# Data Wrangling and Data Analysis Welcoming

**Hakim Qahtan**

Department of Information and Computing Sciences

Utrecht University

Utrecht University

# Data Science - Myth

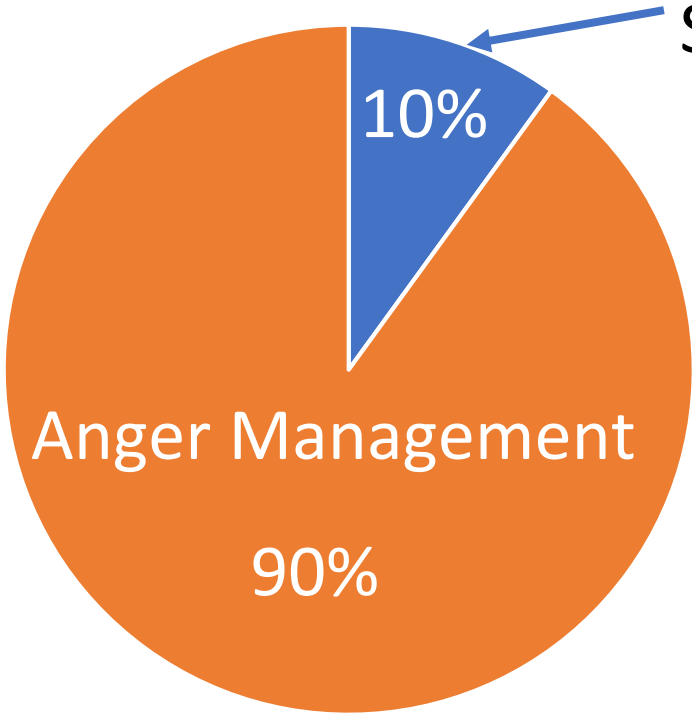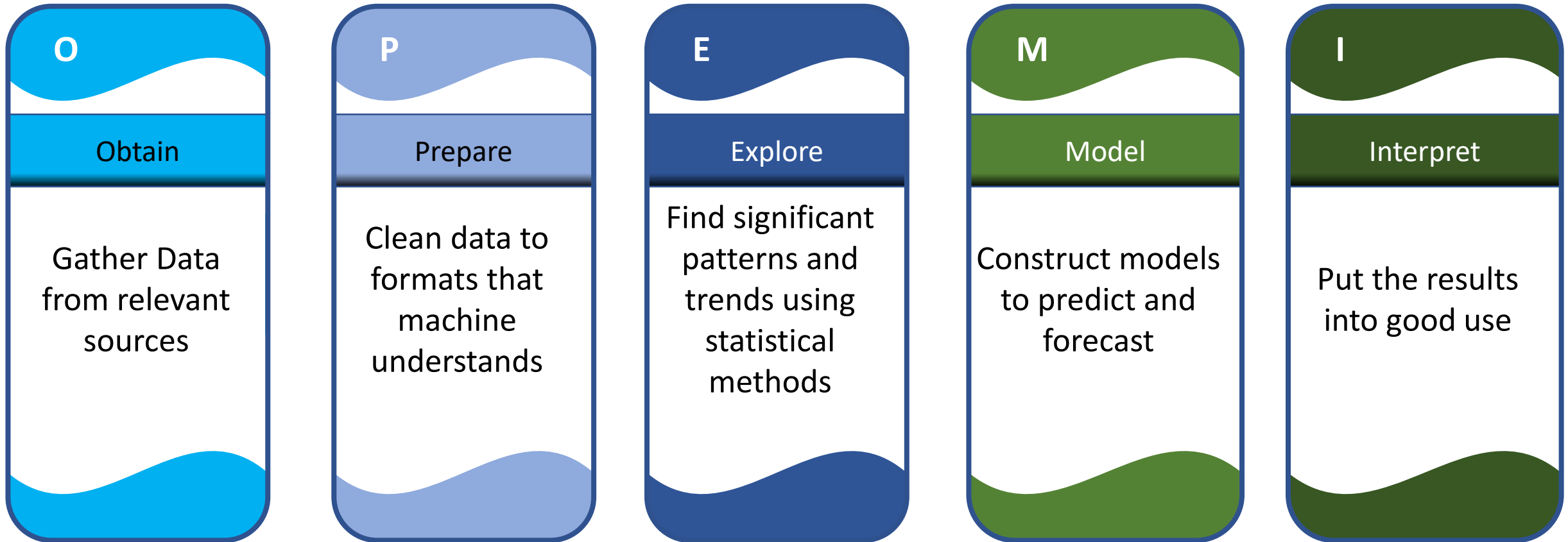# Data Science - Reality



- Data Scientists: 80% of their time looking for and cleaning the data
  - Sometimes, they are called data janitors

# What about Data Wrangling?

# Data Science Process



| O | P | E | M | I |
|---|---|---|---|---|
| Obtain | Prepare | Explore | Model | Interpret |
| Gather Data from relevant sources | Clean data to formats that machine understands | Find significant patterns and trends using statistical methods | Construct models to predict and forecast | Put the results into good use |

Utrecht University

# Data Extraction (Obtain or Gather)



JSON

XML

CSV

**Data Analytics**

Microsoft SQL Server

- It is recommended to install **PostgreSQL**
- In this course, data extraction will be done using **SQL** queries and **Python** code.

- Data is stored in different formats
- First step is to extract the data required for the analysis.

Utrecht University

# Data Preparation

# Data Exploration

Visualize and understand all available data

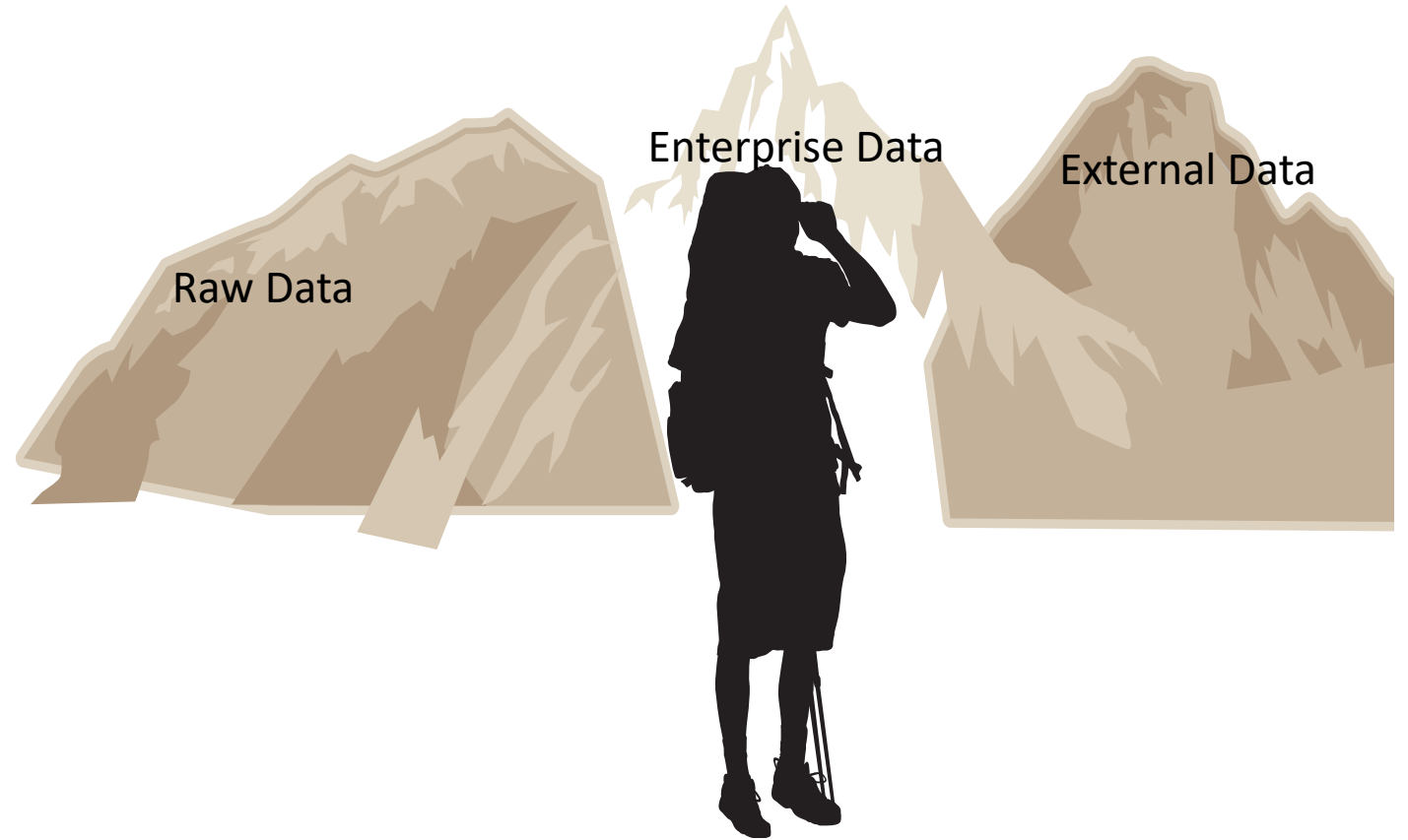Establish connections to access the data

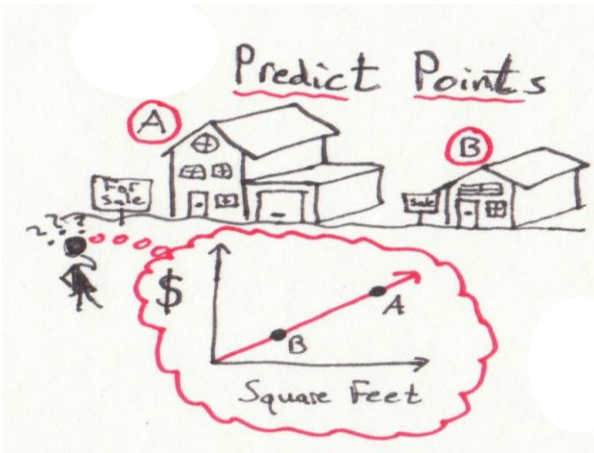Use appropriate tools to separate the most useful content

Discover hidden insights

Raw Data

Enterprise Data
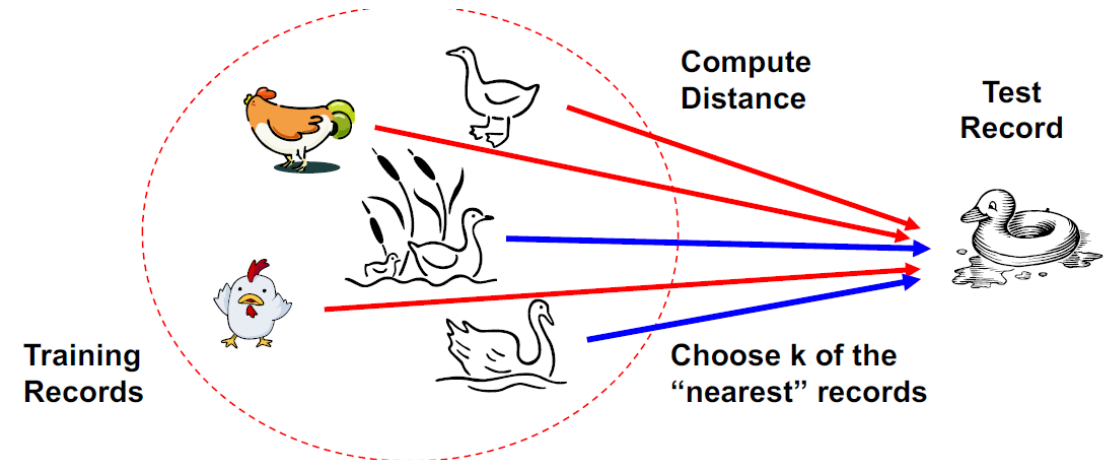
External Data

Utrecht University

# Modeling

- Wide variety of data models
  - Which model better suits your data?
  - What is the problem that you are trying to solve?
- We can use the libraries of Python or R to implement different data analytics models



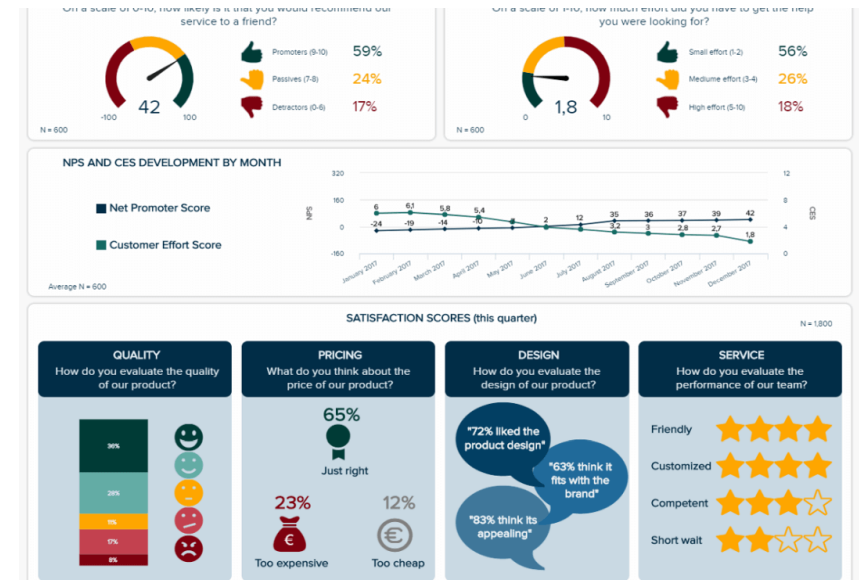Predict unknown values based on a set of known values



Group similar items together



Classify an item: if it walks like a duck, quacks like a duck, then it is probably a duck

Utrecht University

# Interpretation

- A good interpretation of the data and the results can help in:
  - Decision-making
  - Anticipating needs
    - Shazam (music identification application)
  - Cost efficiency: cost-reduction opportunities
    - Intel reduces the 19000 test and saved 3Million USD
  - Clear foresight
    - Companies know about their performance from data they collect

*Thank you*