# Data Cleaning, Transformation and Modelling

<span style="color:red">Question 1 : What is Data Loading in Power BI and why is it considered the first step of analysis?</span>

**Data Loading** in Power BI, often initiated via the "Get Data" feature, is the foundational process of connecting the Power BI Desktop to external data sources and importing that information into the application. It acts as the bridge between raw data storage—whether in Excel files, SQL databases, cloud services, or web pages—and the analytical engine of Power BI.

It is considered the **first and most critical step** of analysis for several reasons:

- **Establishing the Foundation:** We cannot analyze what is not there. Data loading brings the raw material Without this step, the subsequent stages of cleaning, modeling, and visualization simply cannot happen.
- **Defining Scope and Source:** This stage determines *what* data is available for analysis. It involves selecting specific tables or datasets, which sets the boundaries of your entire project.
- **Triggering ETL (Extract, Transform, Load):** Loading data immediately leads into the **Power Query Editor**.

This is where the raw data is cleaned and shaped. Since transformation logic relies on the structure of the incoming data, the loading phase dictates how accurate and efficient our data preparation will be.

Question 2 : Explain the difference between "Load" and "Transform Data" in Power BI.

When you connect to a new data source in Power BI, we are presented with a Navigator window containing two primary action buttons. The choice you make here dictates the immediate next step in your workflow.

## 1. Load

Selecting **Load** tells Power BI to take the data exactly as it currently appears in the source and import it directly into the **Power BI Data Model**.

- The window closes, and Power BI begins processing rows immediately. The data becomes available in the "Data" and "Report" views.
- We have to use this *only* if your data is already perfectly clean, formatted, and requires no changes (e.g., a highly curated SQL view or a verified Excel table).

- If the data has errors, unnecessary columns, or wrong data types, you will clutter your model and potentially slow down performance.

### 2. Transform Data (Recommended)

Selecting **Transform Data** (formerly "Edit") redirects the data to the **Power Query Editor** before it enters the Data Model.

- The Power BI interface opens a separate window (Power Query). Here, we can clean the data—filtering rows, renaming columns, changing data types, or replacing values—without affecting the original source.
- This is the **best practice** for most scenarios. It allows us to shape the data to fit your analytical needs, ensuring that only clean, relevant data is loaded into your report.
- It reduces the file size and improves report performance by filtering out unnecessary data *before* import.

| Feature | Load | Transform Data |
|---|---|---|
| **Destination** | Direct to Data Model & Report View | Power Query Editor |
| **Action** | Imports data "as is" | Allows cleaning & shaping first |
| **Best For** | pre-processed data | Raw, messy, or large datasets |
| **Memory Usage** | High (imports everything) | Optimized (imports only what you keep) |

In Power BI and data modeling (specifically the Star Schema), tables are categorized into two types based on the data they hold:

**1. Fact Table**

A **Fact Table** contains the quantitative data or "metrics" of your business. These are the numbers you want to analyze (sum, average, count). It represents **events** or transactions that happened.

- Long and narrow (many rows, fewer columns), contains numbers and IDs (Foreign Keys) linking to dimension tables.

**Example from the Dataset:** If we were to model our `given dataset`, the **Fact Table** would be the central table containing the daily statistics.

- **Name:** `Fact_Daily_Cases`
- **Columns included:**
  - **Keys (to link data):** `Date` (or DateID), `Location_ID` (a generated ID linking to Country/State).
  - **Metrics (The "Facts"):**

- ■ `Confirmed_Cases`
- ■ `Recovered_Cases`
- ■ `Deaths`
- ■ `Active_Cases`
- ■ `Tests_Conducted`
- ■ `Hospitalized`

These columns contain numbers that you will aggregate (e.g., "Total Deaths in Brazil" or "Average Hospitalized in UK").

## 2. Dimension Table

A **Dimension Table** contains descriptive attributes or "context" related to the facts. These are the text fields you use to slice, dice, and filter your data. Short and wide (fewer rows, more descriptive columns), contains unique values and a Primary Key.

**Example from your Dataset:** From your file, we would extract the descriptive columns to create **Dimension Tables**.

- **Proposed Name:** `Dim_Location`
  - **Columns included:** `Country`, `State`.
  - **Usage:** Used to filter the report (e.g., a Slicer for "Country").
- **Proposed Name:** `Dim_Vaccination_Status`

- **Columns included:** `Vaccination_Status`.

<span style="color:red">Question 4 : Why is Star Schema preferred over Snowflake Schema in Power BI?</span>

In Power BI, the **Star Schema** is universally preferred over the Snowflake Schema primarily due to **performance** and **usability**.

Power BI runs on the **VertiPaq engine**, which is an in-memory, columnar database. It is highly optimized for scanning single tables and compressing data. While Snowflake schemas save storage in traditional databases by normalizing data (splitting tables to avoid repetition), Power BI's compression engine handles data redundancy so efficiently that this storage benefit is negligible.

Instead, the Snowflake schema introduces a major penalty: **Joins**. Every relationship in Power BI has a computational cost. A Star Schema ensures that any dimension is just *one relationship away* from the fact table. A Snowflake schema often requires traversing chains of tables (e.g., Sales > Product > SubCategory > Category) to filter data. This "relationship hopping" significantly slows down query performance and report rendering.

Furthermore, Star Schemas are **simpler for users**. A business user wants to find "Product Category" in the "Product" table, not hunt for it in a separate linked table. This structure also simplifies **DAX formulas**, avoiding the need for complex, performance-killing bi-directional filters often required in Snowflake designs .

<span style="color:red">Question 5 : Identify and remove duplicate records based on Date, Country, and State.</span>

## 1. Open Power Query

- We have  to use the **Home** tab in Power BI Desktop.
- Click **Transform data**.

## 2. Select the Columns

- Hold down the **Ctrl** key on our keyboard.
- Click the headers of the three columns you want to check: `Date`, `Country`, and `State`.

## 3. Removing Duplicates

- With those three columns still highlighted, **Right-click** on one of the headers clicking the **Ctrl** key .
- Select **Remove Duplicates** from the menu.

# 4. Check our Work

- Look at the "Applied Steps" panel on the right. We will be able to see a new step called **"Removed Duplicates"**.
- Click **Close & Apply** (top left) to save the changes.

Note : The date format in the raw data was incorrect. Directly changing the datatype in Power BI resulted in an "Error" for the entire column, even when attempting the change locale.

## Step 1: Open Power Query Editor

1. We have  to use the **Home** tab in Power BI Desktop.
2. Click **Transform data**. This opens the Power Query Editor window.

## Step 2: Identify the Nulls

1. Locate the **Vaccination_Status** column.
2. We will see some cells that are literally empty.
3. To verify, click the small arrow (filter icon) on the column header. We will see (`blank`) in the list of values.
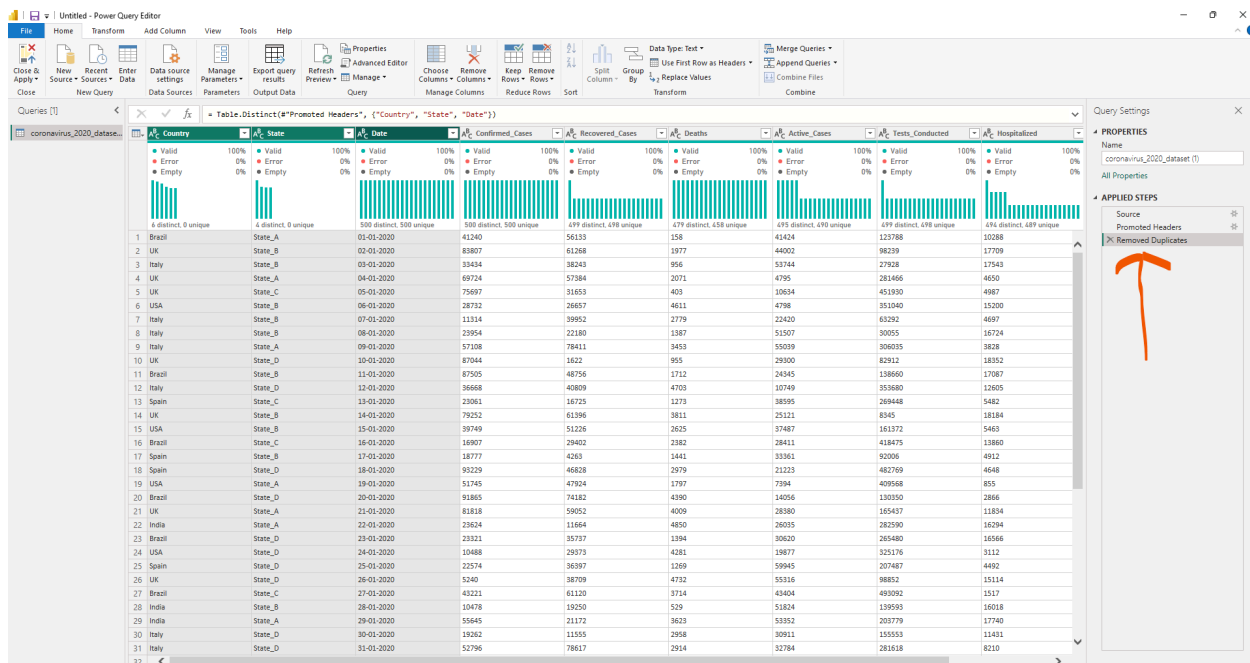
## Step 3: Replace Values

1. **Right-click** the header of the`Vaccination_Status` column.
2. Select **Replace Values** from the menu.
3. In the dialog box that appears:
   - **Value To Find:** Type `null` ( leave this empty as your cells are just blank ).
   - **Replace With:** Type the value you want, such as `Not Mentioned`.

4. Click **OK**.

## Step 4: Verify and Save

1. The *blank* values in the column should now instantly change to "Not Mentioned".
2. Click **Close & Apply** in the top-left corner to save these changes to your report.

Method : Using DAX

**Go to Data View:** Click the Table icon on the left sidebar.

**New Column:**

- Right-click on your table name in the **Data** panel (right side).
- Select **New column**.

**Enter Formula:**

```
Recovery Rate =
DIVIDE('coronavirus_2020_dataset
(1)'[Recovered_Cases],'coronavirus_2020_dat
aset (1)'[Confirmed_Cases],0)
```

Question 8 : Create a summarized table showing total confirmed cases by Country.