

# Data Preparation

Harvesting, Storing, and Retrieving Data





# Data Preparation



Deleting Unnecessary  
Data

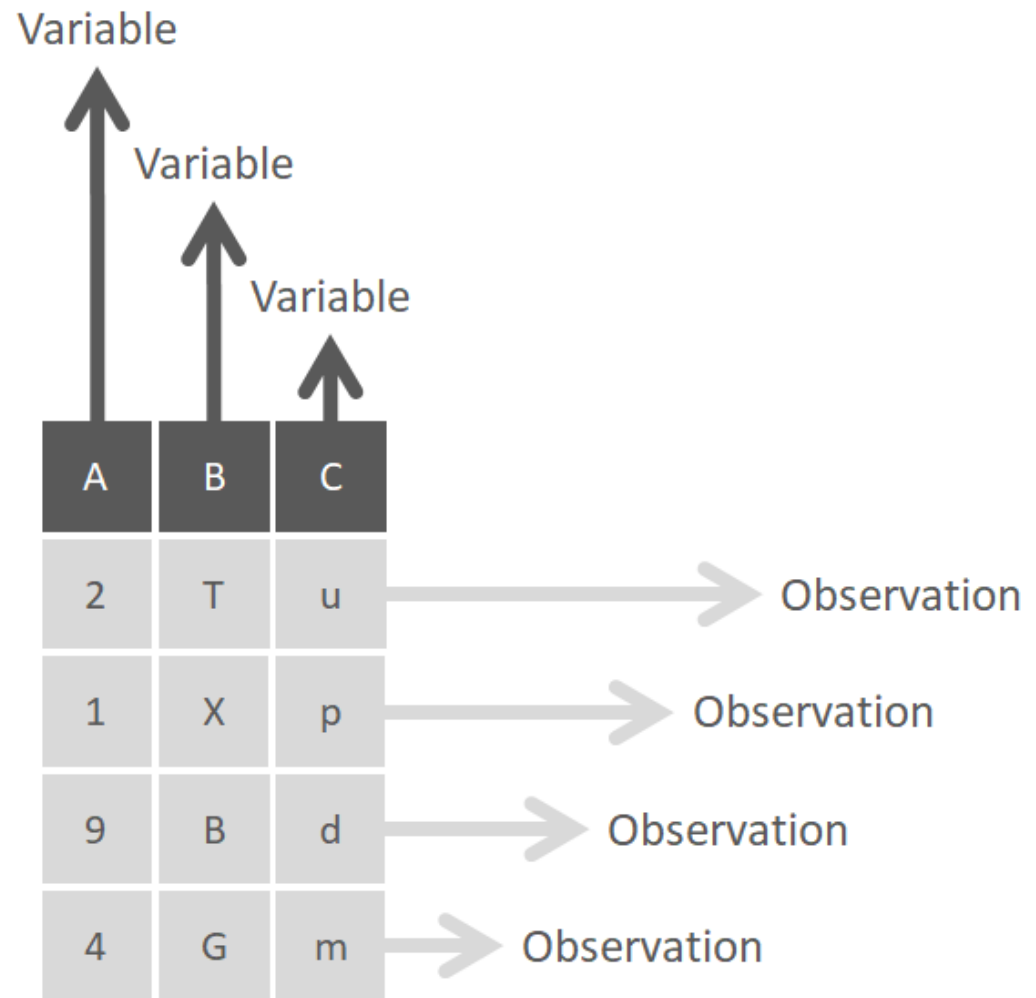
Consolidating/Separating  
Fields

Changing Formats

Transforming Data

Making Corrections

# Quick Review



[Image Source](#)

# Main Tasks for a Single Dataset

## Sort/Arrange

Order rows by value or characters

## Filter

When you only want certain rows/columns

## Transform

Creation on new record fields from existing data

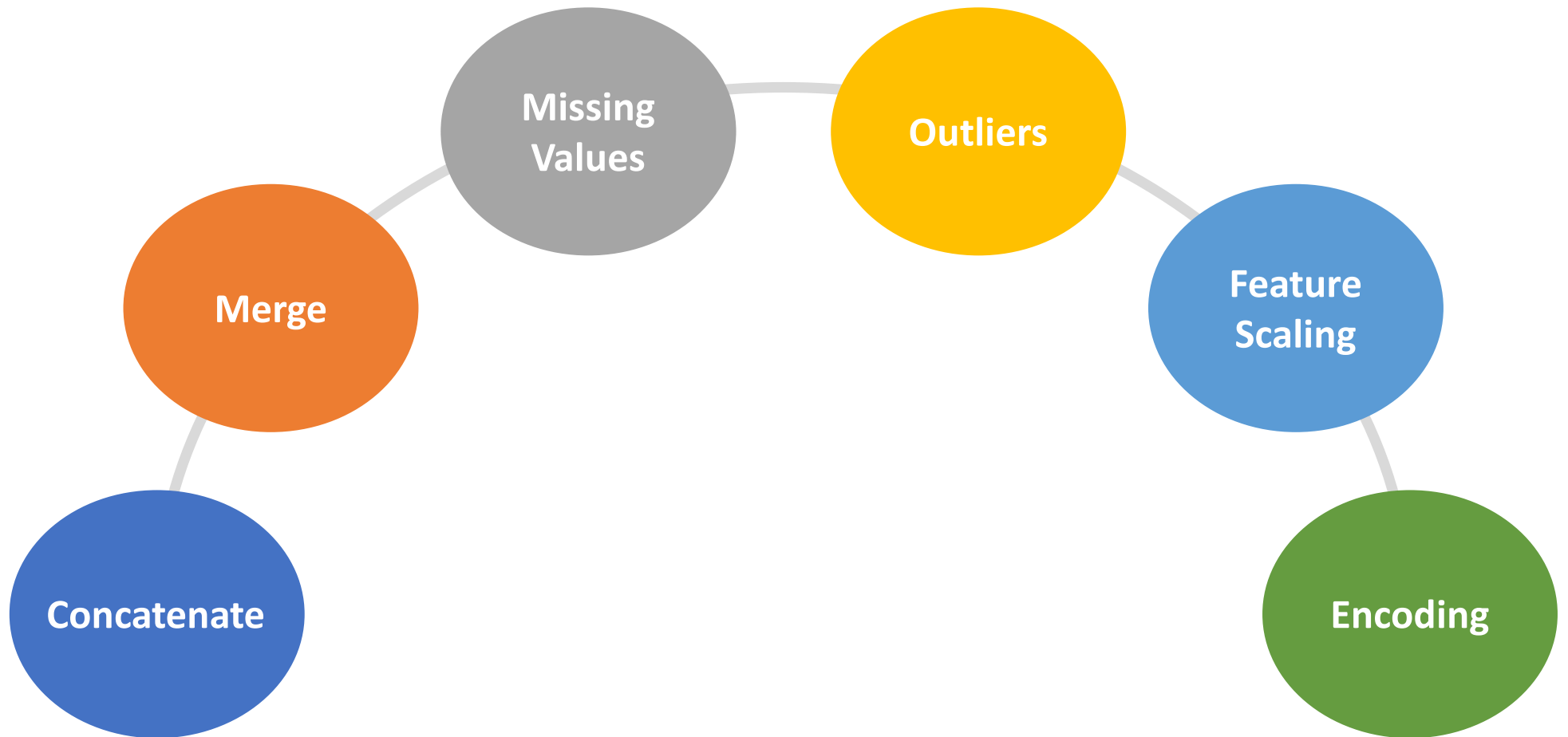
## Select

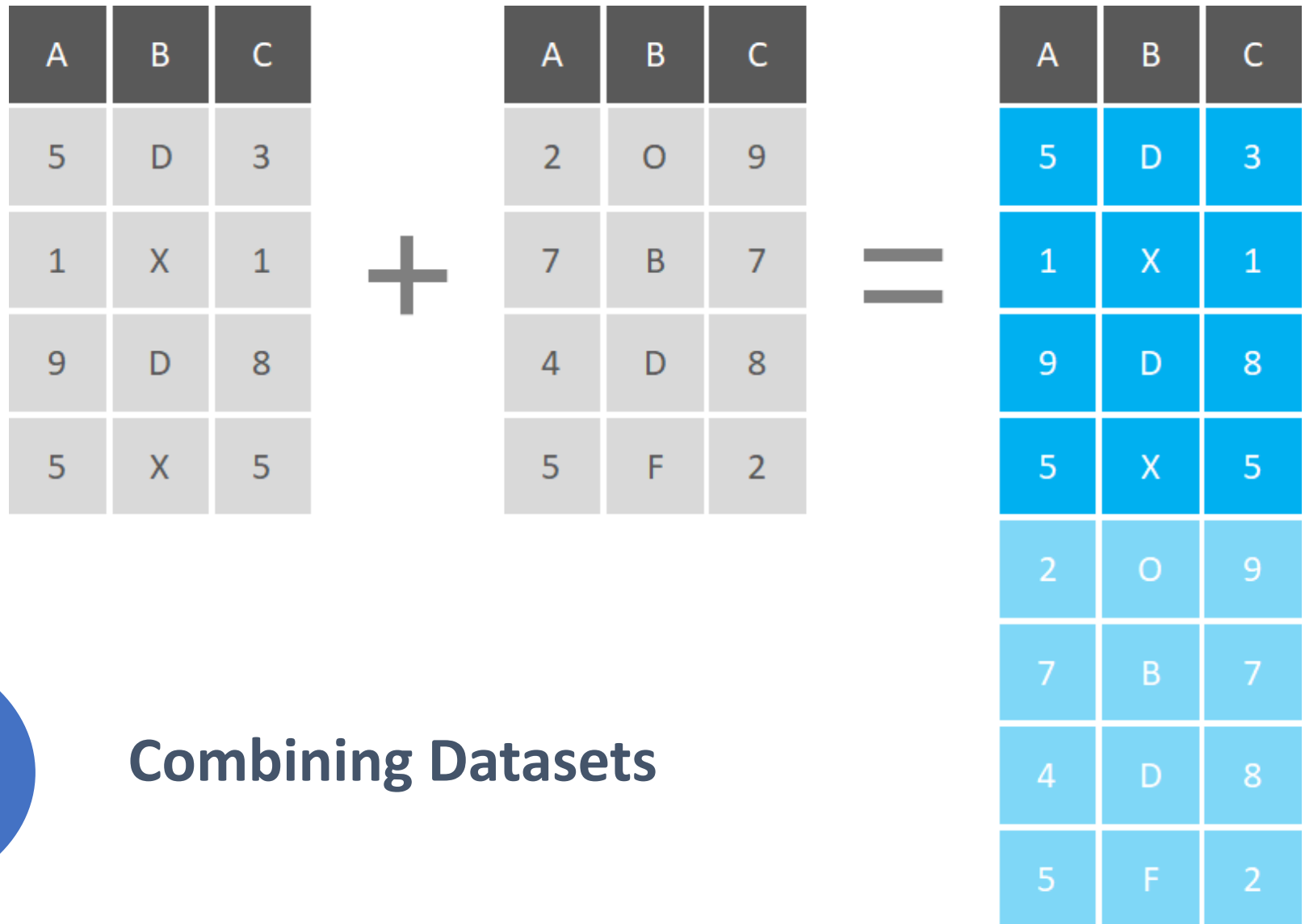
Choosing columns based on a defined criteria

## Aggregate

Gathering and/or summarizing information

## Combining Datasets



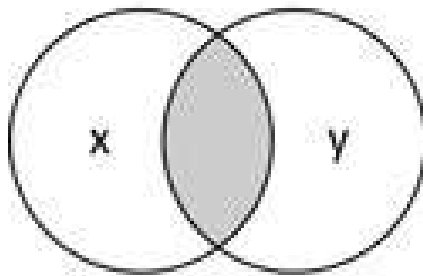


Concatenate

Combining Datasets

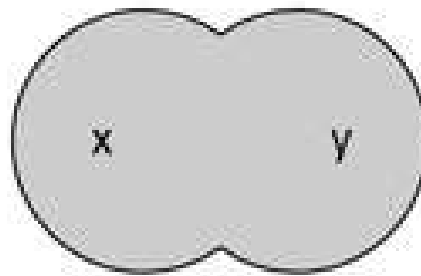
## Combining Datasets

**how='inner'**



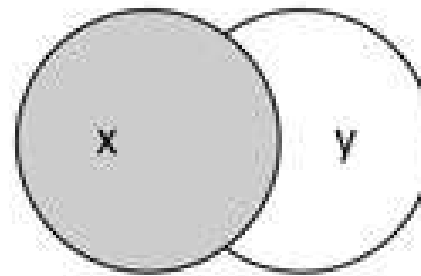
natural join

**how='outer'**



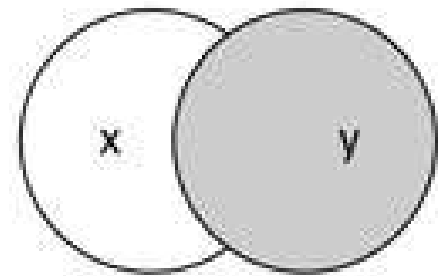
full outer join

**how='left'**



left outer join

**how='right'**

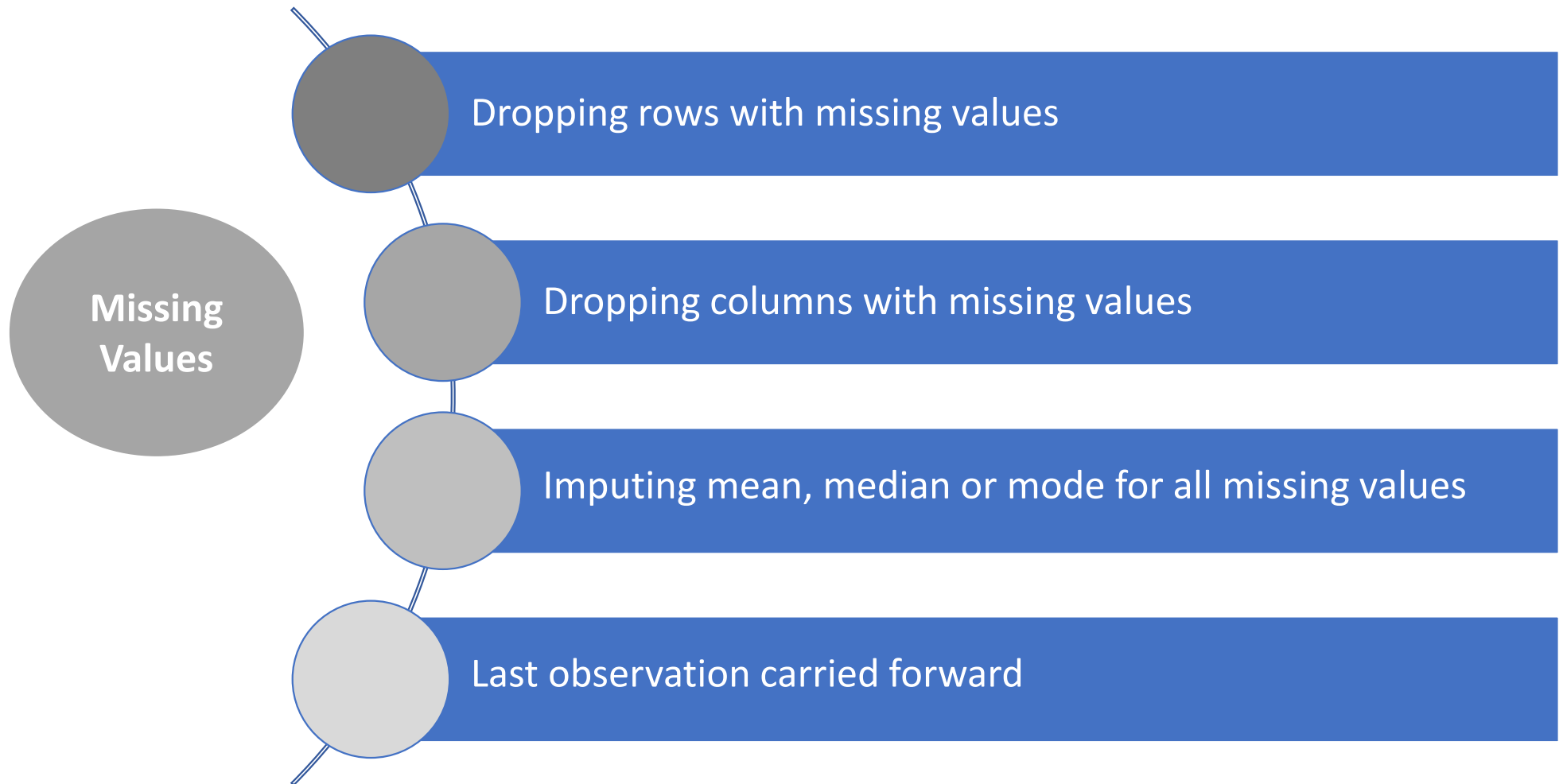


right outer join

Merge

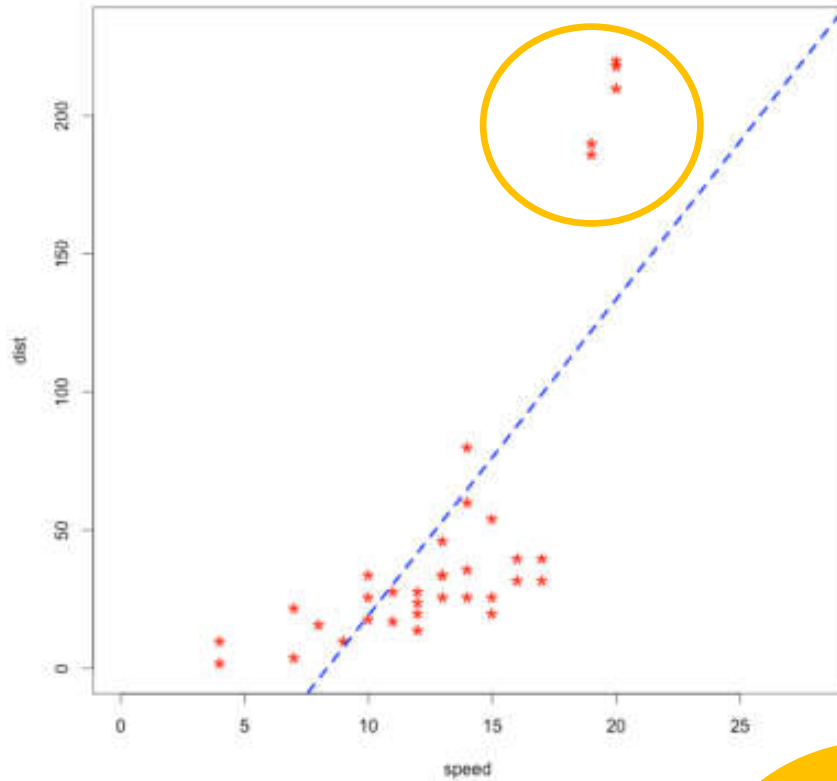


# Combining Datasets

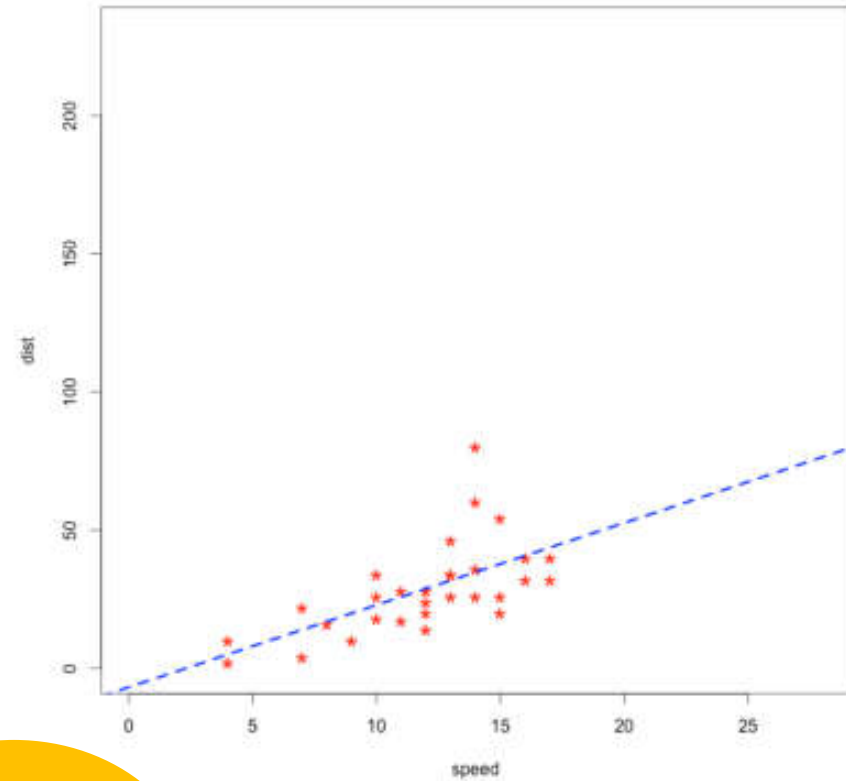


# Combining Datasets

With Outliers



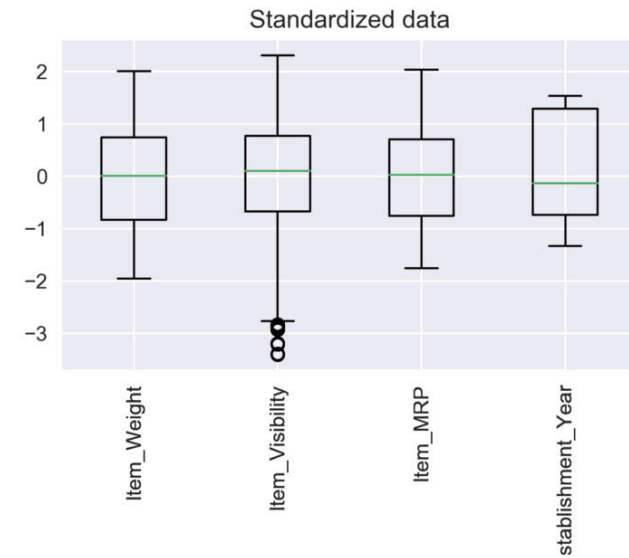
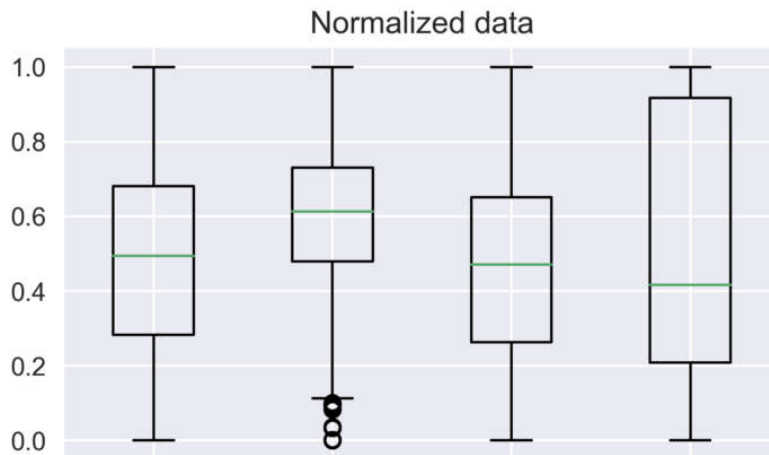
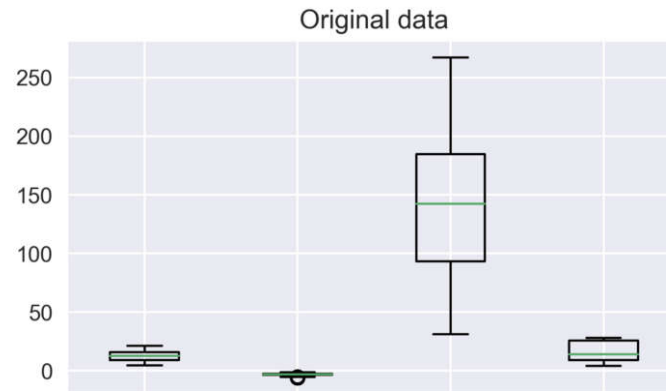
Outliers removed  
A much better fit!



Outliers

# Combining Datasets

Feature  
Scaling





01

You may find data in multiple datasets and in many different formats.



02

You will need to make decisions in the process of data wrangling.



03

Keep in mind that your actions will have consequences.

# SUMMARY

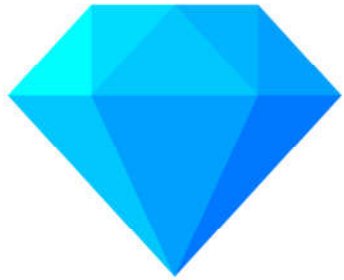
# Data Preparation: Homework

Harvesting, Storing, and Retrieving Data



**OpenRefine**





**OpenRefine**

Open source

Browser-based

Designed for tabular data

Supports CSV, Excel, XML, Google docs and  
RDF files



# OpenRefine

