

Feature Analysis and Visualization

MIE223
Winter 2025

1 Feature Analysis and Visualization

1.1 Exploratory Data Analysis (EDA)

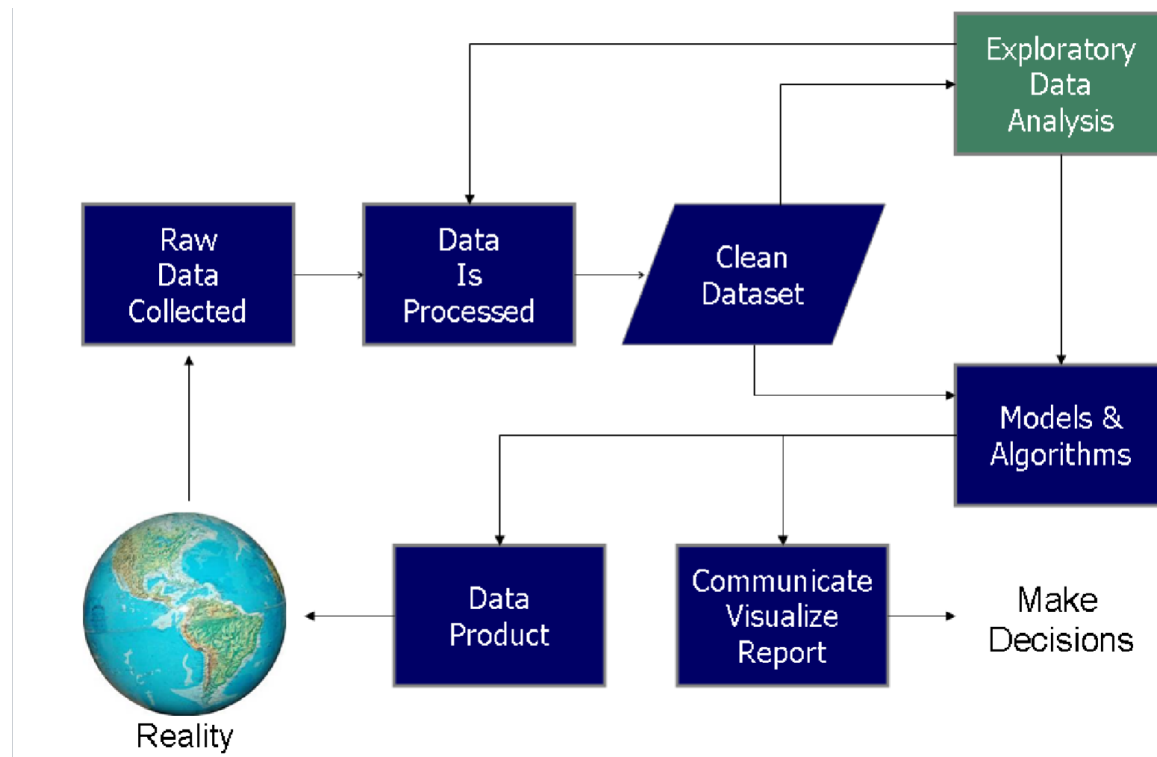
John Tukey (Bell Labs, 1970's) (Inventor of the box-plot)

- The first time corporations had a lot of digital data
- Data on semiconductor processes, networks

Observed that statistics was preoccupied with hypothesis testing

- But there were no established methodologies for **generating (data-driven) hypotheses**
- Espoused a visual methodology and a new language S (R became the free version)

1.2 Data Science Process



3

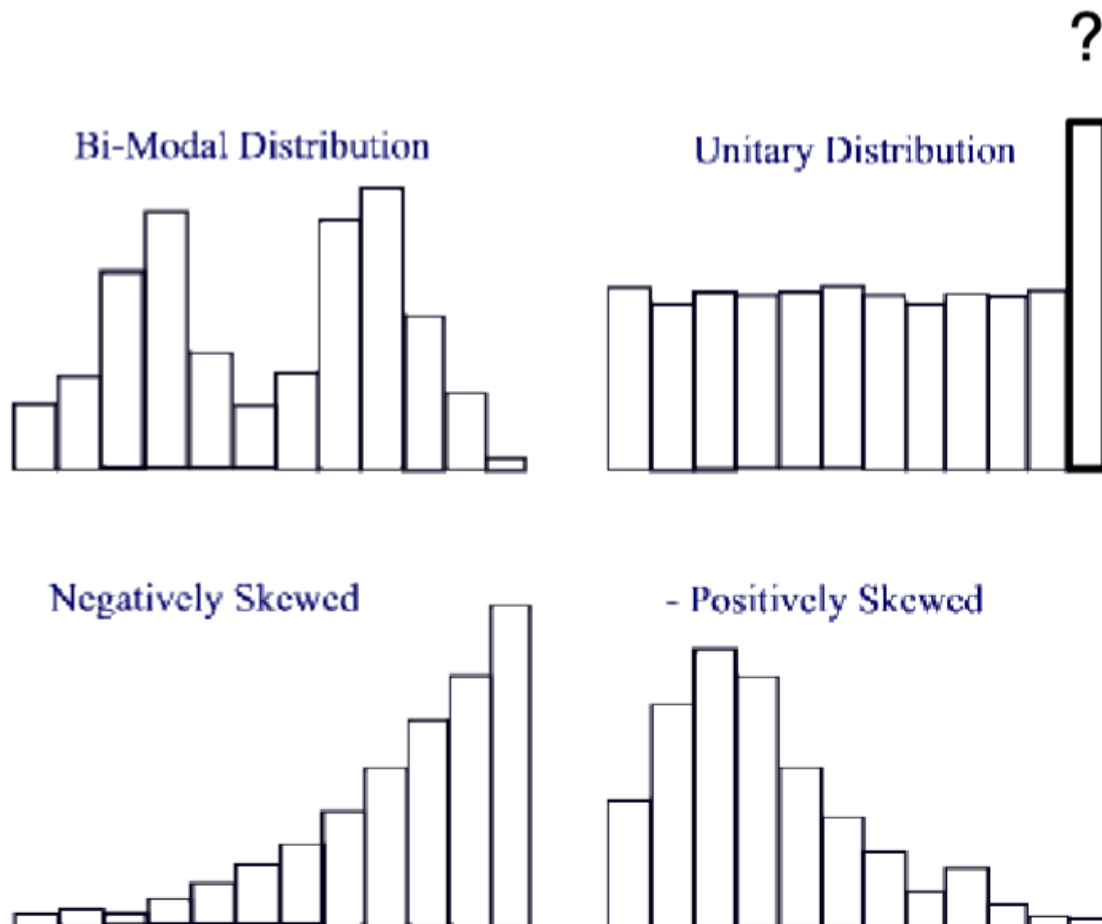
1.3 Data Cleaning: Brief Recap

80% of your Data Science time

Some pointers:

- Missing values – do not replace with 0 or -999

- Treat as missing
- Look for frequencies of values or outliers (-999)
 - **Histograms** invaluable
 - Examine **outliers**



2 Correlation

We'll be look at feature distributions independently (per column), but also relationships between features (columns). Correlation is one of the first tools you think of for looking at the relationship between two features (but has limitations)

2.1 Correlated variables

Two variables are correlated if changes in one variable, correspond to changes in the other
Correlation in general refers to the statistical association between two variables

- This statistical association allows us to make estimates for the one variable based on the value of the other

- While we typically think of linear relationships between variables, nonlinear might exist too

2.2 Linear correlation

- Two variables x and y are said to be linearly correlated if their relationship can be described by the following equation:

$$y = \alpha + \beta x, \beta \neq 0 \quad (1)$$

- This relationship will not be exact
- There will be an error associated with it, and hence, in reality:
- ϵ is the error term for the relationship

$$y = \alpha + \beta x + \epsilon, \beta \neq 0 \quad (2)$$