.DeptName AS Dept

FROM Employee AS e INNER JOIN Department AS d ON e.DeptID = d.DeptID

WHERE e.Salary > 50000;

$$P_{F(\text{EmpName}, \text{Dept})} (\prod_{(e.\text{Name}, d.\text{DeptName})} (\delta_{(e.\text{Salary} > 50000)} (P_e(\text{Employee}) \bowtie_{e.\text{DeptID} = d.\text{DeptID}} P_d(\text{Department})))$$

Types of Normal Forms:



1NF: Atomic (Indivisible) Values for every attribute (No lists/sets)

Primary keys in each table

No repeated columns with Phone, Phone1, ...

Same domain and type for each row

2NF: 1NF + No Nonprime attribute is partially dependent on any candidate key.

3NF: 2NF + Every non-trivial FD ( $x \rightarrow A$ ) either  $\{x\}$  is a superkey of R

BCNF:  $\{x\} \rightarrow X$  or  $X \rightarrow Y$  where X is a superkey of R

Normalization Steps:

0. Start with a canonical cover  $F_c$

1. Compute candidate keys for R.

2. For each functional dependency in  $F_c$ :  $(X \rightarrow A)$

- Generate an  $R_i$

- Compute the  $F_i$  constraint of  $F$  for  $R_i$

- Compute candidate keys for  $R_i$

3. If a dependency is repeated in several  $R_i$ , leave 1.

4. If no candidate key of R appears in a  $R_i$ , create an additional  $R_i$  with  $F_i = \{\}$

5. Mark the normal form (3NF or BCNF) of every  $R_i$

6. Check lossless join (with  $\bowtie$ )

7. Check dependency preserved ( $F \rightarrow F_1 \cup F_2 \cup \dots \cup F_n$ )