# Strategic Data Management (STRADAM)

## Stockholm University

## Department of Computer and Systems Sciences

## Title: Profiling, Cleaning, Transforming and Analyzing Data Using Power Query

## HT2025

By Workneh Y. Ayele (Ph.D.)

## Lab Goals

**After completing this module, students should understand why and how Power Query is used for strategic data management. The goals of this module are therefore:**

- *Familiarize and understand the role of Power Query in data cleaning, profiling, and preparation for analysis and decision making.*
- *Develop practical skills in transforming raw, inconsistent data into a structured, reliable dataset using a smaller dataset, and students should be able to practice on massive datasets after completing this exercise.*
- *Improve data quality and consistency to enable trustworthy business insights.*
- *Prepare datasets for downstream analysis (PivotTables, dashboards, reporting)*
- *Briefly introduce AI-assisted data transformation to improve efficiency and scalability.*

## Learning Outcomes

- *To understand why Power Query is needed in strategic data management (data cleaning, profiling, automation, reproducibility).*
- *Ability to practice data import, query creation, and use the Query Settings pane for data transformation*
- *Ability to do data profiling and quality assessment using Column Quality, Column Distribution, and Column Profile to identify missing values, detect duplicates, spot errors, and inconsistencies*
- *Ability to make data-quality-related decisions*
- *Ability to do data cleaning and standardization ß change data types, identify errors, and clean text by removing leading and trailing spaces, standardize letter cases (e.g., upper case, lower case, sentence case)*
- *Ability to practice data transformation and feature engineering. For example, create new columns using built-in and custom transformations, using calculated fields, and merge columns using custom separators.*
- *Ability to practice data standardization for business logic. For example, standardizing categorical values (e.g., Transfer, Online Transfer, and Wire Money à Transfer).*
- *Ability to practice aggregation and analytical preparation by duplicating queries, grouping using group transformations for analytics tasks, and preparing transformed data suitable for dashboards, strategic decision support, and reporting*
- *Ability to practice AI-assisted data transformation.*

# 1 About Power Query

Power Query is Microsoft's standard data preparation and transformation engine used across Excel, Power BI, Enterprise ETL (Extraction, Transformation, and Loading) in SSIS and Azure, and Low-code platforms (Power Platform). In this tutorial, we demonstrate how data transformation is done using a smaller dataset accessible through the link provided in the instruction. The instructions also apply to larger datasets.

## 1.1 Power Query Supporting Version

Power Query in Microsoft Excel is used mainly to:

- Import data from many sources (CSV, Excel files, Access Files, SharePoint, APIs, etc.)
- Clean and preprocess data without having to code or with low-code/no-code capabilities to filter, merge, split, reshape, change data types, remove duplicates, etc.
- Automate repeatable tasks so that data can be refreshed with a single click.
- Combine multiple data sources into a single, analysis-ready table.

Excel versions that support Power Query:

- Excel 2010 and 2013 on Windows ← Power Query must be downloaded and installed separately from Microsoft.
- Excel 2016, 2019, 2021, Excel 365/Microsoft 365 ← it is built-in in these versions (not a plugin), in these versions, it appears under the ribbon Data → Get & Transform Data

## 1.2 Data pre-processing and profiling

80% of the time spent on data analysis/analytics is allocated to cleaning and pre-processing data. It is therefore essential to learn how to prepare data for analytics. It is possible to use Python for pre-processing data; however, Power Query offers a faster and more user-friendly approach to clean, transform, and reshape data without requiring code. Power Query has an intuitive interface and supports a step-by-step transformation approach, enabling the efficient handling of everyday data preparation tasks, including removing duplicates, correcting missing values, splitting and merging columns, changing data types, filtering rows, and combining data from multiple sources. Additionally, Power Query automatically records every action, making the cleaning process transparent, repeatable, and easy to update when new data is added. Therefore, Power Query is a powerful tool for efficiently and reliably preparing data before performing analytics or building reports and dashboards for businesses.

# 2 Power Query for cleaning, standardization and transforming small dataset

In this lab, you will use Power Query in Excel to clean, standardize, and transform a small dataset.

The goal is to prepare analysis-ready data that can support strategic and managerial decision-making.

Before starting, reflect on the following:

- Why is **data cleaning** important for strategic analysis?
- What problems can inconsistent data cause in dashboards and decisions?
- How does Power Query help automate and document data preparation?

## 2.1 Tools Required

- Microsoft Excel (Windows recommended; Mac has limitations)
- Power Query (built-in in Excel 2016+ / Microsoft 365)

## 2.2 Dataset Description

You will find the dataset (an Excel file) on GitHub → https://github.com/worknehy/Course-Materials---Computer-and-Systems-Sciences-/commit/c19da1b81da774475773ba4acccd3fb3ed17d441

⇨ Download the file

The dataset contains the following columns:

- Order ID ← Unique order identification number
- Order Date ← the date when the order was placed
- Arrival Date ← the date when the ordered items/products arrived
- Name ← contains first and last names, some of the names have leading and trailing spaces
- Department ← the department providing the items/products
- City ← shipment destination city
- Payment Method ← inconsistencies (For example, Transfer and Online Transfer are the same thing)
- Revenue ← revenue from the order
- Profit ← profit earned from the order
- Tracking number ← the first three characters contain the shipping company.

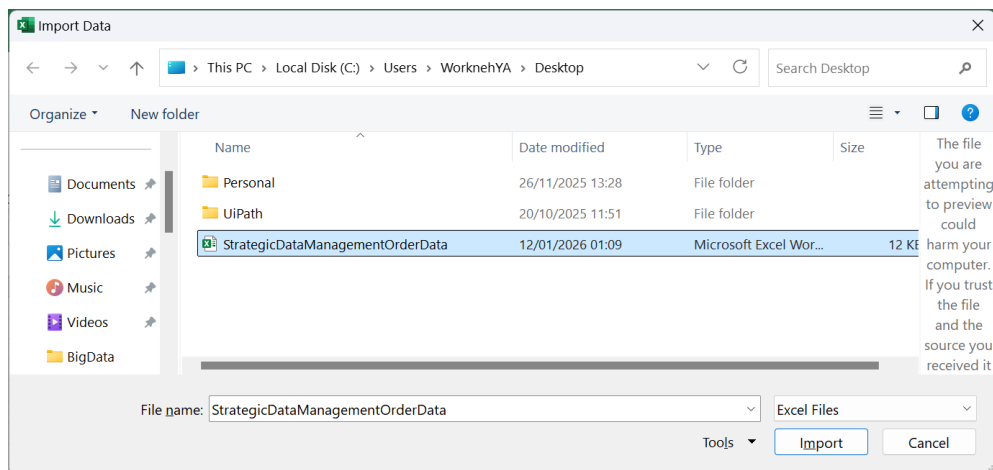| Order ID | Order Date | Arrival Date | Name | Department | City | Payment Method | Revenue | Profit | Tracking number |
|---|---|---|---|---|---|---|---|---|---|
| 73412 | 22/01/2024 | 29/01/2024 | Derek Turner | Quantum Systems | Valencia | DEBIT/CREDIT Card | 44343 | 9798 | FDX7841023956 |
| 84567 | 04/02/2024 | 12/05/2024 | Hazel Kim | Global Markets | Oslo | DEBIT/CREDIT Card | 31578 | 14298 | FDX9034158721 |
| 92304 | 27/03/2024 | 27/04/2024 | Evan Price | Global Markets | Oslo | DEBIT/CREDIT Card | 18821 | 68495 | UPS9482035612 |
| 92304 | 27/03/2024 | 27/04/2024 | Evan Price | Global Markets | Oslo | DEBIT/CREDIT Card | 18821 | 68495 | UPS9482035612 |
| 51239 | 18/04/2024 | 30/05/2024 | Jordan Wells | Supply Network | Cairo | DEBIT/CREDIT Card | 52838 | 20421 | DSV492018573 |
| 68421 | 06/12/2023 | 05/02/2024 | Clara Shaw | Supply Network | Cairo | Transfer | 41831 | 22176 | DHL9812457031 |
| 75932 | 03/06/2024 | 11/06/2024 | Felix Morgan | Quantum Systems | Valencia | Transfer | 38545 | 12048 | FDX6201948753 |
| 80315 | 28/08/2024 | 16/09/2024 | Ivy Malone | Quantum Systems | Valencia | DEBIT/CREDIT Card | 29947 | 25463 | DSV7582031640 |
| 54820 | 25/02/2024 | 09/03/2024 | Julian Frost | Quantum Systems | Valencia | Online Transfer | 46821 | 27518 | FDX9031726485 |
| 69014 | 19/05/2024 | 07/06/2024 | Jamie Abbott | Quantum Systems | Valencia | Online Transfer | 37872 | 21174 | UPS2859477021 |
| 77128 | 05/09/2024 | 28/11/2024 | Chloe Sanders | Quantum Systems | Valencia | Online Transfer | 15978 | 8288 | DHL9518476204 |
| 69305 | 12/01/2024 | 14/02/2024 | Adrian Walsh | Quantum Systems | Valencia | Online Transfer | 39504 | 9181 | FDX4928175036 |

## 2.3  Load Data into Power Query

- Open Excel and load the dataset.
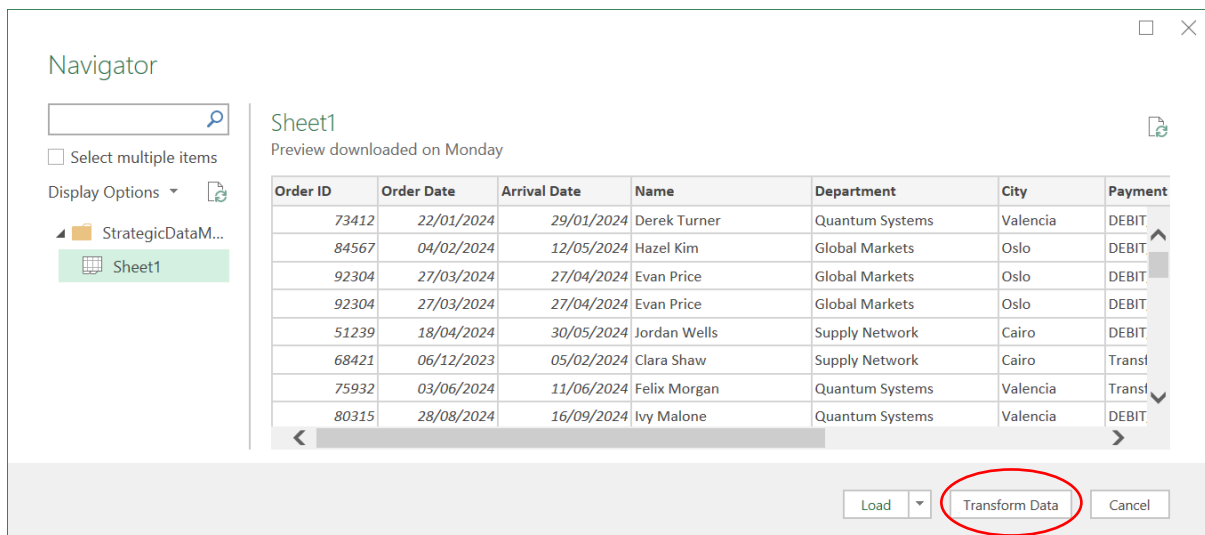  - Click on **Start** Menu



  - Select **Microsoft Excel** from your **Start** menu (or you can type **Excel** as illustrated below under (2) text box to search for Excel and select it when it shows up.
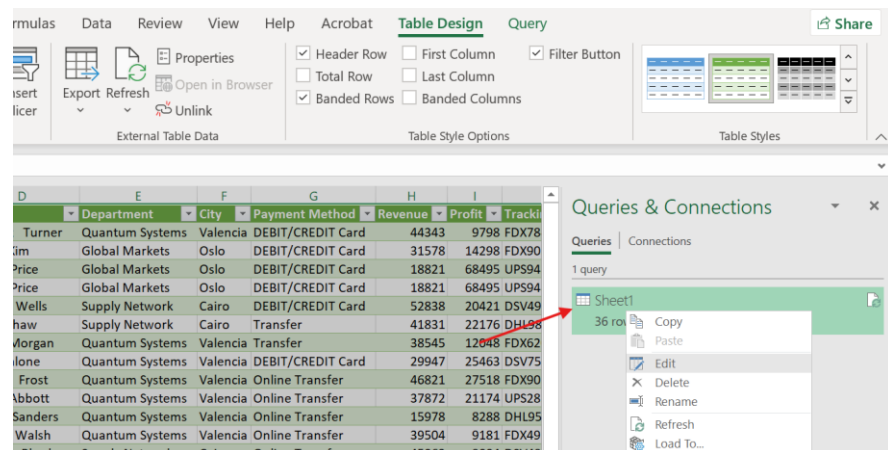


  - 
- Select **Data** to display the **Data** ribbon.
- Select **Get Data**, and the **Get Data** commands will pop up as illustrated below.
  - You can view a list of options to import data from a plethora of sources.
- Select **From File** and you can select **From Workbook** to import from an Excel file, **From Text/C**SV to import from a flat file, **From XML** to import an XML file, etc.
  - Select **From Workbook** ←
- Locate the file, you downloaded

- Click on **Import**
- Select the worksheet that contains the file as illustrated below
- Click on **Transform Data**



- o **Note:** If you mistakenly clicked on **Load,** the **Power Query Window** will not show up, but if you did, right-click on your Worksheet (**Sheet1**) from the **Queries & Connections** panel and select **Edit**, as illustrated below.

- **Power Query Editor** will open; see the illustration below.



## 2.4 Data Profiling (Understand Your Data)

*Purpose:* Understand data completeness, and reliability before cleaning.

### 2.4.1 Column Quality

*Purpose:* check column quality, check valid values, errors and empty values

- In Power Query, go to:
  - Click on the **View** tab → check **Column Quality** from the ribbon



- You can check:
  - Valid values
  - Errors
  - Empty values

### 2.4.2  Column Distribution

*Purpose:* You can check the data distribution in your columns

- In Power Query, go to:
  - Click on the **View** tab → check **Column Distribution** from the ribbon



- You can check:
  - Distinct and unique values under each column

### 2.4.3  Show whitespace

*Purpose:* You can check unnecessary and extra whitespaces before word strings and between them, as illustrated below

- In Power Query, go to:
  - Click on the **View** tab → check **Whitespaces** from the ribbon



- You can check:
  - Extra whitespaces in your columns

### 2.4.4 Column Profile

*Purpose:* Understand data composition, such as column statistics, minimum, maximum, distinct, error, empty, and count. Additionally, you can see the column value distribution.

- In Power Query, go to:
  - ○ Click on the **View** tab → check **Column profile** from the ribbon



- You can check:
  - ○ Column statistics
    - ▪ Count, Error, Empty, Distinct, Unique, Empty string, Min, and Max.
  - ○ Value distribution
    - ▪ The distribution of unique values, as illustrated in the snapshot

### 2.4.5 Remove Empty, Errors, or Duplicates (Use Column Distribution)

*Purpose:* Ensure uniqueness and consistency

- In Power Query, go to:
  - ○ **View → Column Distribution**
- Identify duplicates and unusual values:
  - • Right-click on a **Column Distribution**:
  - • You can select > **Remove Empty**, **Remove Duplicates**, or **Remove Errors** where appropriate

## 2.5 Data Type Validation

### 2.5.1 Change Data Type

*Purpose:* Prevent calculation and aggregation errors later.

- Select a column and change its data type.
    - Click on the top left corner of the columns as illustrated below.
    - Select datatype (e.g., Decimal Number, Currency, Whole Number, Percentage, Date/Time, Date, Time, Date/Time/TimeZone, etc.).

## 2.6  Creating New Columns

### 2.6.1  **Delivery Days** Columns

- Click on the **Add Column (1)** tab,
- Select **Order Date (2)**.
- Click on **Date (3)** Commands
- Select **Day (4)**
- Rename the new column by right-clicking on the new column name and selecting **Rename…**:
    - Type new name: **Delivery Days**



## 2.7  Text Cleaning and Formatting

*Purpose:* Standardize text for grouping and matching.

### 2.7.1  Formatting Name Column

- Select the **Name** column.
- Go to:
    - Select the **Transform** tab,
    - Select **Format** command,
    - Click on **Trim** ← remove Trailing and Leading spaces
- You can also:
    - Change case (e.g., UPPERCASE, lowercase, Capitalize Each Work)
    - Clean
    - Add prefix (e.g., Dr., Mr., Miss, etc.)
    - Add suffix

### 2.7.2  Removing Extra Spaces (spaces between words)

- Select the **Name** column
- Go to:
    - Select the **Transform** tab,
    - **Replace Values**



- Enter two whitespaces under **Value To Find**
- Enter one whitespace under **Replace With**
    - Click on the **OK** button

- Repeat the previous steps:
    - Enter three whitespaces under **Value To Find**
    - Enter one whitespace under **Replace With**
    - Click on the **OK** button
- The previous actions will replace:
    - Double spaces with single spaces
    - Triple spaces with single spaces

## 2.8 Splitting Columns

### 2.8.1 Split Name Column
*Purpose:* Separate first and last names for analysis.

- Select the **Name** column.
- Go to:
    - Select the **Transform** tab**,**
    - Click on the **Split Column** command,
    - Select **By Delimiter,**
    - Select **Space** and check the ,
    - Click OK
- Rename Name.1 to FirstName and Name.2 to LastName.

**Note:** Custom delimiter:

- Or if you have to a use custom delimiter (which includes ",", ";", ":" and space)



- Choose:
    - Custom delimiter
    - Enter Space (which is a whitespace and shows nothing)
- Click **OK**

## 2.9 Standardizing Categorical Data (By Adding a new column)

*Purpose:* Enforce consistent business definitions.

### 2.9.1 Payment Method Standardization

- Select the **Payment Method** column,
- Select the **Add Column** tab,
- Click on **Custom Column**
- Create a new column:
    - Enter **PaymentType** under the **New column name**
- Change:
    - DEBIT/CREDIT Card to Card
    - Online Transfer to Transfer
    - Add the following script as illustrated below:
    - The added code:
        - "(if Text.Contains ([Payment Method], "DEBIT/CREDIT Card") then "Card" else if Text.Contains([Payment Method], "Online Transfer") then "Transfer" else [Payment Method])"

### Custom Column

Add a column that is computed from the other columns.

New column name

PaymentType

Custom column formula ⓘ

```
= (if Text.Contains ([Payment Method], "DEBIT/CREDIT Card")
   then "Card" else if Text.Contains([Payment Method],
   "Online Transfer") then "Transfer" else [Payment Method])
```

Available columns

Order ID
Order Date
Arrival Date
FirstName
LastName
Department
City
~~Payment Method~~

<< Insert

Learn about Power Query formulas

✔ No syntax errors have been detected.

OK        Cancel

## 2.10 Calculated Fields

### 2.10.1 Profit Margin Calculation from Profit and Revenue

*Purpose:* Create metrics needed for strategic analysis.

- Select **Add Column** tab,

- Click on **Custom Column**
- Create a new column:
  - Enter **Profit Margin** under the **New column name**
    - Which will be calculated: **Profit Margin = [Profit]/[Revenue]**, as illustrated below.
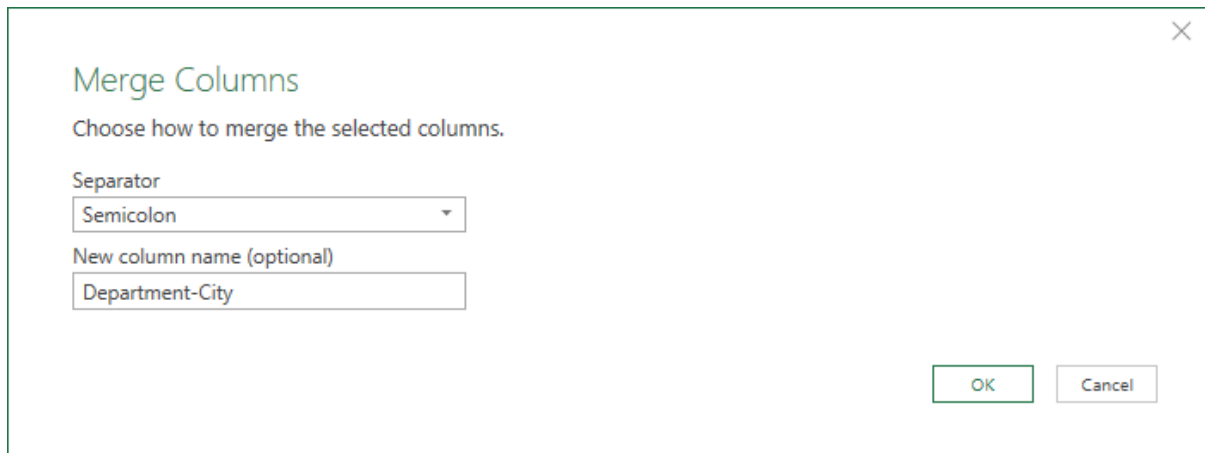


- Click **OK**

## 2.11 Merging Columns

*Purpose:* Create composite dimensions for analysis.

### 2.11.1 Merge Department and City

- Select **Department** and **City** columns (use Shift + clicking).
- Select **Add Column** tab,
- Click on the **Merge Columns** command
- Separator: semicolon (;)
- Under the **New column** name, enter "Department-City"

- Click **OK**

## 2.12 Delivery Time Analysis

### 2.12.1 New Delivery Time Column

Create a column that represents days to delivery, name it **Days to Delivery**.

- Select **Add Column** tab,
- Click on the **Custom Column** command
- Under the **New column** name, enter "**Days to Delivery**"
- Enter [Arrival Date]-[Order Date] under the formula.

## 2.12.2 Analysis by Department and Region

Duplicate the current query

- In **Queries**, right-click your cleaned query:
  - o Click **Duplicate**



- In the duplicated query (Sheet1 (2) → rename it to Delivery Days by Department):
  - o Select the **Transform** tab,
  - o Select **the Group By** command,
  - o Click on the **Advanced** button,
  - o Under Group by Select the **Department-City** column:
    - ▪ Under the **New column name:** enter **Average,** under **Operations:** select **Average,** under **Column: select Days to Delivery**
  - o Click **Add aggregate**
    - ▪ Under the **New column name:** enter **Minimum,** under **Operations:** select **Min,** under **Column:** select **Days to Delivery**
    - ▪ Under **New column name:** enter **Maximum,** under **Operations:** select **Max,** under **Column: select Days to Delivery**
  - o Your dialog box should be similar to the following snapshot:

- Aggregate delivery days by:
  - Department
  - Region
  - The result (which can be used for decision making (dashboard)):



| | ABC Department-City | ⏱ Average | ⏱ Minimum | ⏱ Maximum |
|---|---|---|---|---|
| | • Valid 100%<br>• Error 0%<br>• Empty 0%<br>4 distinct, 4 unique | • Valid 100%<br>• Error 0%<br>• Empty 0%<br>4 distinct, 4 unique | • Valid 100%<br>• Error 0%<br>• Empty 0%<br>4 distinct, 4 unique | • Valid 100%<br>• Error 0%<br>• Empty 0%<br>4 distinct, 4 unique |
| 1 | Quantum Systems;Valencia | 41.17:08:34.2857142 | 7.00:00:00 | 190.00:00:00 |
| 2 | Global Markets;Oslo | 52.04:00:00 | 10.00:00:00 | 98.00:00:00 |
| 3 | Supply Network;Cairo | 49.16:00:00 | -10.00:00:00 | 187.00:00:00 |
| 4 | Predictive Analytics;Lisbon | 32.20:34:17.1428572 | -27.00:00:00 | 86.00:00:00 |

## 2.13 Analysis by Payment Method

*Purpose:* Support strategic comparisons between payment methods.

### 2.13.1 Payment Method Analysis

- Duplicate the first query (the query that was cleaned for analysis), rename it Analysis of Payment.

- o
- Click on the **Transform** ribbon,
- Click on the **Group By** command,
  - o Under the **Group by** list, select: Payment Type
- Analyze metrics such as:
  - o Average delivery days
  - o Revenue or profit
- Your aggregation should look the same as the snapshot below.
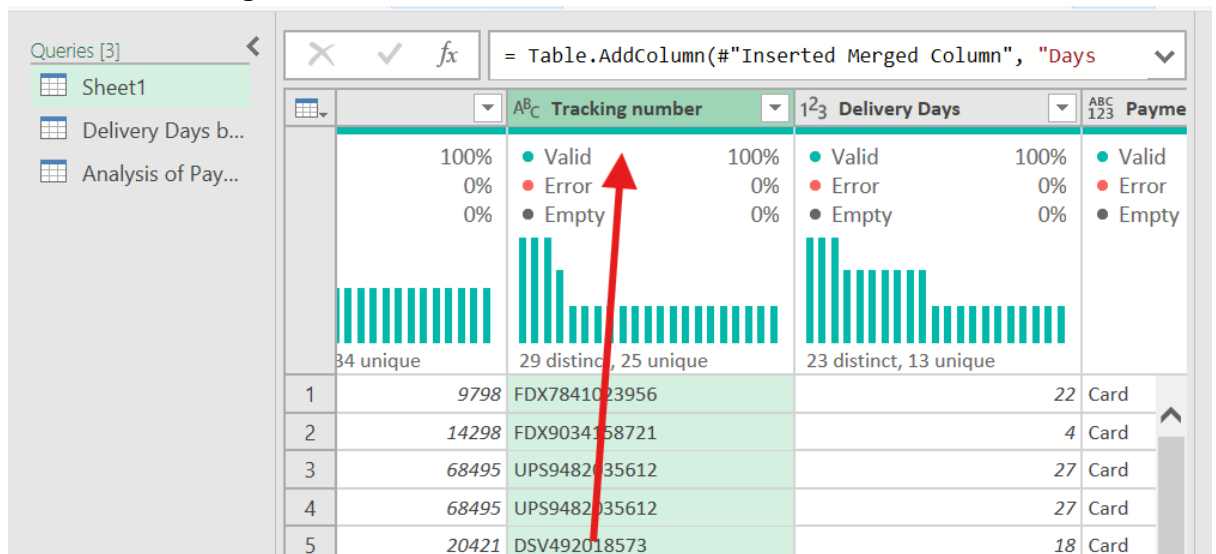
Your result:



## 2.14 AI-Assisted Transformation

### 2.14.1 Extract Shipping Provider (AI Feature)
UPS, DHL, FDX, and DSV provide shipping services. For Example, FDX7841023956, DSV492018573, and DHL5928143 are tracking numbers for FDX, DSV, and DHL, respectively.

Steps:

- Select Sheet1 query,
- Select the **Tracking Number** column.



- Click on the **Add Column** tab,
- Select **Column from Examples**
- Manually type examples (UPS, DHL, FDX, DSV), in the new column corresponding to the first three characters in the Tracking column:

| | 1²₃ Profit | ✓ | ᴬᴮ𝑐 Tracking | First Characters |
|---|---|---|---|---|
| 1 | 9798 | | FDX7841023 | FDX |
| 2 | 14298 | | FDX9034158 | FDX |
| 3 | 68495 | | UPS9482035 | UPS |
| 4 | 68495 | | UPS9482035 | UPS |
| 5 | 20421 | | DSV4920185 | DSV |
| 6 | 22176 | | DHL9812457 | DHL |
| 7 | 12048 | | FDX6201948 | FDX |
| 8 | 25463 | | DSV7582031 | DSV |
| 9 | 27518 | | FDX9031726 | FDX |
| 10 | 21174 | | UPS2859477 | UPS |
| 11 | 8288 | | DHL951847€ | DHL |
| 12 | 9181 | | FDX4928175 | FDX |
| 13 | 9004 | | DSV4920185 | DSV |
| 14 | 15688 | | UPS973204C | UPS |

*(Note to the right of rows 1–3: "Type here")*

- Power Query automatically generates the logic.

When you are done, close the Power Query and keep the changes. It is simple to return to Power Query by double-clicking on your queries, under **Queries & Connections.**

# 3  Final Reflection

Answer briefly:

- How did Power Query improve data quality?
- Which transformations were most important for analysis?
- How can this workflow scale to larger datasets?