

Database Normalization

[Introduction to Normalization Theory](#)

[Database Anomalies](#)

[Insertion Anomaly](#)

[Update anomaly](#)

[Delete anomaly](#)

[Decomposition and spurious tuples](#)

[Functional Dependencies](#)

[Examples](#)

[Normalization](#)

[Prime attribute, full functional dependency and transitive function dependency](#)

[Normal Forms](#)

[Examples](#)

[Example 1 \(2NF\)](#)

[Example 2 \(3NF\)](#)

[Example 3 \(2NF, 3NF\)](#)

[References](#)

Introduction to Normalization Theory

| How to design a “good” DATABASE?

Informal definition:

- The schema should represent a “distinct” entity
- Little or no redundancy.
- Lesser no. of null values
- No modification anomaly
- No spurious tuple

| **Normalization theory helps us in designing a good database in a formal manner.**

Database Anomalies

Consider the following schema: (empid, empname, projid, projname)

Insertion Anomaly

- Inserting a new employee requires us to assign a project and vice versa.

Update anomaly

- Updating the name of a project requires us to update all the employees who are assigned to that project.

Delete anomaly

- Deleting a project might require us to delete all its employees.

Decomposition and spurious tuples

- <https://www.geeksforgeeks.org/spurious-tuples-in-dbms/>
- <https://www.geeksforgeeks.org/difference-between-lossless-and-lossy-join-decomposition/>

Make sure that we preserve the losslessness of the corresponding join while decomposing a table

Functional Dependencies

- **Functional dependencies (FDs)** are *constraints* derived from the meaning of and relationships among attributes
- A set of attributes X **functionally determines** Y, denoted by $X \rightarrow Y$, if the value of X determines a *unique* value of Y
- <https://www.geeksforgeeks.org/introduction-of-database-normalization/>

Try to preserve the FD while splitting

Examples

- `employees(empid , name)` with FD:
 - `(empid) -> name`
- `users(uuid , email , name)` with FDs:
 - `(uuid, email) -> (email)` , **trivial FD**
 - `(uuid, email) -> (email, name)` , **semi non-trivial FD**
 - `(uuid, email) -> (name)` , **completely non-trivial FD**

Normalization

- The process of decomposing relations into smaller relations that conform to certain norms is called **normalization**
- Keys and FDs of a relation determine which normal form a relation is in.
- Different normal forms
 - **1NF**: based on attributes only
 - **2NF, 3NF, BCNF**: based on keys and FDs
 - **4NF**: based on keys and multi-valued dependencies (MVDs)
 - **5NF or PJNF**: based on keys and join dependencies
 - **DKNF**: based on all constraints

For now, we are just concerned about 1NF, 2NF, and 3NF

Prime attribute, full functional dependency and transitive function dependency

- A **prime attribute** must be a member of some candidate key
 - Example: `roll`
- A **non-prime attribute** is not a member of any candidate key
 - Example: `gender`
- An FD `X → Y` is a **full functional dependency** if the FD does not hold when any attribute from X is removed

- Example: $(roll) \rightarrow (name)$
- It is a **partial functional dependency** otherwise
 - $(roll, gender) \rightarrow (name)$
- A FD $X \rightarrow Y$ is a **transitive functional dependency** if it can be derived from two FDs $X \rightarrow Z$ and $Z \rightarrow Y$
 - Example: $(roll) \rightarrow (hod)$ since $(roll) \rightarrow (deptid)$ and $(deptid) \rightarrow (hod)$ hold
- It is **non-transitive** otherwise
 - Example: $(roll) \rightarrow (name)$

Normal Forms

Informal definition:

- **1NF:** All attributes depend on the key
- **2NF:** All attributes depend on the whole key
- **3NF:** All attributes depend on nothing but the key

Tests for normal forms:

- **1NF:** The relation should have no multivalued attributes or nested relations
- **2NF:** For a relation where the candidate key contains multiple attributes, no non-key attribute should be functionally dependent on a part of the candidate key
- **3NF:** The relation should not have a non-key attribute functionally determined by a set of non-key attributes

Remedies for normal forms:

- **1NF:** Form new relations for each multi-valued attribute or nested relation
- **2NF:** Decompose and set up a relation for each partial key with its dependent(s); retain the primary key and attributes fully dependent on it
- **3NF:** Decompose and set up a relation for each non-key attribute with non-key attributes functionally dependent on it

Examples

Example 1 (2NF)

- Consider $(\underline{\text{Id}}, \underline{\text{ProjId}}, \text{Hrs}, \text{Name}, \text{ProjName})$ with FDs:
 - $(\text{Id}, \text{ProjId}) \rightarrow (\text{Hrs})$
 - $(\text{Id}) \rightarrow (\text{Name})$
 - $(\text{ProjId}) \rightarrow (\text{ProjName})$
- It is not in **2NF** since (Name) depends partially on $(\text{Id}, \text{ProjId})$
- After **2NF** normalization,
 - $(\text{Id}, \text{ProjId}, \text{Hrs})$ with FD: $(\text{Id}, \text{ProjId}) \rightarrow (\text{Hrs})$
 - (Id, Name) with FD: $(\text{Id}) \rightarrow (\text{Name})$
 - $(\text{ProjId}, \text{ProjName})$ with FD: $(\text{ProjId}) \rightarrow (\text{ProjName})$

Example 2 (3NF)

- Consider $(\underline{\text{empid}}, \text{empname}, \text{projid}, \text{projname})$ with FDs:
 - $(\text{empid}) \rightarrow (\text{empname}, \text{projid})$
 - $(\text{projid}) \rightarrow (\text{projname})$
- it is not in **3NF** since (projname) depends transitively on (Id) through (ProjId) .
- After **3NF** normalization,
 - $(\underline{\text{empid}}, \text{empname}, \text{projid})$ with FD: $(\text{empid}) \rightarrow (\text{empname}, \text{projid})$
 - $(\underline{\text{projid}}, \text{projname})$ with FD: $(\text{projid}) \rightarrow (\text{projname})$

Example 3 (2NF, 3NF)

- $L = (\text{Id}, \text{Dist}, \text{Lot}, \text{Area}, \text{Price}, \text{Rate})$ with FDs:
 - $(\text{Id}) \rightarrow (\text{Dist}, \text{Lot}, \text{Area}, \text{Price}, \text{Rate})$
 - $(\text{Dist}, \text{Lot}) \rightarrow (\text{Id}, \text{Area}, \text{Price}, \text{Rate})$
 - $(\text{Dist}) \rightarrow (\text{Rate})$
 - $(\text{Area}) \rightarrow (\text{Price})$
- L is not in **2NF** because (Rate) depends partially on (Dist)
- $L_1 = (\text{Id}, \text{Dist}, \text{Lot}, \text{Area}, \text{Price})$ with FDs:
 - $(\text{Id}) \rightarrow (\text{Dist}, \text{Lot}, \text{Area}, \text{Price})$

- $(\text{Dist, Lot}) \rightarrow (\text{Id, Area, Price})$
- $(\text{Area}) \rightarrow (\text{Price})$
- $\text{L2} = (\text{Dist, Rate})$ with FD:
 - $(\text{Dist}) \rightarrow (\text{Rate})$
- L1 is in **2NF** but not **3NF** because (Price) depends on (Id) through (Area)
- L2 is in **2NF** and in **3NF**
- $\text{L11} = (\text{Id, Dist, Lot, Area})$ with FDs:
 - $(\text{Id}) \rightarrow (\text{Dist, Lot, Area})$
 - $(\text{Dist, Lot}) \rightarrow (\text{Id, Area})$
- $\text{L12} = (\text{Area, Price})$ with FD:
 - $(\text{Area}) \rightarrow (\text{Price})$
- L11 and L12 are in **3NF**

References

- <https://www.geeksforgeeks.org/normal-forms-in-dbms/>
- <https://docs.microsoft.com/en-us/office/troubleshoot/access/database-normalization-description>