# Data Cleaning with Power BI

1. Deal with Existing DataSet.
2. .csv vs .xlsx file.
3. Parsing .xml/.json
4. Exploring more features of Power Query Editor.
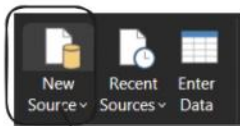5. Scrapping data/table from Website and clean them.

This will allow us to add dataset directly in power query editor



Loading data directly into PowerBI

| 123 ID | | ABC Name | | ABC Details | |
|---|---|---|---|---|---|
| • Valid | 100% | • Valid | 100% | • Valid | 100% |
| • Error | 0% | • Error | 0% | • Error | 0% |
| • Empty | 0% | • Empty | 0% | • Empty | 0% |
| 4 distinct, 4 unique | | 4 distinct, 4 unique | | 4 distinct, 4 unique | |
| 1 | Alice | | {"age": 25, "city": "New York", "skills": ["Python", "SQL"]} | |
| 2 | Bob | | {"age": 30, "city": "San Francisco", "skills": ["Java", "Spring"]} | |
| 3 | Charlie | | {"age": 28, "city": "Los Angeles", "skills": ["JavaScript", "React"]} | |
| 4 | David | | {"age": 35, "city": "Chicago", "skills": ["C#", "Azure"]} | |

Transform    Add Column    View    Tools    Help

Transpose    Data Type: Text    Replace Values    Unpivot Columns

First Row Headers    Reverse Rows    Detect Data Type    Fill    Move    Split Column    Format    Merge Columns    Extract

Count Rows    Rename    Pivot Column    Convert to List    Parse

Table    Any Column    Text

XML

JSON

```
= Table.TransformColumnTypes(#"Promoted Headers",{{"ID", Int64.Ty
```

```
= Table.TransformColumns(#"Changed Type",{{"Details", Json.Document}})
```

| 123 ID | | ABC Name | | Details | |
|---|---|---|---|---|---|
| • Valid | 100% | • Valid | 100% | • Valid | 100% |
| • Error | 0% | • Error | 0% | • Error | 0% |
| • Empty | 0% | • Empty | 0% | • Empty | 0% |
| 4 distinct, 4 unique | | 4 distinct, 4 unique | | | |
| 1 | 1 | Alice | | Record | |
| 2 | 2 | Bob | | Record | |
| 3 | 3 | Charlie | | Record | |
| 4 | 4 | David | | Record | |

Query Settings ✕

▲ PROPERTIES

Name

sample_json_column

All Properties

▲ APPLIED STEPS

Source ⚙

Promoted Headers ⚙

Changed Type

✕ Parsed JSON

ABC 123 Details

☐ (Select All Columns)

☑ age

☑ city

☑ skills

☑ Use original column name as prefix

⚠ List may be incomplete.    Load more

OK    Cancel

| ABC 123 Details.age | | ABC 123 Details.city | | ABC 123 Details.skills | |
|---|---|---|---|---|---|
| • Valid | 100% | • Valid | 100% | • Valid | 100% |
| • Error | 0% | • Error | 0% | • Error | 0% |
| • Empty | 0% | • Empty | 0% | • Empty | 0% |
| 25 | New York | | List | | |
| 30 | San Francisco | | List | | |
| 28 | Los Angeles | | List | | |
| 35 | Chicago | | List | | |

| ABC 123 Details.age | | ABC 123 Details.city | | ABC 123 Details.skills | |
|---|---|---|---|---|---|
| • Valid | 100% | • Valid | 100 | Expand to New Rows | |
| • Error | 0% | • Error | 0 | Extract Values... | |
| • Empty | 0% | • Empty | 0 | | |
| | | | | | |
| | 25 | New York | | List | |
| | 30 | San Francisco | | List | |
| | 28 | Los Angeles | | List | |
| | 35 | Chicago | | List | |

| 1²3 ID | | ABC Name | | ABC 123 Details.age | | ABC 123 Details.city | | ABC 123 Details.skills | |
|---|---|---|---|---|---|---|---|---|---|
| • Valid | 100% | • Valid | 100% | • Valid | 100% | • Valid | 100% | • Valid | 100% |
| • Error | 0% | • Error | 0% | • Error | 0% | • Error | 0% | • Error | 0% |
| • Empty | 0% | • Empty | 0% | • Empty | 0% | • Empty | 0% | • Empty | 0% |
| 4 distinct, 0 unique | | 4 distinct, 0 unique | | | | | | | |
| 1 | Alice | | | 25 | New York | | Python | | |
| 1 | Alice | | | 25 | New York | | SQL | | |
| 2 | Bob | | | 30 | San Francisco | | Java | | |
| 2 | Bob | | | 30 | San Francisco | | Spring | | |
| 3 | Charlie | | | 28 | Los Angeles | | JavaScript | | |
| 3 | Charlie | | | 28 | Los Angeles | | React | | |
| 4 | David | | | 35 | Chicago | | C# | | |
| 4 | David | | | 35 | Chicago | | Azure | | |

| ABC Name | | ABC 123 Details.age | | ABC 123 Details.city | | ABC 123 Details.skills | |
|---|---|---|---|---|---|---|---|
| • Valid | 100% | • Valid | 100% | • Valid | 100% | • Valid | 100% |
| • Error | 0% | • Error | 0% | • Error | 0% | • Error | 0% |
| • Empty | 0% | • Empty | 0% | • Empty | 0% | • Empty | 0% |

**Extract values from list**

Select a delimiter to use for concatenating list values

| None | |
|---|---|
| Colon | |
| Comma | |
| Equals Sign | |
| None | |
| Semicolon | |
| Space | |
| Tab | |
| --Custom-- | |

OK   Cancel

| 1 | Alice |
| 2 | Bob |
| 3 | Charlie |
| 4 | David |

2 new columns using split text to columns

| 1²3 ID | | ABC Name | | ABC 123 Details.age | | ABC 123 Details.city | | ABC Details.skills | |
|---|---|---|---|---|---|---|---|---|---|
| • Valid | 100% | • Valid | 100% | • Valid | 100% | • Valid | 100% | • Valid | 100% |
| • Error | 0% | • Error | 0% | • Error | 0% | • Error | 0% | • Error | 0% |
| • Empty | 0% | • Empty | 0% | • Empty | 0% | • Empty | 0% | • Empty | 0% |
| 4 distinct, 4 unique | | 4 distinct, 4 unique | | | | | | 4 distinct, 4 unique | |
| 1 | Alice | | | 25 | New York | | Python SQL | | |
| 2 | Bob | | | 30 | San Francisco | | Java Spring | | |
| 3 | Charlie | | | 28 | Los Angeles | | JavaScript React | | |
| 4 | David | | | 35 | Chicago | | C# Azure | | |

**PROPERTIES**

Name

sample_json_column

All Properties

**APPLIED STEPS**

| Source | ⚙ |
|---|---|
| Promoted Headers | ⚙ |
| Changed Type | |
| Parsed JSON | |
| Expanded Details | ⚙ |
| ✕ Extracted Values | ⚙ |

Aggregation
[Numerical
Columns]

Categorical
Columns
[Group BY]

## Group By

Specify the column to group by and the desired output.

○ Basic  ● Advanced

Product bought ▾

**New column name**
Total Quantity

**Operation**
Sum ▾

**Column**
Qty bought ▾

[OK]  [Cancel]

---

| ᴬᴮC **Product bought** ▾ | 1.2 **Total Quantity** ▾ |
|---|---|
| ● Valid           100% | ● Valid           100% |
| ● Error             0% | ● Error             0% |
| ● Empty             0% | ● Empty             0% |
| 5 distinct, 5 unique | 5 distinct, 5 unique |
| Pen Set | 271 |
| Binder | 300 |
| Pencil | 123 |
| Desk | 5 |
| Pen | 155 |

---

## Group By

Specify the columns to group by and one or more outputs.

● Basic  ○ Advanced

Product bought ▾

[ Add grouping ]

| **New column name** | **Operation** | **Column** |
|---|---|---|
| Total Quantity | Sum ▾ | Qty bought ▾ |
| Average Price Per Product | Average ▾ | Price per item ▾ |

[ Add aggregation ]

[OK]  [Cancel]

---

**Queries [6]** ‹

- ▦ Customers
- ▦ Sales
- ▦ Sales detail
- ▦ Product Information Po...
- ▦ sample_js...
- ▦ Product Si...

| × ✓ ƒx |
| ▦ _ ᴬᴮC  Product bought |
| ● Valid |
| ● Error |
| ● Empty |

Context menu:
- 🗎 Copy
- 📋 Paste
- × Delete
- ⌨ Rename
- ✓ Enable load
- ✓ Include in report refresh
- 🗎 Duplicate
- Reference
- Move To Group ▸
- Move Up
- Move Down
- Create Function...
- Convert To Parameter
- 🗎 Advanced Editor
- 🗎 Properties...

Query Settings panel:

**PROPERTIES**
Name: Product Summary
All Properties

**APPLIED STEPS**
- Source
- Promoted Headers
- Changed Type
- Removed Columns
- Filtered Rows
- Grouped Rows
- ✕ Rounded Off

Formula bar:
`= Table.TransformColumns(#"Grouped Rows",{{"Average Price Per Product", each Number.Round(_, 2), type number}})`

Queries [6]:
- Customers
- Sales
- Sales detail
- Product Information Po...
- sample_json_column
- Product Summary

| | Product bought | Total Quantity | Average Price Per Product |
|---|---|---|---|
| 1 | Pen Set | 271 | 12.48 |
| 2 | Binder | 300 | 12.66 |
| 3 | Pencil | 123 | 2.22 |
| 4 | Desk | 5 | 200 |
| 5 | Pen | 155 | 10.32 |



Split Column ▾ | Format ▾ | Merge Col... | Extract ▾ | Parse ▾

- By Delimiter
- By Number of Characters
- By Positions
- By Lowercase to Uppercase
- By Uppercase to Lowercase
- By Digit to Non-Digit
- By Non-Digit to Digit

"Coding Ninjas"     6 -> Coding

Pune@Maharashtra@India

codingNinja

12234  Abbas



Home | Transform | Add Colu...

Close & Apply ▾ | New Source ▾ | Recent Sources ▾ | Enter Data

Close | New Query

**Create Table**

| | City@State@... | DigitToNonD... | lowerToUpper | + |
|---|---|---|---|---|
| 1 | Faridabad@Ha... | 101Krishna | krishnaMadan | |
| 2 | SouthDelhi@D... | 199Abhishek | rabindraMurmu | |
| 3 | pune@mahara... | 356rahul | lionelMessi | |
| 4 | Mumbai@Mah... | 8976himesh | mugdhaSurnis | |
| 5 | Balasore@Odis... | 717Nandini | nandiniBhutani | |
| 6 | noida@UttarPr... | 78787naman | namanVerma | |
| 7 | Pune@Mahara... | 8888nitin | himeshHalli | |
| 8 | Kolkata@West... | 343rakesh | rakeshNagda | |
| 9 | kochi@kerala@... | 121rakesh | sohailKhan | |
| 10 | dahod@Gujara... | 121Abbas | abbasRajpur | |
| + | | | | |

## Split Column by Delimiter

Specify the delimiter used to split the text column.

**Select or enter delimiter**

--Custom--

@

**Split at**
- Left-most delimiter
- Right-most delimiter
- Each occurrence of the delimiter

▷ Advanced options

**Quote Character**

☐ Split using special characters

Insert special character

Table.Transform

AᴮC City@State@Country

| | Valid | 100% | | Val |
|---|-------|------|---|-----|
| ● | Error | 0% | ● | Err |
| ● | Empty | 0% | ● | Em |

10 distinct, 10 unique

| 1 | Faridabad@Haryana@India | 101Kris |
| 2 | SouthDelhi@Delhi@India | 199Abh |
| 3 | pune@maharastra@india | 356rah |
| 4 | Mumbai@Maharastra@India | 8976hi |
| 5 | Balasore@Odisha@india | 717Nar |
| 6 | noida@UttarPradesh@india | 78787n |

**Column statistics** ...

Count  10

---

| AᴮC City@State@Country.1 | AᴮC City@State@Country.2 | AᴮC City@State@Country.3 |
|---|---|---|
| ● Valid 100% | ● Valid 100% | ● Valid 100% |
| ● Error 0% | ● Error 0% | ● Error 0% |
| ● Empty 0% | ● Empty 0% | ● Empty 0% |
| 10 distinct, 10 unique | 9 distinct, 8 unique | 3 distinct, 1 unique |
| Faridabad | Haryana | India |
| SouthDelhi | Delhi | India |
| pune | maharashtra | india |
| Mumbai | Maharastra | India |
| Balasore | Odisha | india |
| noida | UttarPradesh | india |
| Pune | Maharastra | IndiA |
| Kolkata | WestBengal | India |
| kochi | kerala | india |
| dahod | Gujarat | india |

---

Merge Columns

Split Column · Format · Extract · Parse · Statis

- By Delimiter
- By Number of Characters  ry.1", type
- By Positions  Digit
- By Lowercase to Uppercase
- By Uppercase to Lowercase
- **By Digit to Non-Digit**  Split values i
- By Non-Digit to Digit

10 distinct, 10 unique

101Krishna
199Abhishek
356rahul
8976himesh
717Nandini
78787naman

---

Merge Columns

Split Column · Format · Extract · Parse · Statistics · Standard · Scientific · Trigonometry · Rounding · Information · Date · Time

Number Column    Date & Tim

- By Delimiter
- By Number of Characters
- By Positions
- **By Lowercase to Uppercase**
- By Uppercase to Lowercase
- By Digit to Non-Digit
- By Non-Digit to Digit

CharacterTransition({"0".."9"}, (c) => not List.Contains(

| Digit.1 | AᴮC DigitToNonDigit.2 | AᴮC lowerToUpper |
|---|---|---|
| | | Split values in the selected column based on transitions from a lowercase le |
| 0% | ● Error 0% | ● Error 0% |
| 0% | ● Empty 0% | ● Empty 0% |
| 9 distinct, 8 unique | 9 distinct, 8 unique | 10 distinct, 10 unique |
| 101 | Krishna | krishnaMadan |
| 199 | Abhishek | rabindraMurmu |
| 356 | rahul | lionelMessi |
| 8976 | himesh | mugdhaSurnis |
| 717 | Nandini | nandiniBhutani |
| 78787 | naman | namanVerma |

## Split Column / Format menu

**Format menu:**
- lowercase
- UPPERCASE
- Capitalize Each Word
- Trim
- Clean
- Add Prefix
- Add Suffix

**Split Column menu:**
- By Delimiter
- By Number of Characters
- By Positions
- By Lowercase to Uppercase
- By Uppercase to Lowercase
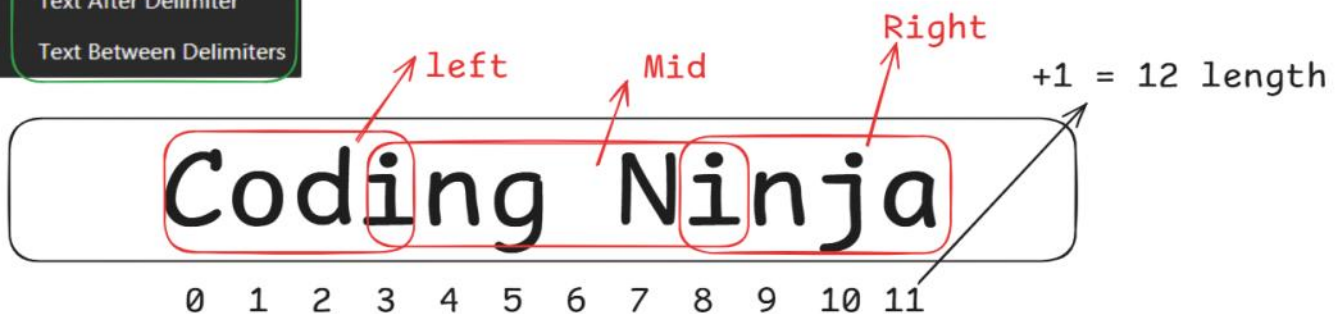- By Digit to Non-Digit
- By Non-Digit to Digit

"City", Text

100%
0%
0%

## Split Column by Number of Characters

- Valid
- Error
- Empty

9 distinct, 8 unique

Faridabad
Southdelhi
Pune
Mumbai
Balasore
Noida
Pune
Kolkata
Kochi
Dahod

Specify the number of characters used to split the text column.

**Number of characters**

5

**Split**
- Once, as far left as possible
- Once, as far right as possible
- Repeatedly

▷ Advanced options

## Extract menu

- Length
- First Characters
- Last Characters
- Range
- Text Before Delimiter
- Text After Delimiter
- Text Between Delimiters

len('Coding Ninja') 12

Left, Right, Mid

Right

left    Mid

+1 = 12 length

## Coding Ninja

```
0 1 2 3 4 5 6 7 8 9 10 11
```

## Pune@Maharastra#India

## Text Between Delimiters

Enter the delimiters that mark the beginning and end of what you would like to extract.

**Start delimiter**

@

**End delimiter**

#

▷ Advanced options

[ OK ]  [ Cancel ]

| Split_Column | Text Before Delimiter | Text After Delimiter | Text Between Delimiters |
|---|---|---|---|
| Valid 100% | Valid 100% | Valid 100% | Valid 100% |
| Error 0% | Error 0% | Error 0% | Error 0% |
| Empty 0% | Empty 0% | Empty 0% | Empty 0% |
| 6 distinct, 6 unique | 6 distinct, 6 unique | 2 distinct, 0 unique | 6 distinct, 6 unique |
| Pune@Maharastra#India | Pune@Maharastra | India | Maharastra |
| Balasore@Odisha#India | Balasore@Odisha | India | Odisha |
| dahod@gujarat#india | dahod@gujarat | india | gujarat |
| noida@utterpradesh#india | noida@utterpradesh | india | utterpradesh |
| mumbai@maharashtra#india | mumbai@maharashtra | india | maharashtra |
| Bengaluru@Karnataka#India | Bengaluru@Karnataka | India | Karnataka |

### Load

⌇ Sales detail
  Creating connection in model...

⌇ Product Information Power BI
  Creating connection in model...

⌇ sample_json_column
  Creating connection in model...

⌇ Product Summary
  Creating connection in model...

⌇ Exploring_Features
  Creating connection in model...

[ Cancel ]

### Data

🔍 Search

> ▦ Customers
> ▦ Exploring_Features
> ▦ Product Information Power BI
> ▦ Product Summary
∨ ▦ Sales
  Σ Customer Age
  Σ Customer ID
    Customer Region
  Σ Sales
> ▦ Sales detail
> ▦ sample_json_column
> ▦ Split_Feature