# DATA WRANGLING

CCDATSCL | COM221-ML

## Data Wrangling

- Data wrangling, sometimes referred to as data munging, is the process of *transforming and mapping data* from one "raw" data form into another format with the intent of making it more appropriate and valuable for a variety of downstream purposes such as analytics.
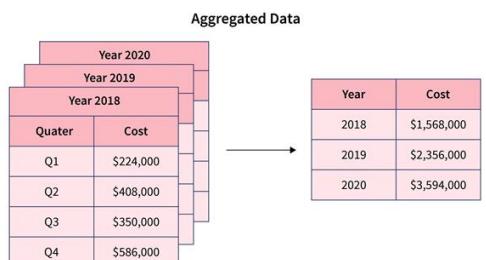


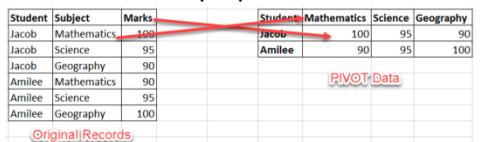## Data Wrangling tasks



1. **Discovering**
   - Discovery refers to the process of *familiarizing yourself with data* so you can conceptualize how you might use it.
   - During discovery, you may identify trends or patterns in the data, along with obvious issues, such as missing or incomplete values that need to be addressed.
   - This is an important step, as it will inform every activity that comes afterward.
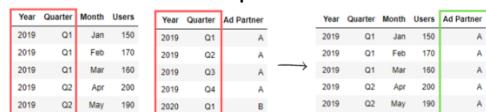
2. **Structuring**
   - The data structuring step, sometimes called data transformation, focuses on *organizing the data into a unified format* so that it is suitable for analysis.
   - **Aggregation**: Combining rows of data by using summary statistics and grouping data based on certain variables.



- **Pivoting**: Shifting data between rows and columns or transforming data into other formats to prepare it for use.



- **Joining**: Combining data from multiple tables and combining related information from disparate sources.
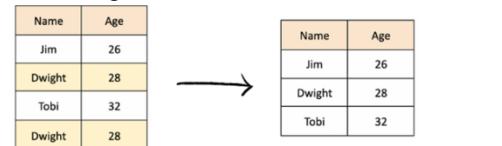


- **Data type conversion**: Changing the data type of a variable to aid in performing calculations and applying statistical methods.
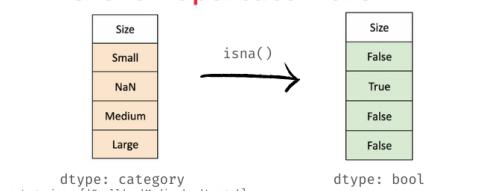


3. **Cleaning**
   - Data cleaning involves handling *missing values, removing duplicates* and correcting errors or inconsistencies.



### Remove Duplicate Rows



### Missing Values in Category Column

**Fillna with another column values**



**Missing Values Count**

## 4. Enriching

- Data enrichment involves adding new information to existing data sets to enhance their value. Sometimes called *data augmentation*
- It involves assessing what additional information is necessary and where it might come from.
- Then, the additional information must be integrated with the existing data set and cleaned in the same ways as the original data.
- *Merging data* from multiple sources to develop a more comprehensive dataset.



- *Creating new features* from existing data that can provide additional insights when analyzed.
- *Feature creation* involves generating new features from domain knowledge or by observing patterns in the data.
- It can be:
  **Domain-specific**: Created based on industry knowledge like business rules.
  **Data-driven**: Derived by recognizing patterns in data.
  **Synthetic**: Formed by combining existing features.

## 5. Validating

- This step involves verifying the accuracy and consistency of the wrangled data. First, validation rules must be established based on business logic, data constraints and other issues.
  - ✓ Data type validation: Helping ensure correct data types.
  - ✓ Range or format checks: To verify values fall within acceptable ranges and adhere to certain formats.
  - ✓ Consistency checks: Making sure that there is a logical agreement between related variables.
  - ✓ Uniqueness checks: Confirming that certain variables (such as customer or product ID numbers) have unique values.
  - ✓ Cross-field validation: Checking for logical relationships between variables (for example, age and birthdate).
  - ✓ Statistical analysis: Identifying outliers or anomalies by using descriptive statistics and visualizations.

## 6. Publishing

- Once your data has been validated, you can publish it. This involves making it available to others within your organization for analysis.