

Data normalisation and ACID properties

Agenda

- Data Anomalies in DBMS
- Data Normalisation
 - 1NF
 - 2NF
 - 3NF
 - Boyce-Codd Normal Form (BCNF)
- Transactions
 - ACID properties

Key Terms

Functional Dependencies

Functional Dependency is when one attribute determines another attribute in a DBMS

Data Normalisation

the process of splitting relations into well-structured relations that allow users to insert, delete, and update tuples without introducing database inconsistencies

Transactions

A set of logically related operation. A transaction is an action or series of actions. It is performed by a single user to perform operations for accessing the contents of the database.

ACID properties

ACID (atomicity, consistency, isolation, durability) is a set of properties of database transactions intended to guarantee data validity despite errors, power failures, and other mishaps

Data Anomalies

An anomaly is something that is unusual or unexpected; an abnormality

A very common source for issues in the database is redundancy. Apart from storage issues, having the same value present in multiple rows can lead to inconsistent data. Have a look at the following **STUDENTS** table.

| <u>ID</u> | NAME | EMAIL | BATCH_ID | BATCH_NAME |
|-----------|------|-------|----------|------------|
|-----------|------|-------|----------|------------|

Let us assume that the above table is the only table in the database. The student and batch entities are tightly coupled. What are some issues that can be caused by this?

- How do we create a new batch without students?

- How do we create a new student without a batch?
- What data do we lose if we delete a student?
- What happens if we modify a batch name but miss a record?

Anomalies are avoided by the process of normalisation. The above issues or anomalies can be categorised into the following categories:

Insertion Anomalies

The inability to insert a new tuple into a table due to missing data is known as insertion anomaly.

An insertion anomaly occurs when data cannot be inserted into a database due to other missing data.

This is most common for fields where a foreign key must not be NULL, but lacks the appropriate data. Adding a student without a batch is not possible in the above schema as a batch_id is required. Whereas creating a new batch without students would require multiple null values and special handling.

Deletion Anomalies

A deletion anomaly occurs when data is unintentionally lost due to the deletion of other data. A deletion anomaly is the unintended loss of data due to deletion of other data.

| <u>ID</u> | NAME | EMAIL | BATCH_ID | BATCH_NAME |
|-----------|-------------|---------------|----------|-------------------|
| 1 | John Watson | j@sherlock.ed | 1 | Sherlock Season 6 |
| 2 | Mary Watson | m@sherlock.ed | 1 | Sherlock Season 6 |
| 3 | Kilvish | kil@vi.sh | 2 | Shaktimaan |

In the above table, just **Kilvish** is associated with the batch **Shaktimaan**. Now, if we delete **Kilvish** from the database, we lose the data associated with the batch. This results in database inconsistencies and is an example of how combining information that does not really belong together into one table can cause problems.

Update Anomalies

An update anomaly occurs when data is only partially updated in a database. An update anomaly is a data inconsistency that results from data redundancy and a partial update.

| <u>ID</u> | NAME | EMAIL | BATCH_ID | BATCH_NAME |
|-----------|----------------|---------------------|----------|-------------------|
| 1 | John Watson | j@sherlock.ed | 1 | Sherlock Season 6 |
| 2 | Mary Watson | m@sherlock.ed | 1 | Sherlock Season 6 |
| 3 | Kilvish | kil@vi.sh | 2 | Shaktimaan |
| 4 | Mycroft Holmes | brother@sherlock.ed | 1 | Sherlock Season 6 |

In the above table, we have three students associated with the batch **Sherlock Season 6**. If we have to update the batch name, each row will have to be updated due to redundancy. This adds an overhead and a

likely source of data inconsistency. If our developer or query misses a record, the database will be in an inconsistent state.

Functional Dependencies

a dependency FD: $X \rightarrow Y$ means that the values of Y are determined by the values of X. Two tuples sharing the same values of X will necessarily have the same values of Y.

- A functional dependency is a constraint that specifies the relationship between two sets of attributes where one set can accurately determine the value of other sets.
- It is denoted as $X \rightarrow Y$, where X is a set of attributes that is capable of determining the value of Y
- The attribute set on the left side of the arrow, X is called Determinant, while on the right side, Y is called the Dependent

Suppose one is designing a system to track vehicles and the capacity of their engines. Each vehicle has a unique vehicle identification number (VIN). One would write **VIN \rightarrow EngineCapacity** because it would be inappropriate for a vehicle's engine to have more than one capacity. On the other hand, EngineCapacity \rightarrow VIN is incorrect because there could be many vehicles with the same engine capacity

ID NAME EMAIL BATCH_ID BATCH_NAME

Using the above schema, the following observations can be made:

- **ID** can be used to derive **NAME** and **EMAIL**. Hence,
 - **ID \rightarrow NAME**
 - **ID \rightarrow EMAIL**
- **BATCH_ID** can be used to derive **BATCH_NAME**
 - **BATCH_ID \rightarrow BATCH_NAME**
- Since an **email** is a unique identifier, it can be used to derive **ID** and **NAME**
 - **EMAIL \rightarrow ID**
 - **EMAIL \rightarrow NAME**
- **ID** can also be used to derive **BATCH_ID**
 - **ID \rightarrow BATCH_ID**

Figure out the functional dependencies for the below table

| MENTOR_ID | STUDENT_ID | SESSION_ID | RATING | FEEDBACK |
|-----------|------------|------------|--------|-----------|
| 1 | 1 | 1 | 5 | Very Good |
| 1 | 2 | 1 | 4 | Good |
| 2 | 3 | 1 | 3 | Average |
| 1 | 1 | 2 | 4 | Good |

Data Normalisation

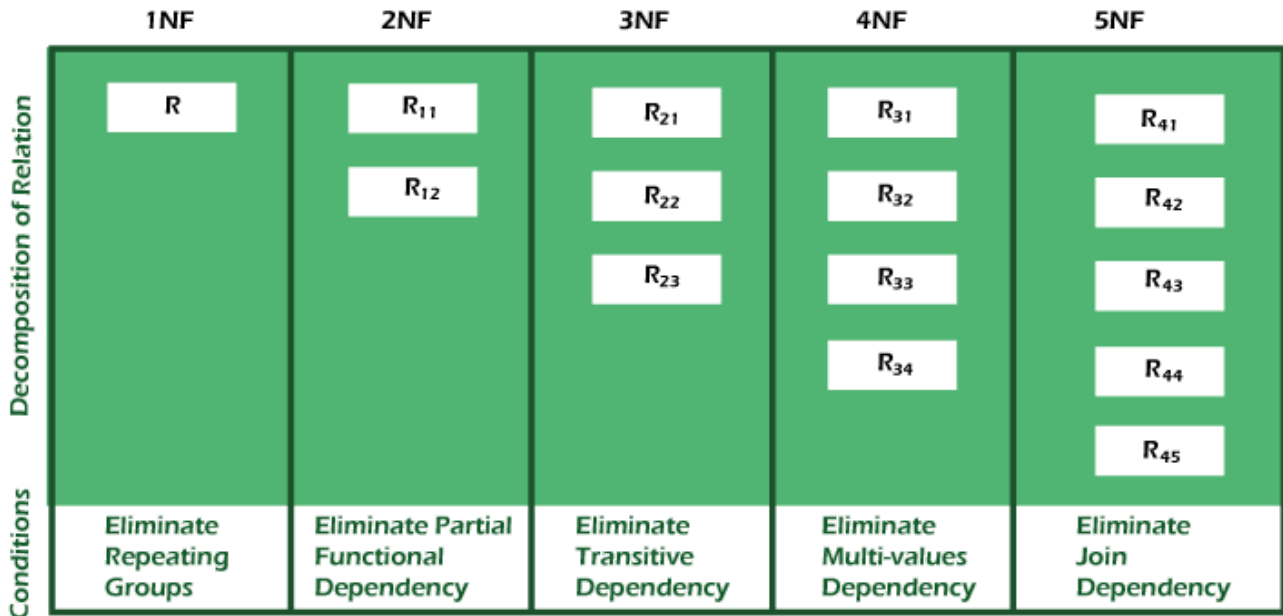
the process of structuring a relational database in accordance with a series of so-called normal forms in order to reduce data redundancy and improve data integrity. Normalization entails

organizing the columns (attributes) and tables (relations) of a database to ensure that their dependencies are properly enforced by database integrity constraints

The goal of normalisation is to produce a set of tables that

- Is a faithful model of the enterprise
- Is highly flexible
- Reduces redundancy-saves space and reduces inconsistency in data
- Is free of update, insertion and deletion anomalies

Following are the various normal forms:



1NF

- A relation will be 1NF if it contains an **atomic** value.
- It states that an attribute of a table cannot hold multiple values. It must hold only **single-valued attribute**.
- First normal form disallows the multi-valued attribute, composite attribute, and their combinations.

| ID | NAME | EMAIL | PHONE_NUMBERS |
|----|-------------|------------------|------------------------|
| 1 | Tantia Tope | tantia@rani.bai | [123456789, 987654321] |
| 2 | Kilvish | kil@vi.sh | [987654321, 123456789] |
| 3 | John Watson | i.am@sherlock.ed | [123456789, 987654321] |

The above table is not 1NF because it contains a multi-valued attribute i.e. phone numbers.

Redundant columns

| NAME | EMAIL | PHONE_NUMBER_01 | PHONE_NUMBER_02 |
|-------------|-----------------|-----------------|-----------------|
| Tantia Tope | tantia@rani.bai | 123456789 | 987654321 |

| NAME | EMAIL | PHONE_NUMBER_01 | PHONE_NUMBER_02 |
|-------------|------------------|-----------------|-----------------|
| Kilvish | kil@vi.sh | 987654321 | 123456789 |
| John Watson | i.am@sherlock.ed | 123456789 | 987654321 |

Cons -

- Wasteful if all rows have mostly one phone number
- Hard to determine upper bound of number of phone numbers
- Querying is not efficient since multiple columns needs to be queried
- Multiple indexes

Redundant rows

| ID | NAME | EMAIL | PHONE_NUMBERS |
|----|-------------|------------------|---------------|
| 1 | Tantia Tope | tantia@rani.bai | 123456789 |
| 1 | Tantia Tope | tantia@rani.bai | 987654321 |
| 2 | Kilvish | kil@vi.sh | 123456789 |
| 2 | Kilvish | kil@vi.sh | 987654321 |
| 3 | John Watson | i.am@sherlock.ed | 987654321 |
| 3 | John Watson | i.am@sherlock.ed | 123456789 |

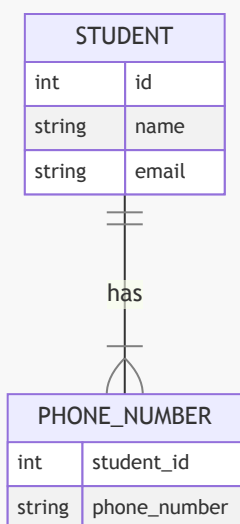
What will be primary key in the above table?

Cons -

- A lot of redundant rows which can lead to anomalies
- Primary key needs to be altered

Separate Mapping Table

The two solutions above are not ideal due to the large amount of redundant data. In order to properly convert the above table to 1NF, we need to create a separate table that maps the redundant data to the primary key. Hence, a **PHONE_NUMBER** table is created with **student_id** and **phone_number** columns. There are multiple rows for each student and the ID is used to map the redundant data. This minimises the amount of redundant data.



2NF

- In the 2NF, relational must be in 1NF.
- There should be no partial dependencies.
- Every non candidate-key attribute must depend on the whole candidate key, not just part of it

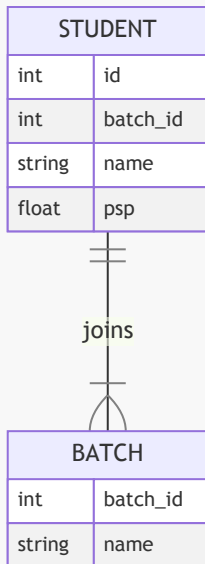
ID NAME BATCH_ID BATCH_NAME PSP

Listing out the dependencies for the above table:

- **ID** and **BATCH_ID** can be used to derive **PSP**
 - **ID, BATCH_ID** → **PSP**
- **BATCH_ID** can be used to derive **BATCH_NAME**
 - **BATCH_ID** → **BATCH_NAME**

In the first dependency, **ID** and **BATCH_ID** can determine a non-candidate key attribute. However, in the second dependency just **BATCH_ID** can determine a non-candidate key. This is an example of a partial dependency and hence violates 2NF.

The above table can be normalised by creating a separate table for batch information.



3NF

- In the 3NF, relational must be in 2NF.
- It should also not contain any transitive dependencies.

ID **NAME** **BATCH_ID** **BATCH_NAME**

Listing out the dependencies for the above table:

- **ID** → **NAME**
- **ID** → **BATCH_ID**
- **BATCH_ID** → **BATCH_NAME**
- **ID** → **BATCH_NAME**

It can be observed in the last three dependencies that **ID** can determine **BATCH_ID** that can be used to determine **BATCH_NAME** i.e. **ID** → **BATCH_ID** → **BATCH_NAME**. This is an example of a **transitive dependency** and hence violates 3NF.

A relation is in third normal form if it holds atleast one of the following conditions for every non-trivial function dependency $X \rightarrow Y$.

- X is a super key.
- Y is a prime attribute, i.e., each element of Y is part of some candidate key.

ID **NAME** **PHONE** **BATCH_ID** **BATCH_NAME**

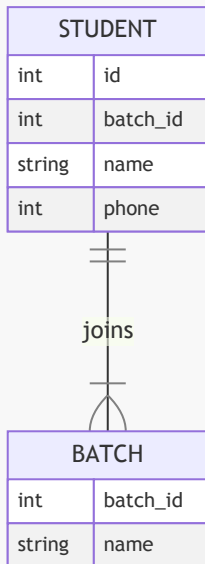
In the above table,

- **ID** → **PHONE**
- **PHONE** → **ID**
- **PHONE** → **NAME**

Do any of the above violate 3NF? **NO** Since **ID** and **PHONE** are both candidate keys. What about **BATCH_ID** → **BATCH_NAME**?

Yes, this violates 3NF since **BATCH_ID** is not a super key and **BATCH_NAME** is not a prime attribute.

Again, the above table can be normalised by creating a separate table for batch information.



Boyce-Codd Normal Form (BCNF)

- A table is in BCNF if every functional dependency $X \rightarrow Y$, **X is the primary key of the table.**

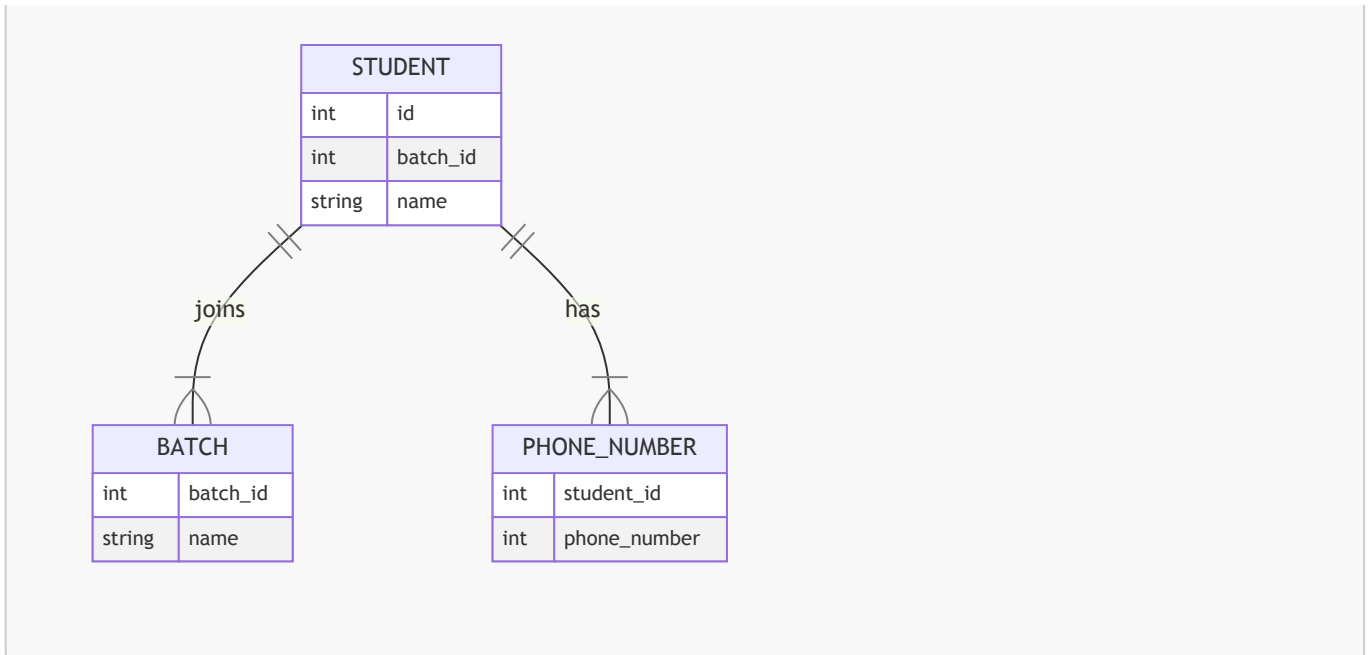
Looking at the normalised table from 3NF

ID NAME PHONE BATCH_ID

We can list the following dependencies for the above table:

- **ID → NAME**
- **ID → PHONE**
- **ID → BATCH_ID**
- **PHONE → NAME**

It can be clearly seen that the last dependency violates BCNF as phone is not a primary key. The above table can be normalised by creating a separate table for phone information.



Transactions

A database transaction symbolizes a unit of work performed against a database, and treated in a coherent and reliable way independent of other transactions.

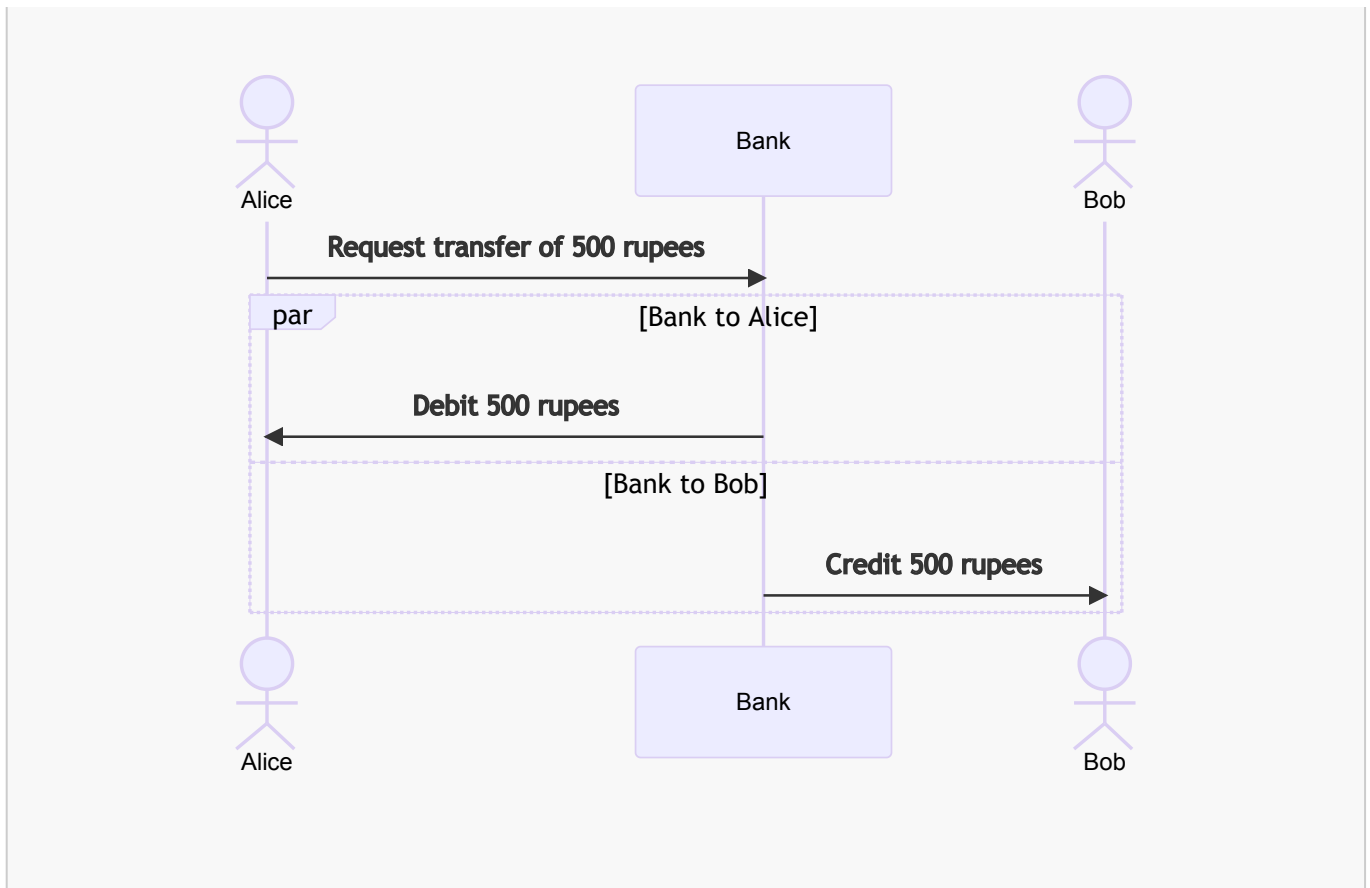
Transactions in a database environment have two main purposes:

- To provide reliable units of work that allow correct recovery from failures and keep a database consistent even in cases of system failure.
- To provide isolation between programs accessing a database concurrently. If this isolation is not provided, the programs' outcomes are possibly erroneous.

Imagine a scenario where Alice has pay Bob 500 rupees. Instead of paying each other directly, they use a bank.

Alice requests a payment to be made to Bob. The bank has to do two operations:

1. Deduce the amount to be paid to Bob from Alice's account.
2. Add the amount to be paid to Bob's account.



The above set of operations need to be performed in a single transaction else, in case of failure, the bank will be left in an inconsistent state. For example, the debit operation succeeds but the credit operation fails. Now Alice has 500 rupees less but Bob does not have 500 rupees more.

ACID properties

A sequence of database operations that satisfies the ACID properties (which can be perceived as a single logical operation on the data) is called a transaction.

Atomicity

means either all successful or none.

Consistency

ensures bringing the database from one consistent state to another consistent state.
ensures bringing the database from one consistent state to another consistent state.

Isolation

ensures that transaction is isolated from other transaction.

Durability

means once a transaction has been committed, it will remain so, even in the event of errors, power loss etc.

Atomicity

Atomicity is the guarantee that series of database operations in an atomic transaction will either all occur (a successful operation), or none will occur (an unsuccessful operation)

If the logical transaction consists of transferring funds from Alice to Bob, this may be composed of

- first removing the amount from account A,
- then depositing the same amount in account B.

We would not want to see the amount removed from Alice before we are sure it has also been transferred into Bob. Then, until both transactions have happened and the amount has been transferred to Bob, the logical transfer has not occurred.

Consistency

- The integrity constraints are maintained so that the database is consistent before and after the transaction.
- The execution of a transaction will leave a database in either its prior stable state or a new stable state.

In the banking example of Alice and Bob, the total value of the accounts is 2000, split equally. Now if Alice transfers 500 rupees to Bob, the total value of the accounts should still be the same.

But in case of credit failure, the total value will now be 1500 and hence the system will be left in an inconsistent state.

Isolation

Ensures that concurrent execution of transactions leaves the database in the same state that would have been obtained if the transactions were executed sequentially

Durability

Guarantees that once a transaction has been committed, it will remain committed even in the case of a system failure (e.g., power outage or crash). This usually means that completed transactions (or their effects) are recorded in non-volatile memory.

References

- [Data Anomalies in DBMS](#)
- [Data Anomalies II](#)
- [Functional Dependency](#)
- [ACID](#)
- [Transactions](#)

Further Reading

- [Isolation](#)