# Business Intelligence
## Practical #2

**VSIT**

| Name | Sanjeev Gupta | Roll Number | 21302B0023 |
|---|---|---|---|
| Class | TYBScIT | Division | C |
| Subject/Course: | Business intelligence | | |
| Topic | Perform Data Wrangling (ETL) | | |

## What is ETL?

ETL, which stands for extract, transform and load, is a data integration process that combines data from multiple data sources into a single, consistent data store that is loaded into a data warehouse or other target system.

As the databases grew in popularity in the 1970s, ETL was introduced as a process for integrating and loading data for computation and analysis, eventually becoming the primary method to process data for data warehousing projects.

ETL provides the foundation for data analytics and machine learning workstreams. Through a series of business rules, ETL cleanses and organizes data in a way which addresses specific business intelligence needs, like monthly reporting, but it can also tackle more advanced analytics, which can improve back-end processes or end user experiences. ETL is often used by an organization to:
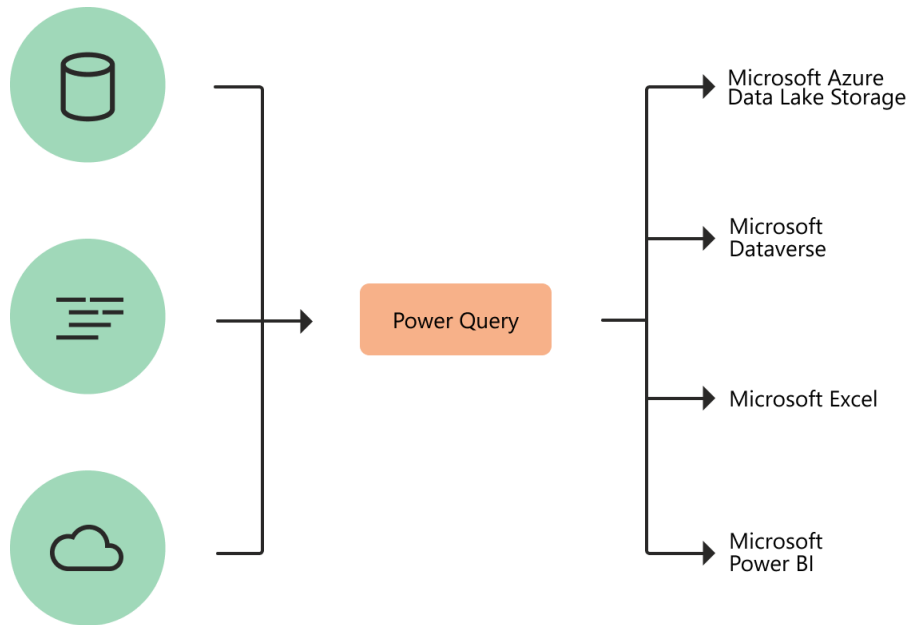
- Extract data from legacy systems
- Cleanse the data to improve data quality and establish consistency
- Load data into a target database

WHY ETL USED IN Data Science PROJECTS:

ETL is essential for organizations that wish to integrate data from multiple sources, such as databases, APIs, and flat files, into a single repository or data warehouse, and transform it into a format that is suitable for analysis or business needs. Overall, ETL is necessary for managing data at scale, ensuring data quality and accuracy, and enabling businesses to make informed decisions based on the data.

## What is Power Query?

Power Query is a data transformation and data preparation engine. Power Query comes with a graphical interface for getting data from sources and a Power Query Editor for applying transformations. Because the engine is available in many products and services, the destination where the data will be stored depends on where Power Query was used. Using Power Query, you can perform the extract, transform, and load (ETL) processing of data.



## What are the components of Power Query Editor?

The Power Query Editor is a versatile tool within Microsoft Excel that allows users to clean, transform, and combine data from various sources. It is an essential part of the Power BI suite and is used to create and manage data for data analysis and visualization. The components of the Power Query Editor can be broadly categorized into four main sections:

**Data Source Navigator:** This is the starting point for importing data into the Power Query Editor. It allows users to browse through the available data sources and import the desired data into the editor. Data sources can include Excel files, text files, databases, web pages, and more.

**Query Settings:** This section provides options to configure the imported data, such as renaming columns, changing data types, and filtering rows. Users can also perform various transformations on the data, such as splitting columns, merging columns, and removing duplicates.

**Transformation and Actions:** The Power Query Editor provides a range of built-in functions and transformations to modify the data. Users can also create custom functions using the M language, which is a functional programming language specifically designed for data manipulation. Some common transformations include sorting, grouping, and calculating aggregates.

**Advanced Editor:** This section allows users to view and edit the underlying M code that represents the data transformation process. Advanced users can modify the M code directly or use it as a reference for understanding the data manipulation process.

In summary, the Power Query Editor is a powerful tool for data manipulation and transformation, consisting of the Data Source Navigator, Query Settings, Transformation and Actions, and the Advanced Editor. These components work together to enable users to clean, transform, and combine data from various sources for analysis and visualization.

## Write the steps to perform ETL in Power BI?

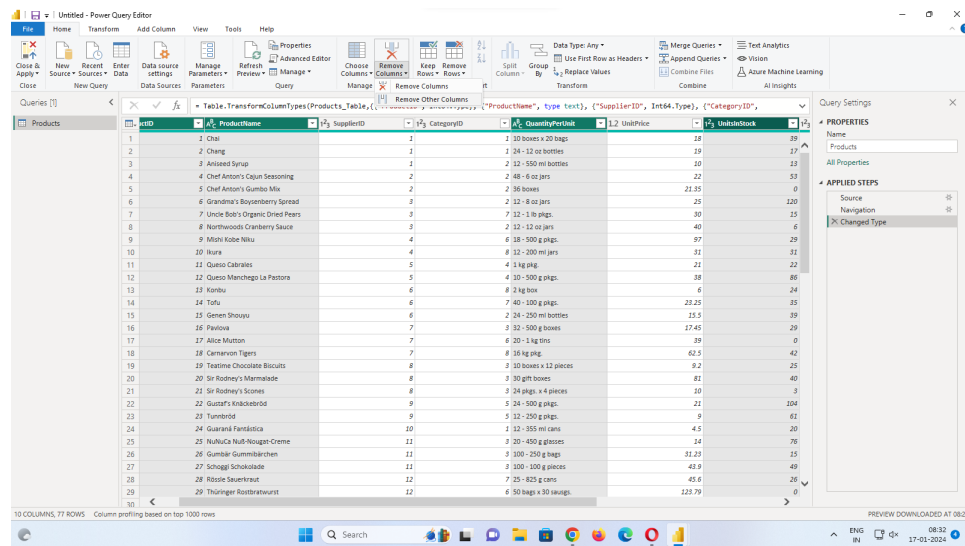1. Remove other columns to only display columns of interest.



2. Change the data type in UnitInStock column.



3. Expand the Order_Details table.

| Sort | Transform | Combine | AI Insights |
|------|-----------|---------|-------------|

| | Employee | | Order_Details | | Shipper | |
|---|---|---|---|---|---|---|

**Search Columns to Expand** [A↓Z]

◉ Expand  ○ Aggregate

- ☑ (Select All Columns)
- ☐ OrderID
- ☑ ProductID
- ☑ UnitPrice
- ☑ Quantity
- ☐ Discount
- ☐ Order
- ☐ Product

☑ Use original column name as prefix

[ OK ]  [ Cancel ]

Record (repeated many times)

| Record | Table | Record |
|--------|-------|--------|
| Record | Table | Record |
| Record | Table | Record |
| Record | Table | Record |
| Record | Table | Record |
| Record | Table | Record |
| Record | Table | Record |
| Record | Table | Record |
| Record | Table | Record |

| Manage Columns | Reduce Rows | Sort | Transform | Combine |
|---|---|---|---|---|

`"Order_Details", {"ProductID", "UnitPrice", "Quantity"}, {"Order_Details.ProductID",`

| Order_Details.ProductID | Order_Details.UnitPrice | Order_Details.Quantity | Shi |
|---|---|---|---|
| 11 | 14 | 12 | Record |
| 42 | 9.8 | 10 | Record |
| 72 | 34.8 | 5 | Record |
| 14 | 18.6 | 9 | Record |
| 51 | 42.4 | 40 | Record |
| 41 | 7.7 | 10 | Record |
| 51 | 42.4 | 35 | Record |
| 65 | 16.8 | 15 | Record |
| 22 | 16.8 | 6 | Record |
| 57 | 15.6 | 15 | Record |
| 65 | 16.8 | 20 | Record |
| 20 | 64.8 | 40 | Record |
| 33 | 2 | 25 | Record |
| 60 | 27.2 | 40 | Record |
| 31 | 10 | 20 | Record |
| 39 | 14.4 | 42 | Record |
| 49 | 16 | 40 | Record |
| 24 | 3.6 | 15 | Record |
| 55 | 19.2 | 21 | Record |
| 74 | 8 | 21 | Record |
| 2 | 15.2 | 20 | Record |
| 16 | 13.9 | 35 | Record |
| 36 | 15.2 | 25 | Record |
| 59 | 44 | 30 | Record |
| 53 | 26.2 | 15 | Record |

4. Calculate the line total for each Order_Details row.

Employee  | 🔢 1²₃ Order_Details.ProductID | 1.2 Order_Details.UnitPrice | 1.2 Order_Details.Quantity

## Custom Column

Add a column that is computed from the other columns.

New column name

Line_Total

Custom column formula ⓘ

= [Order_Details.UnitPrice]*[Order_Details.Quantity]

Available columns

ShipPostalCode
ShipCountry
Customer
Employee
Order_Details.ProductID
Order_Details.UnitPrice
Order_Details.Quantity
Shipper

<< Insert

Learn about Power Query formulas

✔ No syntax errors have been detected.

OK      Cancel

cord | 16 | 13.9 | 35 Rec

| 1.2 Order_Details.UnitPrice | 1.2 Order_Details.Quantity | ABC 123 Line_Total | |
|---|---|---|---|
| 14 | 12 | 168 | Reco |
| 9.8 | 10 | 98 | Reco |
| 34.8 | 5 | 174 | Reco |
| 18.6 | 9 | 167.4 | Reco |
| 42.4 | 40 | 1696 | Reco |
| 7.7 | 10 | 77 | Reco |
| 42.4 | 35 | 1484 | Reco |
| 16.8 | 15 | 252 | Reco |
| 16.8 | 6 | 100.8 | Reco |
| 15.6 | 15 | 234 | Reco |
| 16.8 | 20 | 336 | Reco |
| 64.8 | 40 | 2592 | Reco |
| 2 | 25 | 50 | Reco |
| 27.2 | 40 | 1088 | Reco |
| 10 | 20 | 200 | Reco |
| 14.4 | 42 | 604.8 | Reco |
| 16 | 40 | 640 | Reco |
| 3.6 | 15 | 54 | Reco |
| 19.2 | 21 | 403.2 | Reco |
| 8 | 21 | 168 | Reco |
| 15.2 | 20 | 304 | Reco |
| 13.9 | 35 | 486.5 | Reco |
| 15.2 | 25 | 380 | Reco |
| 44 | 30 | 1320 | Reco |
| 26.2 | 15 | 393 | Reco |
| 10.4 | 12 | 124.8 | Reco |
| 35.1 | 25 | 877.5 | Reco |
| 14.4 | 6 | 86.4 | Reco |
| 10.4 | 15 | 156 | Reco |

5. Rename and Reorder columns in the query

```
ed Columns",{{"Order_Details.ProductID", "ProductID"}, {"Order_Details.UnitPrice",
```

| 1²₃ ProductID | 1.2 PerUnitPrice | 1.2 Quantity | ABC 123 Lin |
|---:|---:|---:|---|
| 11 | 14 | 12 | |
| 42 | 9.8 | 10 | |
| 72 | 34.8 | 5 | |
| 14 | 18.6 | 9 | |
| 51 | 42.4 | 40 | |
| 41 | 7.7 | 10 | |
| 51 | 42.4 | 35 | |
| 65 | 16.8 | 15 | |
| 22 | 16.8 | 6 | |
| 57 | 15.6 | 15 | |
| 65 | 16.8 | 20 | |
| 20 | 64.8 | 40 | |
| 33 | 2 | 25 | |
| 60 | 27.2 | 40 | |
| 31 | 10 | 20 | |
| 39 | 14.4 | 42 | |
| 49 | 16 | 40 | |
| 24 | 3.6 | 15 | |
| 55 | 19.2 | 21 | |
| 74 | 8 | 21 | |
| 2 | 15.2 | 20 | |
| 16 | 13.9 | 35 | |
| 36 | 15.2 | 25 | |
| 59 | 44 | 30 | |
| 53 | 26.2 | 15 | |
| 77 | 10.4 | 12 | |
| 27 | 35.1 | 25 | |